

---

# Unexpected Improvements to Expected Improvement for Bayesian Optimization

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Expected Improvement (EI) is arguably the most popular acquisition function  
2 in Bayesian optimization and has been applied successfully in practice, but its  
3 performance is often exceeded by that of more recent methods. However, EI and its  
4 variants, including for the parallel and multi-objective settings, are very challenging  
5 to optimize due to the fact that their acquisition values are numerically zero in  
6 many regions, resulting in inconsistencies in EI’s performance across the literature.  
7 This difficulty generally increases as the number of observations, dimensionality  
8 of the search space, or the number of constraints grow. We propose a new family  
9 of acquisition functions, LogEI, whose members either have identical or approxi-  
10 mately equal optima as their canonical counterparts, but are substantially easier  
11 to optimize numerically. We demonstrate that numerical pathologies manifest  
12 themselves in “classic” analytic EI, as well as its variants, including constrained  
13 EI, parallel EI, and expected hypervolume improvement, and propose correspond-  
14 ing reformulations that avoid these pathologies. Our empirical results show that  
15 members of the LogEI family of acquisition functions achieve substantially better  
16 performance than their canonical counterparts and surprisingly, are on par with or  
17 exceed the performance of recent state-of-the-art acquisition functions, highlighting  
18 the understated role of numerical optimization in the literature.

## 19 1 Introduction

20 Bayesian Optimization (BO) is a popular framework for sample-efficient black-box optimization  
21 of expensive-to-evaluate functions [15, 18]. BO leverages a probabilistic *surrogate model* in con-  
22 junction with an *acquisition function* to determine where to query the underlying objective function.  
23 Improvement-based acquisition functions, such as Expected Improvement (EI) and Probability Im-  
24 provement (PI), are among the earliest and most widely used acquisition functions for efficient global  
25 optimization of non-convex functions [30, 42]. EI has been extended to the constrained [17, 19],  
26 noisy [38], and multi-objective [13] setting, as well as their respective batch variants [2, 8, 60], and  
27 is a standard baseline in the BO literature [15, 50]. While much of the literature has focused on  
28 developing new sophisticated acquisition functions, subtle yet critical implementation details of  
29 foundational BO methods are often overlooked. Notably, the performance of EI and its variants can  
30 vary greatly from one implementation to another, leading to inconsistent results even for *mathemati-*  
31 *cally identical* formulations. To our knowledge, these inconsistencies have not been detailed in the  
32 literature, although the problem of optimizing EI effectively has been discussed in various works  
33 (e.g., [15, 21, 60]).

34 In this work, we identify pathologies of several improvement-based acquisition functions and propose  
35 reformulations utilizing careful numerical implementations that lead to increases in the optimization  
36 performance of improvement-based acquisition functions which often match or exceed that of recent  
37 methods. Our main contributions are:

1. We introduce LogEI, a new family of acquisition functions whose members either have identical or approximately equal optima as their canonical counterparts, but are substantially easier to optimize numerically. Notably, the analytic variant of LogEI, which results in the same BO policy as EI (mathematically), empirically shows significantly improved optimization performance.
2. We extend the ideas behind analytical LogEI to a number of other members of the EI family, including constrained EI (cEI), Expected Hypervolume Improvement (EHVI), as well as their respective batch variants for parallel BO, qEI and qEHVI, using smooth approximations of the acquisition utilities to obtain non-zero gradients.
3. We demonstrate that our newly proposed acquisition functions substantially outperform their respective analogues on a broad range of benchmarks without incurring meaningful additional computational cost.

While we focus on Expected Improvement, the same issues and similar solutions also apply to Probability of Improvement [35] and its variants (see Appendix A).

## Motivation

Maximizing acquisition functions in BO is a challenging problem. The optimization surface is generally non-convex, and often contains many local maxima. While zeroth-order methods are sometimes used, on continuous domains gradient-based methods tend to be far more effective at optimizing acquisition functions, especially in higher dimensions.

In addition to the challenges stemming from non-convexity that are shared across acquisition functions, improvement-based acquisition functions are particularly challenging to optimize because their value and gradient can be minuscule in large swaths of its domain. Although EI is never *mathematically* zero under a Gaussian posterior distribution (except at previously evaluated points in the case of noiseless observations), it often is *exactly* zero numerically due to floating point precision. The same applies to its gradient, making EI (and PI, see Appendix A) difficult to optimize via gradient-based methods. Figure 1 illustrates this behavior on a simple one-dimensional problem.

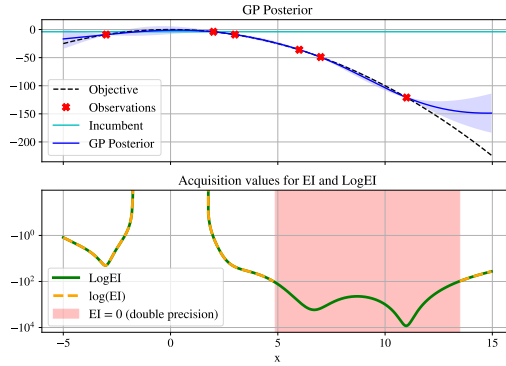


Figure 1: Values of EI and LogEI on a quadratic objective. For points at which the likelihood of improving over the incumbent is small, EI takes on extremely small values. Notably, EI is *exactly* zero numerically (in double floating point precision) in a large part of the domain ( $\approx [5, 13.5]$ ). As Figure 2 to the right shows, this behavior gets worse as the dimension of the problem and the number of observed data points grow, rendering gradient-based optimization of EI futile.

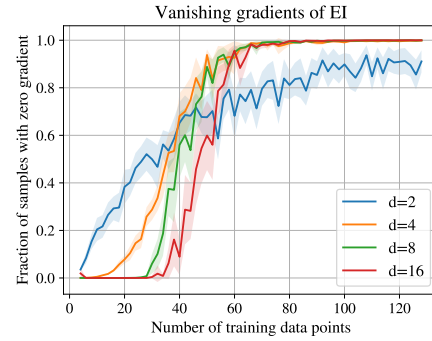


Figure 2: Empirical fraction of points randomly sampled from the domain for which the gradient of EI is approximately numerically zero ( $< 10^{-10}$ ) as a function of the number of (randomly generated) data points  $n$  for different dimensions  $d$  on the Ackley. As the model uncertainty shrinks with more data, EI and its gradients become numerically zero across most of the domain.

To increase the chance of finding the global optimum of non-convex functions, gradient-based optimization is typically performed from multiple starting points [54]. A multi-start gradient-based approach can indeed help avoid getting stuck in local optima. However, for improvement-based acquisition functions, optimization becomes increasingly challenging as more data is collected and the likelihood of improving over the incumbent diminishes (see our theoretical results in Section 3 and the empirical illustration in Figure 2). As a result, gradient-based optimization with multiple random

starting points will eventually degenerate into random search when the gradients at the starting points are numerically zero. This problem is particularly acute in high dimensions and for objectives with a large range.

Various initialization heuristics have been proposed to address this behavior by modifying the random-restart strategy. Rather than starting from random candidates, an alternative naïve approach would be to use initial conditions close to the best previously observed inputs. However, doing that alone inherently limits the initialization heuristic to a type of local search, which cannot have global guarantees. To attain such guarantees, it is necessary to use an asymptotically space-filling heuristic; even if not random, this will entail evaluating the acquisition function in regions where no prior observation lies. Ideally, these regions should permit gradient-based optimization of the objective for efficient acquisition function optimization, which necessitates the gradients to be non-zero. In this work, we show that this can be achieved for a number of improvement-based acquisition functions, and demonstrate empirically how this leads to substantially improved BO performance.

## 2 Background

We consider the problem of maximizing an expensive-to-evaluate black-box function  $f_{\text{true}} : \mathbb{X} \mapsto \mathbb{R}^M$  over some feasible set  $\mathbb{X} \subset \mathbb{R}^d$ . Suppose we have collected data  $\mathcal{D}_n = \{(x_i, y_i)\}_{i=1}^n$ , where  $x_i \in \mathbb{X}$  and  $y_i = f_{\text{true}}(x_i) + v_i(x_i)$  with  $v_i$  some noise corrupting the true function value  $f_{\text{true}}(x_i)$ . The response  $f_{\text{true}}$  may be multi-output (e.g. in the case of multiple objectives or black-box constraints), in which case  $y_i, v_i \in \mathbb{R}^M$ . We use Bayesian optimization (BO), which relies on a surrogate model  $f$  that for any *batch*  $\mathbf{x} := \{x_1, \dots, x_q\}$  of candidate points provides a probability distribution over the outputs  $f(\mathbf{x}) := (f(x_1), \dots, f(x_q))$ . The acquisition function  $\alpha$  then utilizes this posterior prediction to assign an acquisition value to  $\mathbf{x}$  that quantifies the value of evaluating the points in  $\mathbf{x}$ , trading off exploration and exploitation.

### 2.1 Gaussian Processes

Gaussian Processes (GP) models [49] are the most widely used surrogates in BO, due to their high data efficiency and good uncertainty quantification. For our purposes, it suffices to consider a GP as a mapping that provides a multivariate Normal distribution over the outputs  $f(\mathbf{x})$  for any  $\mathbf{x}$ :

$$f(\mathbf{x}) \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}), \boldsymbol{\Sigma}(\mathbf{x})), \quad \boldsymbol{\mu} : \mathbb{X}^q \rightarrow \mathbb{R}^{qM}, \quad \boldsymbol{\Sigma} : \mathbb{X}^q \rightarrow \mathcal{S}_+^{qM}. \quad (1)$$

In the single-outcome ( $M = 1$ ) setting,  $f(\mathbf{x}) \sim \mathcal{N}(\mu(\mathbf{x}), \Sigma(\mathbf{x}))$  with  $\mu : \mathbb{X}^q \rightarrow \mathbb{R}^q$  and  $\Sigma : \mathbb{X}^q \rightarrow \mathcal{S}_+^q$ . In the sequential ( $q = 1$ ) case, this further reduces to a univariate Normal distribution:  $f(x) \sim \mathcal{N}(\mu(x), \sigma^2(x))$  with  $\mu : \mathbb{X} \rightarrow \mathbb{R}$  and  $\sigma : \mathbb{X} \rightarrow \mathbb{R}_+$ .

### 2.2 Expected Improvement and (some of) its variants

**Expected Improvement** In the fully-sequential ( $q = 1$ ), single-outcome ( $M = 1$ ) setting, “classical” EI [43] is defined as

$$\text{EI}_{y^*}(x) = \mathbb{E}_{f(x)}[f(x) - y^*]_+ = \sigma(x) h\left(\frac{\mu(x) - y^*}{\sigma(x)}\right), \quad (2)$$

where  $[\cdot]_+$  denotes the  $\max(0, \cdot)$  operation,  $y^*$  is the *incumbent* (the best function value observed so far),  $h(z) = \phi(z) + z\Phi(z)$ , and  $\phi, \Phi$  are the standard Normal density and distribution functions, respectively. This formulation is arguably the most widely used acquisition function in BO, and the default in many popular software packages.

**Constrained Expected Improvement** The *constrained BO* setting involves one or more black-box constraints; the problem is typically formulated as finding  $\max_{x \in \mathbb{X}} f_{\text{true},1}(x)$  such that  $f_{\text{true},i}(x) \leq 0$  for  $i \in \{2, \dots, M\}$ . Feasibility-weighting the improvement [17, 19] is a natural way to formulate constrained EI:

$$\text{cEI}_{y^*}(x) = \mathbb{E}_{f(x)}[f_1(x) - y^*]_+ \Pi_{i=2}^M \mathbf{1}\{f_i(x) \leq 0\}, \quad (3)$$

where  $\mathbf{1}$  is the indicator function. If the constraints are modeled as conditionally independent of the objective  $f_1$  this can be further simplified as the product of EI and the probability of feasibility.

114 **Parallel Expected Improvement** In many settings, one may evaluate  $f_{\text{true}}$  on  $q > 1$  candidates in  
 115 parallel to increase throughput. The associated batch analogue of EI [20, 58] is given by

$$\text{qEI}_{y^*}(\mathbf{x}) = \mathbb{E} \left[ \max_{j=1, \dots, q} \{ [f(x_j) - y^*]_+ \} \right]. \quad (4)$$

116 Unlike EI, qEI does not admit a closed-form expression and thus is typically computed using Monte  
 117 Carlo sampling [2, 58]:

$$\text{qEI}_{y^*}(\mathbf{x}) \approx \sum_{i=1}^N \max_{j=1, \dots, q} \{ [\xi_j^i(\mathbf{x}) - y^*]_+ \}, \quad (5)$$

119 where  $\xi^i(\mathbf{x}) \sim f(\mathbf{x})$  are random samples drawn from the (joint) model posterior at  $\mathbf{x}$ .  
 118

120 **Expected Hypervolume Improvement** In multi-objective optimization (MOO), there generally is  
 121 no single best solution; instead the goal is to explore the Pareto Frontier between multiple competing  
 122 objectives (the set of mutually-optimal objective vectors). A common measure of the quality of a  
 123 (finitely approximated) Pareto Frontier  $\mathcal{P}$  between  $M$  objectives with respect to a specified reference  
 124 point  $\mathbf{r} \in \mathbb{R}^M$  is its *hypervolume*  $\text{HV}(\mathcal{P}, \mathbf{r}) := \lambda(\bigcup_{i=1}^{|\mathcal{P}|} [\mathbf{r}, y_i])$ , where  $[\mathbf{r}, y_i]$  denotes the hyper-  
 125 rectangle bounded by vertices  $\mathbf{r}$  and  $y_i$ , and  $\lambda$  is the Lebesgue measure. Thus, a natural acquisition  
 126 function to optimize is the expected hypervolume improvement (EHVI)

$$\text{EHVI}(\mathbf{x}) = \mathbb{E} [\text{HV}(\mathcal{P} \cup f(\mathbf{x}), \mathbf{r})] - \text{HV}(\mathcal{P}, \mathbf{r}) \quad (6)$$

127 from obtaining a set  $\mathcal{Y} \sim f(\mathbf{x})$  of new observations. If  $q = 1$  and the objectives are modeled with  
 128 independent GPs, EHVI can be expressed in closed form [63]. In the general case, Monte Carlo  
 129 approximations are used (qEHVI) [8].

## 130 2.3 Optimizing acquisition functions

131 Optimizing an acquisition function (AF) is a challenging task that amounts to solving a non-convex  
 132 optimization problem, to which multiple approaches and heuristics have been applied. These include  
 133 gradient-free methods such as divided rectangles [29], evolutionary methods such as CMA-ES [22],  
 134 first-order methods such as stochastic gradient ascent (see e.g., Daulton et al. [9], Wang et al. [58]),  
 135 and (quasi-)second order methods [15] such as L-BFGS-B [6]. Multi-start optimization is commonly  
 136 employed with gradient-based methods to mitigate the risk of getting stuck in local minima. Initial  
 137 points for optimization are selected via various heuristics with different levels of complexity, ranging  
 138 from simple uniform random selection to BoTorch’s initialization heuristic, which selects initial  
 139 points by performing Boltzmann sampling on a set of random points according to their acquisition  
 140 function value [2]. See Appendix B for a more complete account of initialization strategies and  
 141 optimization procedures used by popular implementations. We focus on gradient-based optimization  
 142 as often leveraging gradients results in faster and more performant optimization [8].

143 Optimizing AFs for parallel BO that quantify the value of a batch of  $q > 1$  points is more challenging  
 144 than optimizing their sequential counterparts due to the higher dimensionality of the optimization  
 145 problem ( $qd$  instead of  $d$ ) and the more challenging optimization surface. A common approach  
 146 to simplify the problem is to use a *sequential greedy* strategy that greedily solves a sequence of  
 147 single point selection problems. For  $i = 1, \dots, q$ , candidate  $x_i$  is selected by optimizing the AF (for  
 148  $q = 1$ ) conditional on the previously selected designs  $x_1, \dots, x_{i-1}$  and their unknown observations  
 149 (e.g. by fantasizing the values at those designs) [60]. For submodular AFs, including EI, PI, and  
 150 EHVI, using a sequential greedy strategy will incur no more than  $1/e$  regret compared to joint  
 151 optimization, and previous works have found that sequential greedy optimization yields *improved*  
 152 BO performance compared to joint optimization [8, 60]. Contrary to these findings, we show that  
 153 with our reformulations joint batch optimization is indeed superior to the sequential greedy strategy.

## 154 2.4 Related Work

155 While there is a substantial body of work introducing a large variety of different AFs, much less focus  
 156 has been on the question of how to effectively implement and optimize these AFs.

157 Zhan and Xing [64] provide a comprehensive review of a large number of different variants of the  
 158 EI family, but do not discuss any numerical or optimization challenges. Zhao et al. [65] propose to  
 159 combine a variety of different initialization strategies to select initial conditions for optimization of

acquisition functions and show empirically that this improves optimization performance. However, they do not address any potential issues or degeneracies with the acquisition functions themselves. Recent works have considered effective gradient-based approaches for acquisition optimization. Wilson et al. [60] demonstrates how stochastic first-order methods can be leveraged for optimizing Monte Carlo acquisition functions. Balandat et al. [2] build on this work and propose to use a sample average approximation to MC acquisition functions that admits gradient-based optimization using deterministic higher-order optimizers such as L-BFGS-B.

Another line of work proposes to switch from BO to local optimization based on some stopping criterion to achieve faster local convergence, using either zeroth order [44] or gradient-based [41] optimization. While McLeod et al. [41] are also concerned with numerical issues, we emphasize that those issues arise due to ill-conditioned covariance matrices and are orthogonal to the numerical pathologies of improvement-based acquisition functions.

### 3 Theoretical Analysis of Expected Improvement’s Vanishing Gradients

In this section, we shed light on the conditions on the objective function and surrogate model that give rise to the numerically vanishing gradients in EI, as seen in Figure 2. In particular, we show that as a BO algorithm closes the optimality gap  $f^* - y_n^*$  (here  $f^*$  is the global maximum of the function  $f_{\text{true}}$ ) and the associated GP surrogate’s uncertainty decreases, it becomes exceedingly likely for EI to exhibit numerically vanishing gradients.

Given a distribution  $P_x$  over the inputs  $x$ , the probability that the argument  $(\mu(x) - y^*)/\sigma(x)$  to  $h$  in Eq. (2) is smaller than a threshold  $B$  exceeds  $P_x(f(x) < f^* - \epsilon_n)$ , with high probability, where  $\epsilon_n$  depends on the optimality gap  $f^* - y_n^*$  and the maximum posterior uncertainty  $\max_x \sigma_n(x)$ . This pertains to the problem of numerically vanishing values and gradients of EI, since the numerical support  $\mathcal{S}_\eta(h) = \{x : |h(x)| > \eta\}$  of a naïve implementation of  $h$  in (2) is limited by a lower bound  $B(\eta)$  that depends on the floating point precision  $\eta$ . Formally,  $\mathcal{S}_\eta(h) \subset [B(\eta), \infty)$  even though  $\mathcal{S}_0(h) = \mathbb{R}$  mathematically. As a consequence, the following result can be seen as a bound on the probability of encountering numerically vanishing values and gradients in EI using samples from the distribution  $P_x$  to initialize the acquisition function optimization.

**Lemma 1.** *Suppose  $f$  is drawn from a Gaussian process prior  $P_f$ ,  $y^* \leq f^*$ ,  $\mu_n, \sigma_n$  are the mean and standard deviation of the posterior  $P_f(f|\mathcal{D}_n)$  and  $B \in \mathbb{R}$ . Then with probability  $1 - \delta$ ,*

$$P_x \left( \frac{\mu_n(x) - y_n^*}{\sigma_n(x)} < B \right) \geq P_x(f(x) < f^* - \epsilon_n) \quad (7)$$

where  $\epsilon_n = (f^* - y_n^*) + (\sqrt{-2 \log(2\delta)} - B) \max_x \sigma_n(x)$ .

For any given – and especially early – iteration,  $\epsilon_n$  does not have to be small, as both the optimality gap and the maximal posterior standard deviation can be large initially. Note that under certain technical conditions on the kernel function and the asymptotic distribution of the training data  $\mathcal{D}_n$ , the maximum posterior variance vanishes guaranteeably as  $n$  increases, see [36, Corollary 3.2]. On its own, Lemma 1 gives insight into the non-asymptotic behavior by exposing a dependence to the distribution of objective values  $f$ . In particular, if the set of inputs that give rise to high objective values ( $\approx f^*$ ) is concentrated,  $P(f(x) < f^* - \epsilon)$  will decay very slowly as  $\epsilon$  increases, thereby maintaining a lower bound on the probability of close to 1. As an example, this is the case for the Ackley function, especially as the dimensionality increases, which explains the behavior in Figure 2.

## 4 Improving Expected Improvement (and its cousins)

In this section, we propose modifications to and re-formulations of analytic and MC-based improvement-based acquisition functions that render them significantly easier to optimize. In the following, we will use differing fonts, e.g.  $\log$  and  $\text{Log}$ , to differentiate between the mathematical functions and their numerical implementations.

### 4.1 Analytic LogEI

Implementations of “classic” analytic EI exhibit numerically zero values and gradients even though EI and its gradient are mathematically nonzero on the entire real line, except in the noiseless case

at the previously evaluated points. However, if implemented naïvely,  $h$  is numerically zero when  $(\mu(x) - y^*)/\sigma(x)$  is small, which happens when the model has high confidence that little improvement can be achieved at  $x$ .

We propose an implementation of  $\log \circ h$  that can be accurately computed for a much larger range of inputs than a naïve implementation  $h$  or  $\log \circ h$ . Specifically, we compute analytic

$$\text{LogEI}_{y^*}(x) = \text{log\_h}((\mu(x) - y^*)/\sigma(x)) + \log(\sigma(x)), \quad (8)$$

where  $\text{log\_h}$  is mathematically equivalent to  $\log \circ h$  and can be stably and accurately computed by

$$\text{log\_h}(z) = \begin{cases} \log(\phi(z) + z\Phi(z)) & z > -1 \\ -z^2/2 - c_1 + \text{log1mexp}(\text{logerfcx}(-z/\sqrt{2})|z| + c_2) & z < -1 \end{cases} \quad (9)$$

where  $c_1 = \log(2\pi)/2$ , and  $c_2 = \log(\pi/2)/2$ , and  $\text{log1mexp}$ ,  $\text{logerfcx}$  are numerically stable implementations of  $\log(1 - \exp(x))$  and  $\log(\exp(x^2)\text{erfc}(x))$ , respectively (see Appendix A for details). Notably, the asymptotically quadratic behavior of  $\text{log\_h}$  becomes apparent in the second case, making the function particularly amenable to gradient-based optimization. This has *significant* practical implications for BO using EI, as evidenced by the empirical results in Section 5. Numerically vanishing values and gradients affect – as far as we are aware – all public implementations of EI.

## 4.2 Monte Carlo Parallel LogEI

For Monte Carlo formulations of Parallel EI that perform differentiation on the level of MC samples, the situation is worse in that they exhibit not just numerically, but mathematically zero gradients for a significant proportion of practically relevant inputs. For qEI, the primary issue is taking the discrete maximum over the  $q$  outcomes for each MC sample in (5). In particular, the acquisition utility of expected improvement in Eq. 4 on a single sample  $\xi_i$  of  $f$  is  $\max_j [\xi_i(x_j) - y^*]_+$ . Mathematically, we smoothly approximate the acquisition utility in two stages: 1)  $u_{ij} = \text{softplus}_{\tau_0}(\xi_i(x_j) - y^*) \approx [\xi_i(x_j) - y^*]_+$  and 2)  $\|u_{i\cdot}\|_{1/\tau_{\max}} \approx \max_j u_{ij}$ . Since the resulting quantities are strictly positive, they can be transformed to log-space permitting an implementation of qLogEI that is numerically stable and can be optimized effectively, similar to the analytic case. In particular,

$$\begin{aligned} \text{qLogEI}_{y^*}(\mathbf{x}) &= \log \int \left( \sum_j \text{softplus}_{\tau_0}(f(x_j) - y^*)^{1/\tau_{\max}} \right)^{\tau_{\max}} df \\ &\approx \text{logsumexp}_i(\tau_{\max} \text{logsumexp}_j(\text{logsoftplus}_{\tau_0}(\xi^i(x_j) - y^*)/\tau_{\max})), \end{aligned} \quad (10)$$

where  $i$  is the index of the Monte Carlo draws from the GP posterior,  $j = 1, \dots, q$  is the index for the candidate in the batch, and  $\text{logsoftplus}$  is a numerically stable implementation of  $\log(\log(1 + \exp(x)))$ . See Appendix A for additional details.

While the smoothing in (10) approximates the original qEI formulation, the following result shows that the resulting approximation error can be quantified and tightly bounded as a function of the temperature parameters  $\tau_0$ ,  $\tau_{\max}$  and the batch size  $q$ . See Appendix C for the proof.

**Lemma 2.** [Approximation Guarantee] *Given the temperature parameters  $\tau_0$  and  $\tau_{\max}$ , the approximation error of qLogEI to qEI is bounded by*

$$|\exp(\text{qLogEI}(x)) - \text{qEI}(x)| \leq \log(2) \tau_0 + q^{\tau_{\max}} - 1. \quad (11)$$

In Appendix D, we show the importance of setting the temperatures sufficiently low for qLogEI to achieve good optimization characteristics, something that only becomes possible by transforming all involved computations to log-space. Otherwise, the smoothed approximation to the acquisition utility (e.g., using a regular  $\text{softplus}$  function) would similarly exhibit numerically vanishing gradients, as is the case mathematically for the discrete  $\max$  operator.

## 4.3 Constrained EI

Both analytic and Monte Carlo variants of LogEI can be extended for optimization problems with black-box constraints. For analytic cEI with independent constraints of the form  $f_i(x) \leq 0$ , the constrained formulation in Eq. (3) simplifies to  $\text{logcEI}(x) = \text{LogEI}(x) + \sum_i \log(P(f_i(x) \leq 0))$ , which can be readily and stably computed using LogEI in Eq. (8) and, if  $f_i$  is modelled by a GP, a stable implementation of the Gaussian log cumulative distribution function. For the Monte Carlo variant, we apply a similar strategy as for Eq. (10) to the constraint indicators in Eq. (3): 1) a smooth approximation via a sigmoid and 2) an accurate and stable implementation of its log value (see Appendix A for details).

#### 4.4 Monte Carlo Parallel LogEHVI

The numerical difficulties of qEHVI in (6) are similar to those of qEI, and the basic ingredients of smoothing and log-transformations still apply, but the details are significantly more complex since qEHVI uses many operations that have mathematically zero gradients with respect to some of the inputs. Our implementation is based on the differentiable inclusion-exclusion formulation of the hypervolume improvement [8]. As a by-product, the implementation also readily allows for the differentiable computation of the expected log hypervolume, instead of the log expected hypervolume, note the order, which can be preferable in certain applications of multi-objective optimization [16].

### 5 Empirical Results

We compare standard versions of analytic EI (EI) and constrained EI (cEI), Monte Carlo parallel EI (qEI), as well as Monte Carlo EHVI (qEHVI). We also compare with other state-of-the-art baselines, namely lower-bound Max-Value Entropy Search (GIBBON) [45] and single- and multi-objective Joint Entropy Search (JES) [25, 55]. All experiments are implemented using BoTorch [2] and utilize multi-start optimization of the AF with `scipy`'s L-BFGS-B optimizer. In order to avoid conflating the effect of BoTorch's default initialization strategy with those of our contributions, we use 16 initial points chosen uniformly at random from which to start the L-BFGS-B optimization (for a comparison with other initialization strategies, see the Appendix D). We run multiple replicates and report mean and error bars of  $\pm 2$  standard errors of the mean. Additional details can be found in Appendix D.1.

**Single-objective, sequential BO** We compare EI and LogEI on the 10-dimensional convex Sum-of-Squares (SoS) function  $f(x) = \sum_{i=1}^{10} (x_i - 0.5)^2$ , using 20 restarts seeded from 1024 pseudo-random samples through BoTorch's default initialization heuristic. Figure 3 shows that due to vanishing gradients, EI is unable to make progress even on this trivial problem.

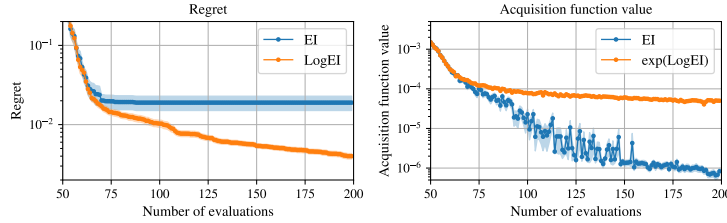


Figure 3: Regret and EI acquisition value for the candidates selected by maximizing EI and LogEI on the convex Sum-of-Squares problem. Optimization stalls out for EI after about 75 observations due to vanishing gradients (indicated by the jagged behavior of the acquisition value), while LogEI continues to make steady progress.

In Figure 4, we compare performance on the Ackley and Michalewicz test functions [51]. Notably, LogEI substantially outperforms EI on Ackley as the dimensionality increases. Ackley is a challenging multimodal function for which it is critical to trade off local exploitation with global exploration, a task made exceedingly difficult by the numerically vanishing gradients of EI in a large fraction of the search space. We see a similar albeit less pronounced behavior on Michalewicz, which reflects the fact that Michalewicz is a somewhat less challenging problem than Ackley.

**BO with Black Box Constraints** Figure 5 shows results on four engineering design problems with black box constraints, where LogcEI massively outperforms the naive cEI implementation. Similar to the unconstrained problems, the performance gains of LogcEI over cEI grow with increasing dimensionality of the problems, as well as with the number of constraints. Notably, while running these benchmarks we found that for some problems, LogcEI in fact improved upon some of the best results quoted in the original literature that used up to almost three orders of magnitude more function evaluations (see Appendix D).

**Parallel Expected Improvement with qLogEI** Figure 6 reports optimization performance of parallel BO on the 16-dimensional Ackley and Levy functions for both sequential greedy and joint

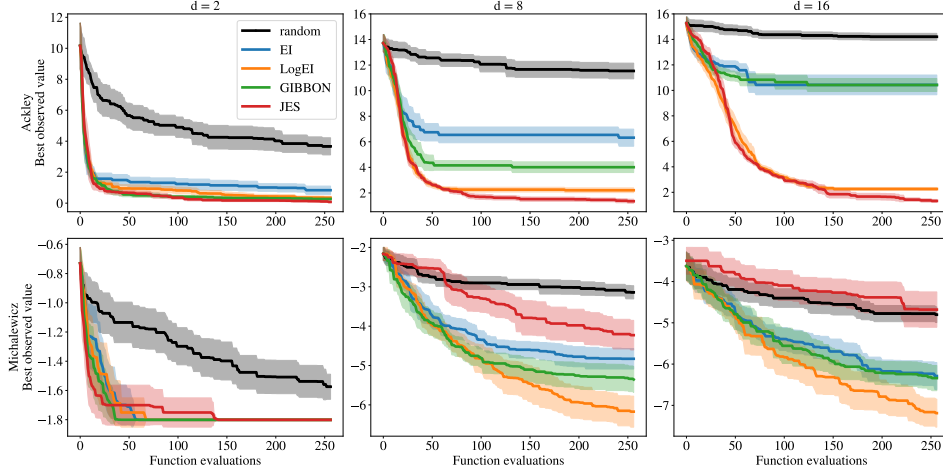


Figure 4: Best objective value as a function of iterations on the moderately and severely non-convex Michalewicz and Ackley problems for varying numbers of input dimensions. LogEI substantially outperforms both EI and GIBBON, and this gap widens as the problem dimensionality increases. JES performs slightly better than LogEI on Ackley, but for some reason fails on Michalewicz. Notably, JES is almost two orders of magnitude slower than the other acquisition functions (see Appendix D).

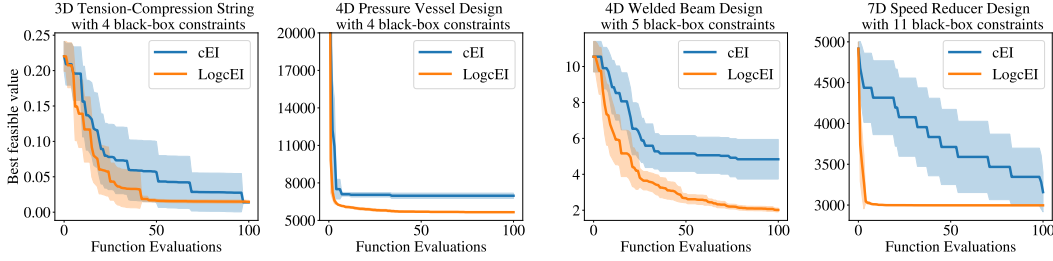


Figure 5: Best feasible objective value as a function of number of function evaluations (iterations) on four engineering design problems with black-box constraints after an initial  $2d$  pseudo-random evaluations.

batch optimization. Besides the apparent substantial advantages of qLogEI over qEI, they key observation here is that, in contrast to previous findings in the literature [60], jointly optimizing the candidates of batch acquisition functions can yield highly competitive optimization performance.

**Multi-Objective optimization with qLogEHVI** Figure 7 compares qLogEHVI and qEHVI on 6 different test problems with 2 or 3 objectives, and ranging from 2-30 dimensions. This includes 3 real world inspired problems: cell network design for optimizing coverage and capacity [12], laser plasma acceleration optimization [26], and vehicle design optimization [40, 52]. The results are consistent with our findings in the single-objective and constrained cases: qLogEHVI consistently outperforms qEHVI, and the gap is larger on higher dimensional problems. See the Appendix D for problem details.

## 6 Conclusion

Our results demonstrate that the problem of vanishing gradients is a major source of the difficulty of optimizing improvement-based acquisition functions and that we can mitigate this issue through careful reformulations and implementations. As a result, we see substantially improved optimization performance across a variety of modified EI variants across a broad range of problems. In particular, contrary to previous findings, we demonstrate that joint batch optimization for parallel BO outperforms the sequential greedy approach typically used in practice (which also benefits from our modifications). Besides the convincing performance gains, one of the key advantages of our



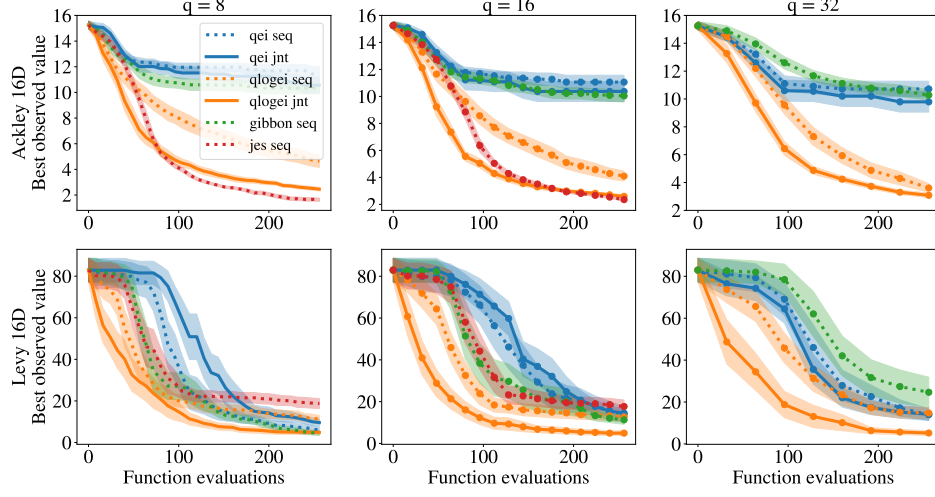


Figure 6: Best objective value for parallel BO as a function of the number evaluations for single-objective optimization in 16 dimensions with varying batch sizes  $q$ . The advantage of qLogEI over qEI grows as the dimensionality increases and, in accordance with Lem. 1, the differences are less pronounced on the Levy function for which the distribution of inputs with good objective values is less concentrated in the search space. Notably and surprisingly, joint optimization of the batch outperforms sequential greedy optimization, overturning previous findings for canonical qEI [60].

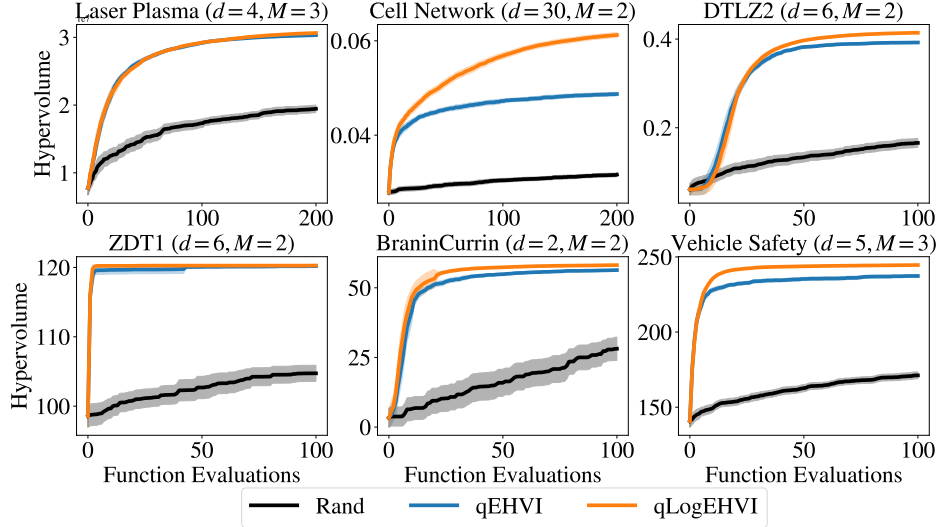


Figure 7: Sequential ( $q = 1$ ) optimization performance on multi-objective problems, as measured by the hypervolume of the Pareto frontier across observed points.

306 modified acquisition functions is that they are much less dependent on heuristic (and potentially  
307 brittle) initialization strategies. Moreover, our proposed modifications do not meaningfully increase  
308 the computational complexity of the respective original acquisition function.

309 While our contributions may not apply verbatim to other classes of acquisition functions, our key  
310 insights and strategies do translate and could help e.g. with improving information-based [24, 59],  
311 cost-aware [37, 50], and other types of acquisition functions that are prone to similar numerical  
312 challenges. Overall, we hope that our findings will increase awareness in the community for the  
313 importance of optimizing acquisition functions well, and for the required care that is needed regarding  
314 the involved numerics.

## References

- [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. URL <https://www.tensorflow.org/>. Software available from tensorflow.org.
- [2] Maximilian Balandat, Brian Karrer, Daniel R. Jiang, Samuel Daulton, Benjamin Letham, Andrew Gordon Wilson, and Eytan Bakshy. BoTorch: A Framework for Efficient Monte-Carlo Bayesian Optimization. In *Advances in Neural Information Processing Systems 33*, 2020.
- [3] Ricardo Baptista and Matthias Poloczek. Bayesian optimization of combinatorial structures, 2018.
- [4] Syrine Belakaria, Aryan Deshwal, and Janardhan Rao Doppa. Max-value entropy search for multi-objective bayesian optimization with constraints, 2020.
- [5] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- [6] Richard H Byrd, Peihuang Lu, Jorge Nocedal, and Ciyu Zhu. A limited memory algorithm for bound constrained optimization. *SIAM Journal on Scientific Computing*, 16(5):1190–1208, 1995.
- [7] Carlos A Coello Coello and Efrén Mezura Montes. Constraint-handling in genetic algorithms through the use of dominance-based tournament selection. *Advanced Engineering Informatics*, 16(3):193–203, 2002.
- [8] Samuel Daulton, Maximilian Balandat, and Eytan Bakshy. Differentiable expected hypervolume improvement for parallel multi-objective bayesian optimization. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 9851–9864. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/6fec24eac8f18ed793f5ead3dd7977c-Paper.pdf>.
- [9] Samuel Daulton, Xingchen Wan, David Eriksson, Maximilian Balandat, Michael A. Osborne, and Eytan Bakshy. Bayesian optimization over discrete and mixed spaces via probabilistic reparameterization. In *Advances in Neural Information Processing Systems 35*, 2022.
- [10] Kalyan Deb, L. Thiele, Marco Laumanns, and Eckart Zitzler. Scalable multi-objective optimization test problems. volume 1, pages 825–830, 06 2002. ISBN 0-7803-7282-4. doi: 10.1109/CEC.2002.1007032.
- [11] Aryan Deshwal, Sebastian Ament, Maximilian Balandat, Eytan Bakshy, Janardhan Rao Doppa, and David Eriksson. Bayesian optimization over high-dimensional combinatorial spaces via dictionary-based embeddings. In Francisco Ruiz, Jennifer Dy, and Jan-Willem van de Meent, editors, *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, volume 206 of *Proceedings of Machine Learning Research*, pages 7021–7039. PMLR, 25–27 Apr 2023.
- [12] Ryan M. Dreifuerst, Samuel Daulton, Yuchen Qian, Paul Varkey, Maximilian Balandat, Sanjay Kasturia, Anoop Tomar, Ali Yazdan, Vish Ponnampalam, and Robert W. Heath. Optimizing coverage and capacity in cellular networks using machine learning, 2021.
- [13] M. T. M. Emmerich, K. C. Giannakoglou, and B. Naujoks. Single- and multiobjective evolutionary optimization assisted by gaussian random field metamodels. *IEEE Transactions on Evolutionary Computation*, 10(4):421–439, 2006.

- [14] David Eriksson, Michael Pearce, Jacob Gardner, Ryan D Turner, and Matthias Poloczek. Scalable global optimization via local Bayesian optimization. In *Advances in Neural Information Processing Systems 32*, NeurIPS, 2019.
- [15] Peter I Frazier. A tutorial on bayesian optimization. *arXiv preprint arXiv:1807.02811*, 2018.
- [16] Tobias Friedrich, Karl Bringmann, Thomas Voß, and Christian Igel. The logarithmic hypervolume indicator. In *Proceedings of the 11th workshop proceedings on Foundations of genetic algorithms*, pages 81–92, 2011.
- [17] 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 of *Proceedings of Machine Learning Research*, pages 937–945, Beijing, China, 22–24 Jun 2014. PMLR.
- [18] Roman Garnett. *Bayesian Optimization*. Cambridge University Press, 2023. to appear.
- [19] 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.
- [20] David Ginsbourger, Rodolphe Le Riche, and Laurent Carraro. A Multi-points Criterion for Deterministic Parallel Global Optimization based on Gaussian Processes. Technical report, March 2008. URL <https://hal.science/hal-00260579>.
- [21] Robert B Gramacy, Annie Sauer, and Nathan Wycoff. Triangulation candidates for bayesian optimization. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 35933–35945. Curran Associates, Inc., 2022.
- [22] Nikolaus Hansen, Sibylle D. Müller, and Petros Koumoutsakos. Reducing the time complexity of the derandomized evolution strategy with covariance matrix adaptation (cma-es). *Evolutionary Computation*, 11(1):1–18, 2003. doi: 10.1162/106365603321828970.
- [23] Tim Head, Manoj Kumar, Holger Nahrstaedt, Gilles Louppe, and Iaroslav Shcherbatyi. scikit-optimize/scikit-optimize, October 2021. URL <https://doi.org/10.5281/zenodo.5565057>.
- [24] José Miguel Henrández-Lobato, Matthew W. Hoffman, and Zoubin Ghahramani. Predictive entropy search for efficient global optimization of black-box functions. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 1*, NIPS’14, pages 918–926, Cambridge, MA, USA, 2014. MIT Press.
- [25] Carl Hvarfner, Frank Hutter, and Luigi Nardi. Joint entropy search for maximally-informed bayesian optimization. In *Advances in Neural Information Processing Systems 35*, 2022.
- [26] F. Irshad, S. Karsch, and A. Döpp. Multi-objective and multi-fidelity bayesian optimization of laser-plasma acceleration. *Phys. Rev. Res.*, 5:013063, Jan 2023. doi: 10.1103/PhysRevResearch.5.013063. URL <https://link.aps.org/doi/10.1103/PhysRevResearch.5.013063>.
- [27] Faran Irshad, Stefan Karsch, and Andreas Doepp. Reference dataset of multi-objective and multi-fidelity optimization in laser-plasma acceleration, January 2023. URL <https://doi.org/10.5281/zenodo.7565882>.
- [28] Shali Jiang, Daniel Jiang, Maximilian Balandat, Brian Karrer, Jacob Gardner, and Roman Garnett. Efficient nonmyopic bayesian optimization via one-shot multi-step trees. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 18039–18049. Curran Associates, Inc., 2020.
- [29] Donald Jones, C. Perttunen, and B. Stuckman. Lipschitzian optimisation without the lipschitz constant. *Journal of Optimization Theory and Applications*, 79:157–181, 01 1993. doi: 10.1007/BF00941892.

- [30] 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.
- [31] Kirthivasan Kandasamy, Karun Raju Vysyaraju, Willie Neiswanger, Biswajit Paria, Christopher R. Collins, Jeff Schneider, Barnabás Póczos, and Eric P. Xing. Tuning hyperparameters without grad students: Scalable and robust bayesian optimisation with dragonfly. *J. Mach. Learn. Res.*, 21(1), jan 2020.
- [32] Jungtaek Kim, Seungjin Choi, and Minsu Cho. Combinatorial bayesian optimization with random mapping functions to convex polytopes. In James Cussens and Kun Zhang, editors, *Proceedings of the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence*, volume 180 of *Proceedings of Machine Learning Research*, pages 1001–1011. PMLR, 01–05 Aug 2022.
- [33] Diederik P Kingma and Max Welling. Auto-Encoding Variational Bayes. *arXiv e-prints*, page arXiv:1312.6114, Dec 2013.
- [34] Dieter Kraft. A software package for sequential quadratic programming. *Forschungsbericht-Deutsche Forschungs- und Versuchsanstalt für Luft- und Raumfahrt*, 1988.
- [35] H. J. Kushner. A New Method of Locating the Maximum Point of an Arbitrary Multipipe Curve in the Presence of Noise. *Journal of Basic Engineering*, 86(1):97–106, 03 1964.
- [36] Armin Lederer, Jonas Umlauft, and Sandra Hirche. Posterior variance analysis of gaussian processes with application to average learning curves. *arXiv preprint arXiv:1906.01404*, 2019.
- [37] Eric Hans Lee, Valerio Perrone, Cedric Archambeau, and Matthias Seeger. Cost-aware Bayesian Optimization. *arXiv e-prints*, page arXiv:2003.10870, March 2020.
- [38] Benjamin Letham, Brian Karrer, Guilherme Ottoni, and Eytan Bakshy. Constrained bayesian optimization with noisy experiments. *Bayesian Analysis*, 14(2):495–519, 06 2019. doi: 10.1214/18-BA1110.
- [39] Qiaohao Liang, Aldair E. Gongora, Zekun Ren, Armi Tiisonen, Zhe Liu, Shijing Sun, James R. Deneault, Daniil Bash, Flore Mekki-Berrada, Saif A. Khan, Kedar Hippalgaonkar, Benji Maruyama, Keith A. Brown, John Fisher III, and Tonio Buonassisi. Benchmarking the performance of bayesian optimization across multiple experimental materials science domains. *npj Computational Materials*, 7(1):188, 2021.
- [40] Xingtao Liao, Qing Li, Xujing Yang, Weigang Zhang, and Wei Li. Multiobjective optimization for crash safety design of vehicles using stepwise regression model. *Structural and Multidisciplinary Optimization*, 35:561–569, 06 2008. doi: 10.1007/s00158-007-0163-x.
- [41] Mark McLeod, Stephen Roberts, and Michael A. Osborne. Optimization, fast and slow: optimally switching between local and Bayesian optimization. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 3443–3452. PMLR, 10–15 Jul 2018.
- [42] Jonas Moćkus. On bayesian methods for seeking the extremum. In *Optimization Techniques IFIP Technical Conference: Novosibirsk, July 1–7, 1974*, pages 400–404. Springer, 1975.
- [43] Jonas Mockus. The application of bayesian methods for seeking the extremum. *Towards global optimization*, 2:117–129, 1978.
- [44] Hossein Mohammadi, Rodolphe Le Riche, and Eric Touboul. Making ego and cma-es complementary for global optimization. In Clarisse Dhaenens, Laetitia Jourdan, and Marie-Éléonore Marmion, editors, *Learning and Intelligent Optimization*, pages 287–292, Cham, 2015. Springer International Publishing.
- [45] Henry B. Moss, David S. Leslie, Javier González, and Paul Rayson. Gibbon: General-purpose information-based bayesian optimisation. *J. Mach. Learn. Res.*, 22(1), jan 2021.

- [46] Changyong Oh, Jakub Tomczak, Efstratios Gavves, and Max Welling. Combinatorial bayesian optimization using the graph cartesian product. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 2914–2924. Curran Associates, Inc., 2019.
- [47] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [48] Victor Picheny, Joel Berkeley, Henry B. Moss, Hrvoje Stojic, Uri Granta, Sebastian W. Ober, Artem Artemev, Khurram Ghani, Alexander Goodall, Andrei Paleyes, Sattar Vakili, Sergio Pascual-Diaz, Stratis Markou, Jixiang Qing, Nasrulloh R. B. S Loka, and Ivo Couckuyt. Trieste: Efficiently exploring the depths of black-box functions with tensorflow, 2023. URL <https://arxiv.org/abs/2302.08436>.
- [49] Carl Edward Rasmussen. *Gaussian Processes in Machine Learning*, pages 63–71. Springer Berlin Heidelberg, Berlin, Heidelberg, 2004.
- [50] 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.
- [51] S. Surjanovic and D. Bingham. Virtual library of simulation experiments: Test functions and datasets. Retrieved May 14, 2023, from <http://www.sfu.ca/~ssurjano>.
- [52] Ryoji Tanabe and Hisao Ishibuchi. An easy-to-use real-world multi-objective optimization problem suite. *Applied Soft Computing*, 89:106078, 2020. ISSN 1568-4946. doi: <https://doi.org/10.1016/j.asoc.2020.106078>.
- [53] The GPyOpt authors. GPyOpt: A bayesian optimization framework in python. <http://github.com/SheffieldML/GPyOpt>, 2016.
- [54] Aimo Törn and Antanas Zilinskas. *Global optimization*, volume 350. Springer, 1989.
- [55] Ben Tu, Axel Gandy, Nikolas Kantas, and Behrang Shafei. Joint entropy search for multi-objective bayesian optimization. In *Advances in Neural Information Processing Systems 35*, 2022.
- [56] Andreas Wächter and Lorenz T. Biegler. On the implementation of an interior-point filter line-search algorithm for large-scale nonlinear programming. *Mathematical Programming*, 106(1):25–57, 2006.
- [57] Xingchen Wan, Vu Nguyen, Huong Ha, Binxin Ru, Cong Lu, and Michael A. Osborne. Think global and act local: Bayesian optimisation over high-dimensional categorical and mixed search spaces. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 10663–10674. PMLR, 18–24 Jul 2021.
- [58] Jialei Wang, Scott C. Clark, Eric Liu, and Peter I. Frazier. Parallel bayesian global optimization of expensive functions, 2016.
- [59] Zi Wang and Stefanie Jegelka. Max-value Entropy Search for Efficient Bayesian Optimization. *ArXiv e-prints*, page arXiv:1703.01968, March 2017.
- [60] James Wilson, Frank Hutter, and Marc Deisenroth. Maximizing acquisition functions for bayesian optimization. In *Advances in Neural Information Processing Systems 31*, pages 9905–9916. 2018.

- 506 [61] Jian Wu and Peter I. Frazier. The parallel knowledge gradient method for batch bayesian  
507 optimization. In *Proceedings of the 30th International Conference on Neural Information*  
508 *Processing Systems*, NIPS'16, page 3134–3142, Red Hook, NY, USA, 2016. Curran Associates  
509 Inc. ISBN 9781510838819.
- 510 [62] Jian Wu, Matthias Poloczek, Andrew Gordon Wilson, and Peter I Frazier. Bayesian optimization  
511 with gradients. In *Advances in Neural Information Processing Systems*, pages 5267–5278, 2017.
- 512 [63] Kaifeng Yang, Michael Emmerich, André Deutz, and Thomas Bäck. Multi-objective bayesian  
513 global optimization using expected hypervolume improvement gradient. *Swarm and Evolution-*  
514 *ary Computation*, 44:945 – 956, 2019.
- 515 [64] Dawei Zhan and Huanlai Xing. Expected improvement for expensive optimization: a review.  
516 *Journal of Global Optimization*, 78(3):507–544, 2020.
- 517 [65] Jiayu Zhao, Renyu Yang, Shenghao Qiu, and Zheng Wang. Enhancing high-dimensional  
518 bayesian optimization by optimizing the acquisition function maximizer initialization. *arXiv*  
519 *preprint arXiv:2302.08298*, 2023.
- 520 [66] Eckart Zitzler, Kalyanmoy Deb, and Lothar Thiele. Comparison of multiobjective evolutionary  
521 algorithms: Empirical results. *Evol. Comput.*, 8(2):173–195, jun 2000. ISSN 1063-6560. doi:  
522 10.1162/106365600568202. URL <https://doi.org/10.1162/106365600568202>.

## 523 A Acquisition Function Details

### 524 A.1 Analytic Expected Improvement

525 Recall that the main challenge with computing analytic LogEI is to accurately compute  $\log h$ , where  
 526  $h(z) = \phi(z) + z\Phi(z)$ , with  $\phi(z) = \exp(-z^2/2)/\sqrt{2\pi}$  and  $\Phi(z) = \int_{-\infty}^z \phi(u)du$ . To express  $\log h$   
 527 in a numerically stable form as  $z$  becomes increasingly negative, we first take the log and multiply  $\phi$   
 528 out of the argument to the logarithm:

$$\log h(z) = z^2/2 - \log(2\pi)/2 + \log\left(1 + z \frac{\Phi(z)}{\phi(z)}\right). \quad (12)$$

529 Fortunately, this form exposes the quadratic factor, and  $\Phi(z)/\phi(z)$  can be computed via standard  
 530 implementations of the scaled complementary error function `erfcx`. However, even `erfcx` can give  
 531 rise to numerical underflow (though after a significantly larger range than the original formulation).  
 532 To further increase numerical stability, we instead compute the  $\log(\Phi(z)/\phi(z))$  using

$$\text{logerfcx}(x) = \begin{cases} \log(\text{erfc}(x)) + x^2 & x < 0 \\ \log(\text{erfcx}(x)) & x \geq 0 \end{cases} \quad (13)$$

533 and afterward use

$$\text{log1mexp}(x) = \begin{cases} \log(-\text{expm1}(x)) & -\log 2 < x \\ \log1p(-\exp(x)) & -\log 2 \geq x \end{cases} \quad (14)$$

534 to compute the last term of  $\log h$  in Eq. (12). In particular, we compute

$$\text{log\_h}(z) = \begin{cases} \log(\phi(z) + z\Phi(z)) & z > -1 \\ -z^2/2 - c_1 + \text{log1mexp}(\text{logerfcx}(-z/\sqrt{2})|z| + c_2) & z < -1 \end{cases} \quad (15)$$

535 where  $c_1 = \log(2\pi)/2$ , and  $c_2 = \log(\pi/2)/2$ .

536 Figure 8 shows both the numerical failure mode of a naïve implementation of EI, which becomes  
 537 *exactly* zero numerically for moderately small  $z$ , while the evaluation via  $\text{log\_h}$  in Eq. (15) exhibits  
 538 quadratic asymptotic behavior that is particularly amenable to numerical optimization routines.

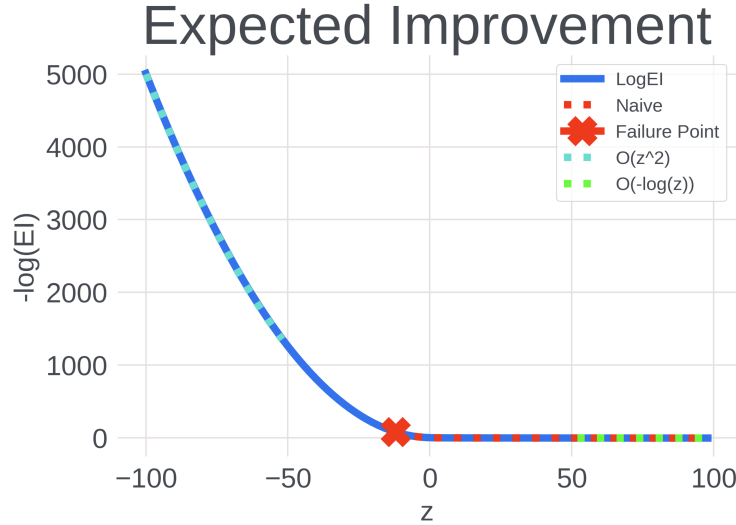


Figure 8: Plot of the  $\log h$ , computed via  $\log \circ h$  and  $\text{log\_h}$  in Eq. (15). Crucially, the naïve implementation fails as  $z = (\mu(x) - f^*)/\sigma(x)$  becomes increasingly negative, due to being exactly numerically zero, while our proposed implementation exhibits quadratic asymptotic behavior.

## 539 A.2 Monte-Carlo Expected Improvement

540 For Monte-Carlo, we cannot directly apply similar numerical improvements as for the analytical  
 541 version, because the utility values (integrand of Eq. (4)) on the sample level are likely to be *mathemat-*  
 542 *ically* zero. For this reason, we first smoothly approximate the acquisition utility and subsequently  
 543 apply log transformations to the approximate acquisition function.

544 To this end, a natural choice is  $\text{softplus}_{\tau_0}(x) = \tau_0 \log(1 + \exp(x/\tau_0))$  for smoothing the  $\max(0, x)$ ,  
 545 where  $\tau_0$  is a temperature parameter governing the approximation error. Further, we approximate the  
 546  $\max_i$  over the  $q$  candidates by the norm  $\|\cdot\|_{1/\tau_{\max}}$  and note that the approximation error introduced  
 547 by both smooth approximations can be bound tightly as a function of two “temperature” parameters  
 548  $\tau_0$  and  $\tau_{\max}$ , see Lemma 2.

549 Importantly, the smoothing alone only solves the problem of having mathematically zero gradients,  
 550 not that of having numerically vanishing gradients, as we have shown for the analytical case above.  
 551 For this reason, we transform all smoothed computations to log space and thus need the following  
 552 special implementation of  $\log \circ \text{softplus}$  that can be evaluated stably for a very large range of inputs:

$$\text{logsoftplus}_{\tau}(x) = \begin{cases} [\log \circ \text{softplus}_{\tau}](x) & x/\tau > l \\ x/\tau + \log(\tau) & x/\tau \leq l \end{cases}$$

553 where  $\tau$  is a temperature parameter and  $l$  depends on the floating point precision used, around  $-35$   
 554 for double precision in our implementation.

555 Note that the lower branch of  $\text{logsoftplus}$  is approximate. Using a Taylor expansion of  $\log(1+z) =$   
 556  $z - z^2/2 + \mathcal{O}(z^3)$  around  $z = 0$ , we can see that the approximation error is  $\mathcal{O}(z^2)$ , and therefore,  
 557  $\log(\log(1 + \exp(x))) = x + \mathcal{O}(\exp(x)^2)$ , which converges to  $x$  exponentially quickly.  $l$  is chosen  
 558 in our implementation so that no significant digit is lost in dropping the second order term from the  
 559 lower branch.

560 Having defined  $\text{logsoftplus}$ , we further note that

$$\begin{aligned} \log \|\mathbf{x}\|_{1/\tau_{\max}} &= \log \left( \sum_i x_i^{1/\tau_{\max}} \right)^{\tau_{\max}} \\ &= \tau_{\max} \log \left( \sum_i \exp(\log(x_i)/\tau_{\max}) \right) \\ &= \tau_{\max} \text{logsumexp}_i(\log(x_i)/\tau_{\max}) \end{aligned}$$

561 Therefore, we can express the logarithm of the smoothed acquisition utility per sample as

$$\tau_{\max} \text{logsumexp}_i(\text{logsoftplus}_{\tau_0}(z_i)/\tau_{\max})$$

562 Applying another  $\text{logsumexp}$  to compute the logarithm of the mean of acquisition utilities over a set  
 563 of Monte Carlo samples gives rise to the expression in Eq. (10).

564 In particular for large batches (large  $q$ ), this expression can still give rise to vanishing gradients for  
 565 some candidates, which is due to the large dynamic range of the outputs of the  $\text{logsoftplus}$  when  
 566  $x \ll 0$ . To solve this problem, we propose a new class of smooth approximations to the “hard”  
 567 non-linearities that decay as  $\mathcal{O}(1/x^2)$  as  $x \rightarrow -\infty$  in the next section.

## 568 A.3 A Class of Smooth Approximations with Fat Tails for Larger Batches

569 A regular  $\text{softplus}(x) = \log(1 + \exp(x))$  function smoothly approximates the ReLU non-linearity  
 570 and – in conjunction with the log transformations – is sufficient to achieve good numerical behavior  
 571 for small batches of the Monte Carlo acquisition functions. However, as more candidates are added,  
 572  $\log \text{softplus}(x) = \log(\log(1 + \exp(x)))$  is increasingly likely to have a high dynamic range as for  
 573  $x \ll 0$ ,  $\log \text{softplus}_{\tau}(x) \sim -x/\tau$ . If  $\tau > 0$  is chosen to be small,  $(-x/\tau)$  can vary orders of  
 574 magnitude within a single batch. This becomes problematic when we approximate the maximum  
 575 utility over the batch of candidates, since  $\text{logsumexp}$  only propagates numerically non-zero gradients  
 576 to inputs that are no smaller than approximately  $(\max_j x_j - 700)$  in double precision, another source  
 577 of vanishing gradients.



To solve this problem, we propose a new smooth approximation to the ReLU, maximum, and indicator functions that decay only polynomially as  $x \rightarrow -\infty$ , instead of exponentially, like the canonical softplus. The high level idea is to use  $(1 + x^2)^{-1}$ , which is proportional to the Cauchy density function (and is also known as a Lorentzian), in ways that maintain key properties of existing smooth approximations – convexity, positivity, etc – while changing the asymptotic behavior of the functions from exponential to  $\mathcal{O}(1/x^2)$  as  $x \rightarrow -\infty$ , also known as a “fat tail”. Further, we will show that the proposed smooth approximations satisfy similar maximum error bounds as their exponentially decaying counterparts, thereby permitting a similar approximation guarantee as Lemma 2 with minor adjustments to the involved constants.

**Fat Softplus** We define

$$\varphi_+(x) = \alpha(1 + x^2)^{-1} + \log(1 + \exp(x)), \quad (16)$$

for a positive scalar  $\alpha$ . The following result shows that we can ensure the monotonicity and convexity – both important properties of the ReLU that we would like to maintain in our approximation – of  $g$  by carefully choosing  $\alpha$ .

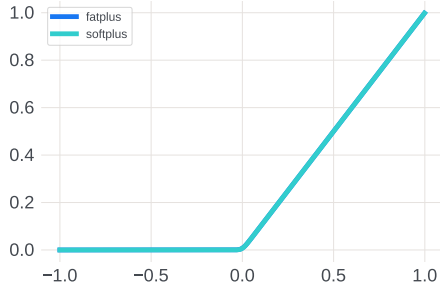


Figure 9: The fat softplus approximates  $\max(x, 0)$  similarly tightly as the regular softplus and is also monotonic, convex, and positive. The plot used a temperature of  $\tau_0 = 0.01$ .

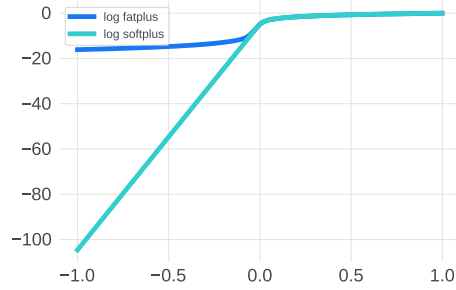


Figure 10: The fat softplus has an  $\mathcal{O}(1/x^2)$  asymptotic decay, versus the  $\mathcal{O}(\exp(x))$  decay as  $x \rightarrow -\infty$ , moderating the dynamic range of the quantities involved in parallel LogEI.

**Lemma 3** (Monotonicity and Convexity).  $\varphi_+(x)$  is positive, monotonically increasing, and strictly convex for  $\alpha$  satisfying

$$0 \leq \alpha < \frac{e^{1/\sqrt{3}}}{2(1 + e^{1/\sqrt{3}})}.$$

*Proof.* Positivity follows due to  $\alpha \geq 0$ , and both sumands being positive. Monotonicity and convexity can be shown via canonical differential calculus and bounding relevant quantities.

In particular, regarding monotonicity, we want to select  $\alpha$  so that the first derivative is bounded below by zero:

$$\partial_x \varphi_+(x) = \frac{e^x}{1 + e^x} - \alpha \frac{2x}{(1 + x^2)^2}$$

First, we note that  $\partial_x \varphi_+(x)$  is positive for  $x < 0$  and any  $\alpha$ , since both terms are positive in this regime. For  $x \geq 0$ ,  $\frac{e^x}{1 + e^x} = (1 + e^{-x})^{-1} \geq 1/2$ , and  $-1/(1 + x^2)^2 \geq -1/(1 + x^2)$ , so that

$$\partial_x \varphi_+(x) \geq \frac{1}{2} - \alpha \frac{2x}{(1 + x^2)}$$

Forcing  $\frac{1}{2} - \alpha \frac{2x}{(1 + x^2)} > 0$ , and multiplying by  $(1 + x^2)^2$  gives rise to a quadratic equation whose roots are  $x = 2\alpha \pm \sqrt{4\alpha^2 - 1}$ . Thus, there are no real roots for  $\alpha < 1/2$ . Since the derivative is certainly positive for the negative reals and the guaranteed non-existence of roots implies that the derivative cannot cross zero elsewhere,  $0 \leq \alpha < 1/2$  is a sufficient condition for monotonicity of  $\varphi_+$ .

Regarding convexity, our goal is to prove a similar condition on  $\alpha$  that guarantees the positivity of the second derivative:

$$\partial_x^2 \varphi_+(x) = \alpha \frac{6x^2 - 2}{(1+x^2)^3} + \frac{e^{-x}}{(1+e^{-x})^2}$$

Note that  $\frac{6x^2-2}{(1+x^2)^3}$  is symmetric around 0, is negative in  $(-\sqrt{1/3}, \sqrt{1/3})$  and has a minimum of  $-2$  at 0.  $\frac{e^{-x}}{(1+e^{-x})^2}$  is symmetric around zero and decreasing away from zero. Since the rational polynomial is only negative in  $(-\sqrt{1/3}, \sqrt{1/3})$ , we can lower bound  $\frac{e^{-x}}{(1+e^{-x})^2} > \frac{e^{-\sqrt{1/3}}}{(1+e^{-\sqrt{1/3}})^2}$  in  $(-\sqrt{1/3}, \sqrt{1/3})$ . Therefore,

$$\partial_x^2 \varphi_+(x) \geq \frac{e^{-x}}{(1+e^{-x})^2} - 2\alpha$$

Forcing  $\frac{e^{-\sqrt{1/3}}}{(1+e^{-\sqrt{1/3}})^2} - 2\alpha > 0$  and rearranging yields the result. Since  $\frac{e^{-\sqrt{1/3}}}{(1+e^{-\sqrt{1/3}})^2} / 2 \sim 0.115135$ , the convexity condition is stronger than the monotonicity condition and therefore subsumes it.  $\square$

Importantly  $\varphi$  decays only polynomially for increasingly negative inputs, and therefore  $\log \varphi$  only logarithmically, which keeps the range of  $\varphi$  constrained to values that are more manageable numerically. Similar to Lemma 5, one can show that

$$|\tau \varphi_+(x/\tau) - \text{ReLU}(x)| \leq (\alpha + \log(2)) \tau. \quad (17)$$

There are a large number of approximations or variants of the ReLU that have been proposed as activation functions of artificial neural networks, but to our knowledge, none satisfy the properties that we seek here: (1) smoothness, (2) positivity, (3) monotonicity, (4) convexity, and (5) polynomial decay. For example, the leaky ReLU does not satisfy (1) and (2), and the ELU does not satisfy (5).

**Fat Maximum** The canonical `logsumexp` approximation to  $\max_i x_i$  suffers from numerically vanishing gradients if  $\max_i x_i - \min_j x_j$  is larger a moderate threshold, around 760 in double precision, depending on the floating point implementation. In particular, while elements close to the maximum receive numerically non-zero gradients, elements far away are increasingly likely to have a numerically zero gradient. To fix this behavior for the smooth maximum approximation, we propose

$$\varphi_{\max}(\mathbf{x}) = \max_j x_j + \tau \log \sum_i \left[ 1 + \left( \frac{x_i - \max_j x_j}{\tau} \right)^2 \right]^{-1}. \quad (18)$$

This approximation to the maximum has the same error bound to the true maximum as the `logsumexp` approximation:

**Lemma 4.** Given  $\tau > 0$

$$\max_i x_i \leq \tau \phi_{\max}(x/\tau) \leq \max_i x_i + \tau \log(d). \quad (19)$$

*Proof.* Regarding the lower bound, by definition there is an index  $i$  such that  $x_i = \max_j x_j$ . For this index, the associated summand in (18) is 1. Since all summands are positive, the entire sum is lower bounded by 1, hence

$$\tau \log \sum_i \left[ 1 + \left( \frac{x_i - \max_j x_j}{\tau} \right)^2 \right]^{-1} > \tau \log(1) = 0$$

Adding  $\max_j x_j$  to the inequality finishes the proof for the lower bound.

Regarding the upper bound, (18) can be maximized when  $x_i = \max_j x_j$  for all  $i$ , in which case each  $(x_i - \max_j x_j)^2$  is minimized, and hence each summand is maximized. In this case,

$$\tau \log \sum_i \left[ 1 + \left( \frac{x_i - \max_j x_j}{\tau} \right)^2 \right]^{-1} \leq \tau \log \left( \sum_i 1 \right) = \tau \log(d).$$

Adding  $\max_j x_j$  to the inequality finishes the proof for the upper bound.  $\square$

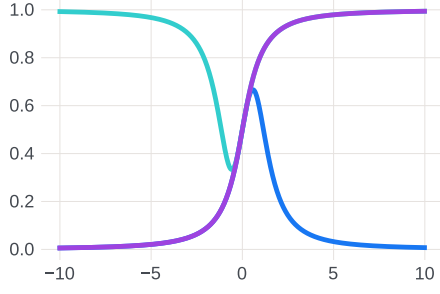


Figure 11: We construct the fat sigmoid approximation (purple) by splicing together two Lorentzians (blue and teal) at one of their inflection points.

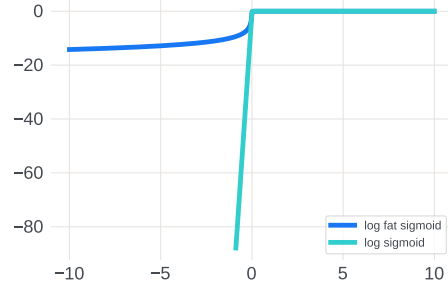


Figure 12: The fat sigmoid approximation ( $\tau = 0.01$ ) decays as  $\mathcal{O}(1/x^2)$  instead of  $\mathcal{O}(\exp(x))$  as  $x \rightarrow -\infty$ , minimizing the dynamic range of the numerical quantities in constrained qLogEI.

**Fat Sigmoid** Notably, we encountered a similar problem with using regular (log)-sigmoids to smooth the constraint indicators for EI with black-box constraints. Here, we want the smooth approximation  $\iota$  to satisfy 1) positivity, 2) monotonicity, 3) polynomial decay, and 4)  $\iota(x) = 1/2 - \iota(-x)$ . Let  $\gamma = \sqrt{1/3}$ , then we define

$$\iota(x) = \begin{cases} \frac{2}{3} \left(1 + (x - \gamma)^2\right)^{-1} & x < 0, \\ 1 - \frac{2}{3} \left(1 + (x + \gamma)^2\right)^{-1} & x \geq 0. \end{cases}$$

$\iota$  is monotonically increasing, satisfies  $\iota(x) \rightarrow 1$  as  $x \rightarrow \infty$ ,  $\iota(0) = 1/2$ , and  $\iota(x) = \mathcal{O}(1/x^2)$  as  $x \rightarrow -\infty$ . Further, we note that the asymptotics are primarily important here, but that we can also make the approximation tighter by introducing a temperature parameter  $\tau$ , and letting  $\iota_\tau(x) = \iota(x/\tau)$ . The approximation error of  $\iota_\tau(x)$  to the Heaviside step function becomes tighter point-wise as  $\tau \rightarrow 0+$ , except for at the origin where  $\iota_\tau(x) = 1/2$ , similar to the canonical sigmoid.

#### A.4 Constrained Expected Improvement

For the analytical case, many computational frameworks already provide a numerically stable implementation of the logarithm of the Gaussian cumulative distribution function, in the case of PyTorch, `torch.special.log_ndtr`, which can be readily used in conjunction with our implementation of LogEI, as described in Sec. 4.3.

For the case of Monte-Carlo parallel EI, we implemented the fat-tailed  $\iota$  function from Sec. A.3 to approximate the constraint indicator and compute the per-candidate, per-sample acquisition utility using

$$(\text{logsoftplus}_{\tau_0}(\xi_i(x_j) - y^*) + \sum_k \log \circ \iota \left( -\frac{\xi_i^{(k)}}{\tau_{\text{cons}}} \right),$$

where  $\xi_i^{(k)}$  is the  $i$ th sample of the  $k$ th constraint model, and  $\tau_{\text{cons}}$  is the temperature parameter controlling the approximation to the constraint indicator. While this functionality is in our implementation, our benchmark results use the analytical version.

#### A.5 Parallel Expected Hypervolume Improvement

The hypervolume improvement can be computed via the inclusion-exclusion principle, see [8] for details, we focus on the numerical issues concerning qEHVI here. To this end, we define

$$z_{k,i_1,\dots,i_j}^{(m)} = \min [\mathbf{u}_k, \mathbf{f}(\mathbf{x}_{i_1}), \dots, \mathbf{f}(\mathbf{x}_{i_j})],$$

where  $\mathbf{f}$  is the vector-valued objective function, and  $\mathbf{u}_k$  is the vector of upper bounds of one of  $K$  hyper-rectangles that partition the non-Pareto-dominated space, see [8] for details on the partitioning.

Letting  $l_k$  be the corresponding lower bounds of the hyper-rectangles, the hypervolume improvement can then be computed as

$$\text{HVI}(\{\mathbf{f}(\mathbf{x}_i)\}_{i=1}^q) = \sum_{k=1}^K \sum_{j=1}^q \sum_{X_j \in \mathcal{X}_j} (-1)^{j+1} \prod_{m=1}^M [z_{k,X_j}^{(m)} - l_k^{(m)}]_+, \quad (20)$$

where  $\mathcal{X}_j = \{X_j \subset \mathcal{X}_{\text{cand}} : |X_j| = j\}$  is the superset of all subsets of  $\mathcal{X}_{\text{cand}}$  of size  $j$  and  $z_{k,X_j}^{(m)} = z_{k,i_1,\dots,i_j}^{(m)}$  for  $X_j = \{\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_j}\}$ .

To find a numerically stable formulation of the logarithm of this expression, we first re-purpose the  $\varphi_{\max}$  function to compute the minimum in the expression of  $z_{k,i_1,\dots,i_j}^{(m)}$ , like so  $\varphi_{\min}(x) = -\varphi_{\max}(-x)$ . Further, we use the  $\varphi_+$  function of Sec. A.3 as for the single objective case to approximate  $[z_{k,X_j}^{(m)} - l_k^{(m)}]_+$ . We then have

$$\log \prod_{m=1}^M \varphi_+[z_{k,X_j}^{(m)} - l_k^{(m)}] = \sum_{m=1}^M \log \varphi_+[z_{k,X_j}^{(m)} - l_k^{(m)}] \quad (21)$$

Since we can only transform positive quantities to log space, we split the sum in Eq. (20) into positive and negative components, depending on the sign of  $(-1)^{j+1}$ , and compute the result using a numerically stable implementation of  $\log(\exp(\log \text{ of positive terms}) - \exp(\log \text{ of negative terms}))$ . The remaining sums over  $k$  and  $q$  can be carried out by applying `logsumexp` to the resulting quantity. Finally, applying `logsumexp` to reduce over an additional Monte-Carlo sample dimension yields the formulation of `qLogEHVI` that we use in our multi-objective benchmarks.

## B Strategies for Optimizing Acquisition Functions

As discussed in Section 2.3, a variety of different approaches and heuristics have been applied to the problem of optimizing acquisition functions. For the purpose of this work, we only consider continuous domains  $\mathbb{X}$ . While discrete and/or mixed domains are also relevant in practice and have received substantial attention in recent years (see e.g. Baptista and Poloczec [3], Daulton et al. [9], Deshwal et al. [11], Kim et al. [32], Oh et al. [46], Wan et al. [57]), our work here on improving acquisition functions is largely orthogonal to this (though the largest gains should be expected when using gradient-based optimizers, as is done in mixed-variable BO when conditioning on discrete variables, or when performing discrete or mixed BO using continuous relaxations, probabilistic reparameterization, or straight-through estimators [9]).

Arguably the simplest approach to optimizing acquisition functions is by grid search or random search. While variants of this combined with local descent can make sense in the context of optimizing over discrete or mixed spaces and when acquisition functions can be evaluated efficiently in batch (e.g. on GPUs), this clearly does not scale to higher-dimensional continuous domains due to the exponential growth of space to cover.

Another relatively straightforward approach is to use zeroth-order methods such as DIRECT [29] (used e.g. by Dragonfly [31]) or the popular CMA-ES [22]. These approaches are easy to implement as they avoid the need to compute gradients of acquisition functions. However, not relying on gradients is also what renders their optimization performance inferior to gradient based methods, especially for higher-dimensional problems and/or joint batch optimization in parallel Bayesian optimization.

The most common approach to optimizing acquisition functions on continuous domains is using gradient descent-type algorithms. Gradients are either computed based on analytically derived closed-form expressions, or via auto-differentiation capabilities of modern ML systems such as PyTorch [47], Tensorflow [1], or JAX [5].

For analytic acquisition functions, a common choice of optimizer is L-BFGS-B [6], a quasi-second order method that uses gradient information to approximate the Hessian and supports box constraints. If other, more general constraints are imposed on the domain, other general purpose nonlinear optimizers such as SLSQP [34] or IPOPT [56] are used (e.g. by BoTorch). For Monte Carlo (MC) acquisition functions, Wilson et al. [60] proposes using stochastic gradient ascent (SGA) based on stochastic gradient estimates obtained via the reparameterization trick [33]. Stochastic first-order

algorithms are also used by others, including e.g. Wang et al. [58] and Daulton et al. [9]. Balandat et al. [2] build on the work by Wilson et al. [60] and show how sample average approximation (SAA) can be employed to obtain deterministic gradient estimates for MC acquisition functions, which has the advantage of being able to leverage the improved convergence rates of optimization algorithms designed for deterministic functions such as L-BFGS-B. This general approach has since been used for a variety of other acquisition functions, including e.g. Daulton et al. [8] and Jiang et al. [28].

Very few implementations of Bayesian Optimization actually use higher-order derivative information, as this either requires complex derivations of analytical expressions and their custom implementation, or computation of second-order derivatives via automated differentiation, which is less well supported and computationally much more costly than computing only first-order derivatives. One notable exception is `Corne11-MOE` [61, 62], which supports Newton’s method (though this is limited to the acquisition functions implemented in C++ within the library and not easily extensible to other acquisition functions).

## B.1 Common initialization heuristics for multi-start gradient-descent

One of the key issues to deal with gradient-based optimization in the context of optimizing acquisition functions is the optimizer getting stuck in local optima due to the generally highly non-convex objective. This is typically addressed by means of restarting the optimizer from a number of different initial conditions distributed across the domain.

A variety of different heuristics have been proposed for this. The most basic one is to restart from random points uniformly sampled from the domain (for instance, `scikit-optimize` [23] uses this strategy). However, as we have argued in this paper, acquisition functions can be (numerically) zero in large parts of the domain, and so purely random restarts can become ineffective, especially in higher dimensions and with more data points. A common strategy is therefore to either augment or bias the restart point selection to include initial conditions that are closer to “promising points”. `GPYOpt` [53] augments random restarts with the best points observed so far, or alternatively points generated via Thompson sampling. `SpearMint` [50] initializes starting points based on Gaussian perturbations of the current best point. `BoTorch` [2] selects initial points by performing Boltzmann sampling on a set of random points according to their acquisition function value; the goal of this strategy is to achieve a biased random sampling across the domain that is likely to generate more points around regions with high acquisition value, but remains asymptotically space-filling. The initialization strategy used by `Trieste` [48] works similarly to the one in `BoTorch`, but instead of using soft-randomization via Boltzmann sampling, it simply selects the top- $k$  points. Most recently, Gramacy et al. [21] proposed distributing initial conditions using a Delaunay triangulation of previously observed data points. This is an interesting approach that generalizes the idea of initializing “in between” observed points from the single-dimensional case. However, this approach does not scale well with the problem dimension and the number of observed data points due to the complexity of computing the triangulation (with wall time empirically found to be exponential in the dimension, see [21, Fig. 3] and worst-case quadratic in the number of observed points).

However, while these initialization strategies can help substantially with better optimizing acquisition functions, they ultimately cannot resolve foundational issues with acquisition functions themselves. Ensuring that acquisition functions provides enough gradient information (not just mathematically but also numerically) is therefore key to be able to optimize it effectively, especially in higher dimensions and with more observed data points.

## C Proofs

**Lemma 1.** *Suppose  $f$  is drawn from a Gaussian process prior  $P_f$ ,  $y^* \leq f^*$ ,  $\mu_n, \sigma_n$  are the mean and standard deviation of the posterior  $P_f(f|\mathcal{D}_n)$  and  $B \in \mathbb{R}$ . Then with probability  $1 - \delta$ ,*

$$P_x \left( \frac{\mu_n(x) - y_n^*}{\sigma_n(x)} < B \right) \geq P_x (f(x) < f^* - \epsilon_n) \quad (7)$$

where  $\epsilon_n = (f^* - y_n^*) + (\sqrt{-2 \log(2\delta)} - B) \max_x \sigma_n(x)$ .

*Proof.*

$$\frac{\mu_n(x) - y^*}{\sigma_n(x)} = \frac{\mu_n(x) - f(x)}{\sigma_n(x)} + \frac{f(x) - f^*}{\sigma_n(x)} \quad (22)$$

We proceed by bounding the first term on the right hand side. Note that by assumption,  $f(x) \sim \mathcal{N}(\mu_n(x), \sigma_n(x)^2)$  and thus  $(\mu_n(x) - f(x))/\sigma_n(x) \sim \mathcal{N}(0, 1)$ . For a positive  $C > 0$  then, we use a standard bound on the Gaussian tail probability to attain

$$P\left(\frac{\mu_n(x) - f(x)}{\sigma_n(x)} > C\right) \leq e^{-C^2/2}/2. \quad (23)$$

Therefore,  $(\mu(x) - f(x))/\sigma_n(x) < C$  with probability  $1 - \delta$  if  $C = \sqrt{-2\log(2\delta)}$ .

Using the bound just derived, and forcing the resulting upper bound to be less than  $B$  yields a sufficient condition to imply  $\mu_n(x) - y_n^* < B\sigma_n(x)$ :

$$\frac{\mu_n(x) - y^*}{\sigma_n(x)} \leq C + \frac{f(x) - y^*}{\sigma_n(x)} < B \quad (24)$$

Re-arranging and using  $y^* = f^* + (y^* - f^*)$  we get with probability  $1 - \delta$ ,

$$f(x) \leq f^* - (f^* - y_n^*) - (\sqrt{-2\log(2\delta)} - B)\sigma_n(x). \quad (25)$$

Thus, we get

$$\begin{aligned} P_x\left(\frac{\mu_n(x) - y_n^*}{\sigma_n(x)} < B\right) &\geq P_x\left(f(x) \leq f^* - (f^* - y_n^*) - (\sqrt{-2\log(2\delta)} - B)\sigma_n(x)\right) \\ &\geq P_x\left(f(x) \leq f^* - (f^* - y_n^*) - (\sqrt{-2\log(2\delta)} - B)\max_x \sigma_n(x)\right). \end{aligned} \quad (26)$$

Note that the last inequality gives a bound that is not directly dependent on the evaluation of the posterior statistics of the surrogate at any specific  $x$ . Rather, it is dependent on the optimality gap  $f^* - y_n^*$  and the maximal posterior standard deviation, or a bound thereof. Letting  $\epsilon_n = (f^* - y_n^*) - (\sqrt{-2\log(2\delta)} - B)\max_x \sigma_n(x)$  finishes the proof.  $\square$

**Lemma 2.** [Approximation Guarantee] Given the temperature parameters  $\tau_0$  and  $\tau_{\max}$ , the approximation error of qLogEI to qEI is bounded by

$$|\exp(\text{qLogEI}(x)) - \text{qEI}(x)| \leq \log(2) \tau_0 + q^{\tau_{\max}} - 1. \quad (11)$$

*Proof.* Let  $z_{iq} = \xi(x_q) - y^*$ , where  $i \in \{1, \dots, n\}$ , and for brevity of notation, and let  $\text{lse}$ ,  $\text{lsp}$  refer to the `logsumexp` and `logsoftplus` functions, respectively, and  $\text{ReLU}(x) = [x]_+$ . We then bound  $n|e^{\text{qLogEI}(x)} - \text{qEI}(x)|$  by

$$\begin{aligned} &\left| \exp(\text{lse}_i(\tau_{\max} \text{lse}_q(\text{lsp}_{\tau_0}(z_{iq})/\tau_{\max}))) - \sum_i \max_q \text{ReLU}(z_{iq}) \right| \\ &\leq \sum_i \left| \exp(\tau_{\max} \text{lse}_q(\text{lsp}_{\tau_0}(z_{iq})/\tau_{\max})) - \max_q \text{ReLU}(z_{iq}) \right| \\ &= \sum_i \left| \|\text{softplus}_{\tau_0}(z_{i\cdot})\|_{1/\tau_{\max}} - \max_q \text{ReLU}(z_{iq}) \right| \\ &\leq \sum_i \left| \|\text{softplus}_{\tau_0}(z_{i\cdot})\|_{1/\tau_{\max}} - \max_q \text{softplus}_{\tau_0}(z_{iq}) \right| \\ &\quad + \left| \max_q \text{softplus}_{\tau_0}(z_{iq}) - \max_q \text{ReLU}(z_{iq}) \right| \end{aligned} \quad (27)$$

First and second inequalities are due to the triangle inequality, where for the second we used  $|a - c| \leq |a - b| + |b - c|$  with  $b = \max_q \text{softplus}(z_{iq})$ .

768 To bound the first term in the sum, note that  $\|x\|_\infty \leq \|x\|_q \leq \|x\|_\infty d^{1/q}$ , thus  $|\|x\|_q - \|x\|_\infty| \leq$   
769  $d^{1/q} - 1$ , and therefore

$$\left| \|\text{softplus}_{\tau_0}(z_{i\cdot})\|_{1/\tau_{\max}} - \max_q \text{softplus}_{\tau_0}(z_{iq}) \right| \leq d^{\tau_{\max}} - 1.$$

770 The second term in the sum can be bound due to  $|\text{softplus}_{\tau_0}(x) - \text{ReLU}(x)| \leq \log(2)\tau_0$  (see  
771 Lemma 5 below) and therefore,

$$\left| \max_q \text{softplus}_{\tau_0}(z_{iq}) - \max_q \text{ReLU}_{\tau_0}(z_{iq}) \right| \leq \log(2)\tau_0.$$

772 Dividing Eq. (27) by  $n$  to compute the sample mean finishes the proof for the Monte-Carlo approx-  
773 imations to the acquisition value. Taking  $n \rightarrow \infty$  further proves the result for the mathematical  
774 definitions of the parallel acquisition values, i.e. Eq. (4).  $\square$

775 Approximating the ReLU using the  $\text{softplus}_\tau(x) = \tau \log(1 + \exp(x/\tau))$  function leads to an  
776 approximation error that is at most  $\tau$  in the infinity norm, i.e.  $\|\text{softplus}_\tau - \text{ReLU}\|_\infty = \log(2)\tau$ .  
777 The following lemma formally proves this.

778 **Lemma 5.** *Given  $\tau > 0$ , we have for all  $x \in \mathbb{R}$ ,*

$$|\text{softplus}_\tau(x) - \text{ReLU}(x)| \leq \log(2) \tau. \quad (28)$$

779 *Proof.* Taking the (sub-)derivative of  $\text{softplus}_\tau - \text{ReLU}$ , we get

$$\partial_x \text{softplus}_\tau(x) - \text{ReLU}(x) = (1 + e^{-x/\tau})^{-1} - \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases}$$

780 which is positive for all  $x < 0$  and negative for all  $x > 0$ , hence the extremum must be at  
781  $x$ , at which point  $\text{softplus}_\tau(0) - \text{ReLU}(0) = \log(2)\tau$ . Analyzing the asymptotic behavior,  
782  $\lim_{x \rightarrow \pm\infty} (\text{softplus}_\tau(x) - \text{ReLU}(x)) = 0$ , and therefore  $\text{softplus}_\tau(x) > \text{ReLU}(x)$  for  $x \in \mathbb{R}$ .  $\square$

## 783 D Additional Empirical Details and Results

### 784 D.1 Experimental details

785 All algorithms are implemented in BoTorch. The analytic EI, qEI, cEI utilize the standard BoTorch  
786 implementations. We utilize the original authors' implementations of single objective JES [25],  
787 GIBBON [45], and multi-objective JES [55], which are all available in the main BoTorch repository.  
788 All simulations are ran with 32 replicates and error bars represent  $\pm 2$  times the standard error of  
789 the mean. We use a Matern-5/2 kernel with automatic relevance determination (ARD), i.e. separate  
790 length-scales for each input dimension, and a top-hat prior on the length-scales in  $[0.01, 100]$ . The  
791 input spaces are normalized to the unit hyper-cube and the objective values are standardized during  
792 each optimization iteration.

### 793 D.2 Combining LogEI with TuRBO for High-Dimensional Bayesian Optimization

794 In the main text, we show how LogEI performs particularly well relative to other baselines in high  
795 dimensional spaces. Here, we show how LogEI can work synergistically with trust-region based  
796 methods for high-dimensional BO, such as TuRBO [14].

797 Fig. 13 compares the performance of LogEI, TuRBO-1 + LogEI, TuRBO-1 + EI, as well as the  
798 original Thompson-sampling based implementation for the 50d Ackley test problem. Combining  
799 TuRBO-1 with LogEI results in substantially better performance than the baselines. Since we  
800 optimize batches of  $q = 50$  candidates jointly, we also increase the number of Monte-Carlo samples  
801 from the Gaussian process from 128, the BoTorch default, to 512, and use the fat-tailed smooth  
802 approximations of Sec. A.3 to ensure a strong gradient signal to all candidates of the batch.

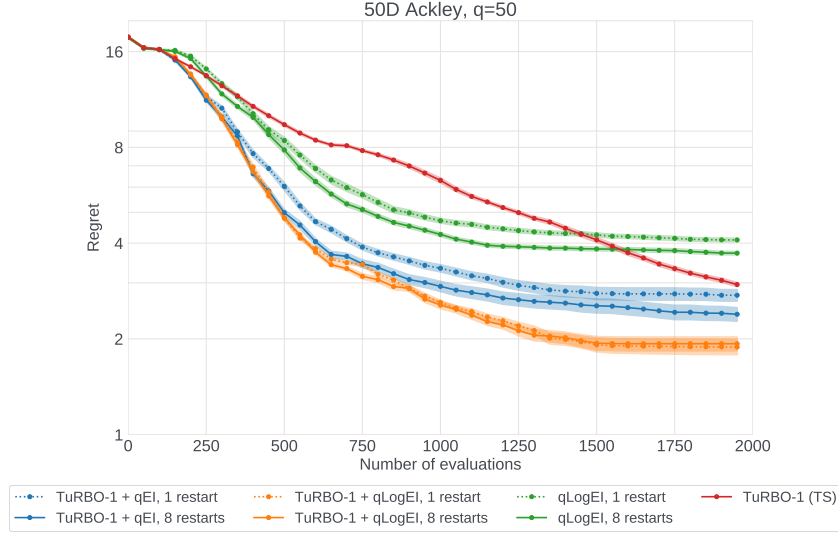


Figure 13: Combining LogEI with TuRBO on the high-dimensional on the 50d Ackley problem yields significant improvement in sample complexity. Unlike qEI, no random restarts are necessary to achieve good performance when performing joint optimization of the batch ( $q = 50$ ). Notably, LogEI without TuRBO is comparable to the performance of TuRBO with Thompson sampling.

### 803 D.3 Constrained Problems

804 While running the benchmarks using cEI in section 5, we found that we in fact improved upon a best  
805 known result from the literature. We compare with the results in Coello and Montes [7], which are  
806 generated using 30 runs of **80,000 function evaluations** each.

- 807 • For the pressure vessel design problem, Coello and Montes [7] quote a best-case feasible objective  
808 of 6059.946341. Out of just 16 different runs, LogEI achieves a worst-case feasible objective  
809 of 5659.1108 **after only 110 evaluations**, and a best case of 5651.8862, a notable reduction in  
810 objective value using almost three orders of magnitude fewer function evaluations.
- 811 • For the welded beam problem, Coello and Montes [7] quote 1.728226, whereas LogEI found a  
812 best case of 1.7496 after 110 evaluations, which is lightly worse, but we stress that this is using  
813 three orders of magnitude fewer evaluations.
- 814 • For the tension-compression problem, LogEI found a feasible solution with value 0.0129 after  
815 110 evaluations compared to the 0.012681 reported in in [7].

816 We emphasize that genetic algorithms and BO are generally concerned with distinct problem classes:  
817 BO focuses heavily on sample efficiency and the small-data regime, while genetic algorithms often  
818 utilize a substantially larger number of function evaluations. The results here show that in this case  
819 BO is competitive with and can even outperforms a genetic algorithm, using only a tiny fraction of  
820 the sample budget. Sample efficiency is particularly relevant for physical simulators whose evaluation  
821 takes significant computational effort, often rendering several tens of thousands of evaluations  
822 infeasible.

### 823 D.4 Parallel Bayesian Optimization with cross-batch constraints

824 In some parallel Bayesian optimization settings, batch optimization is subject to non-trivial constraints  
825 across the batch elements. A natural example for this are budget constraints. For instance, in the  
826 context of experimental material science, consider the case where each manufactured compound  
827 requires a certain amount of different materials (as described by its parameters), but there is only  
828 a fixed total amount of material available (e.g., because the stock is limited due to cost and/or  
829 storage capacity). In such a situation, batch generation will be subject to a budget constraint that  
830 is not separable across the elements of the batch. Importantly, in that case sequential greedy batch



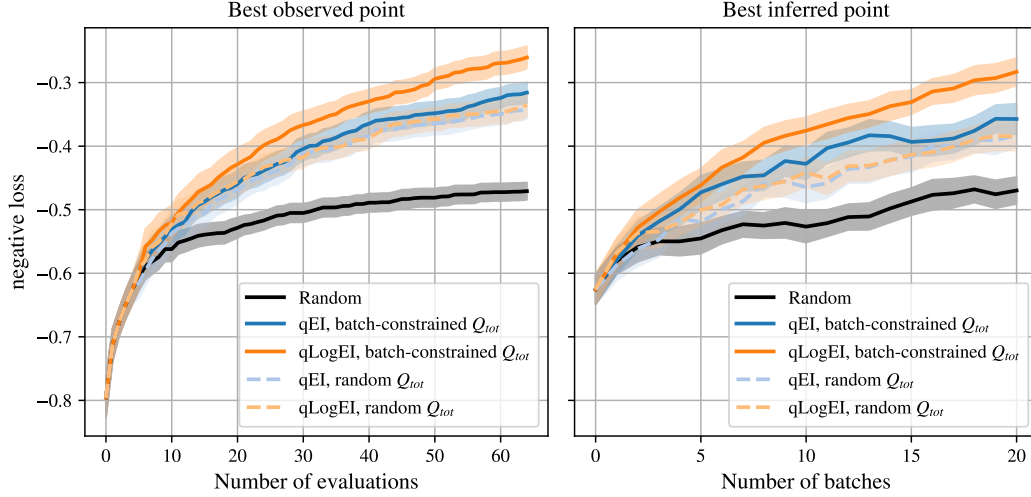


Figure 14: Optimization results on the nanomaterial synthesis material science problem with cross-batch constraints. While qLogEI outperforms qEI under the proper constrained (“batch-constrained  $Q_{tot}$ ”) optimization, this is not the case for the heuristic (“random  $Q_{tot}$ ”), demonstrating the value of both joint batch optimization with constraints and LogEI.

831 generation is not an option since it is not able to incorporate the budget constraint. Therefore, joint  
832 batch optimization is required.

833 Here we give one such example in the context of Bayesian Optimization for sequential experimental  
834 design. We consider the five-dimensional silver nanoparticle flow synthesis problem from Liang  
835 et al. [39]. In this problem, to goal is to optimize the absorbance spectrum score of the synthesized  
836 nanoparticles over five parameters: four flow rate ratios of different components (silver, silver nitrate,  
837 trisodium citrate, polyvinyl alcohol) and a total flow rate  $Q_{tot}$ .

838 The original problem was optimized over a discrete set of parameterizations. For our purposes we  
839 created a continuous surrogate model based on the experimental dataset (available from <https://github.com/PV-Lab/Benchmarking>) by fitting an RBF interpolator (smoothing factor of 0.01)  
840 in `scipy` on the (negative) loss. We use the same search space as Liang et al. [39], but in addition  
841 to the box bounds on the parameters we also impose an additional constraint on the total flow rate  
842  $Q_{tot}^{max} = 2000 \mu\text{L}/\text{min}$  across the batch:  $\sum_{i=1}^q Q_{tot}^i \leq Q_{tot}^{max}$  (the maximum flow rate per syringe  
843 / batch element is  $1000 \mu\text{L}/\text{min}$ ). This constraint expresses the maximum throughput limit of the  
844 microfluidic experimentation setup. The result of this constraint is that we cannot consider the batch  
845 elements (in this case automated syringe pumps) have all elements of a batch of experiments operate  
846 in the high-flow regime at the same time.

848 In our experiment, we use a batch size of  $q = 3$  and start the optimization from 5 randomly sampled  
849 points from the domain. We run 75 replicates with random initial conditions (shared across the  
850 different methods), error bars show  $\pm$  two times the standard error of the mean. Our baseline is  
851 uniform random sampling from the domain (we use a hit-and-run sampler to sample uniformly from  
852 the constraint polytope  $\sum_{i=1}^q Q_{tot}^i \leq Q_{tot}^{max}$ ). We compare qEI vs. qLogEI, and for each of the  
853 two we evaluate (i) the version with the batch constraint imposed explicitly in the optimizer (the  
854 optimization in this case uses `scipy`’s SLSQP solver), and (ii) a heuristic that first samples the total  
855 flow rates  $\{Q_{tot}^i\}_{i=1}^q$  uniformly from the constraint set, and then optimizes the acquisition function  
856 with the flow rates fixed to the sampled values.

857 The results in Figure 14 show that while both the heuristic (“random  $Q_{tot}$ ”) and the proper constrained  
858 optimization (“batch-constrained  $Q_{tot}$ ”) substantially outperform the purely random baseline, it  
859 requires using both LogEI and proper constraints to achieve additional performance gains over the  
860 other 3 combinations. Importantly, this approach is only possible by performing joint optimization of  
861 the batch, which underlines the importance of qLogEI and its siblings being able to achieve superior  
862 joint batch optimization in settings like this.

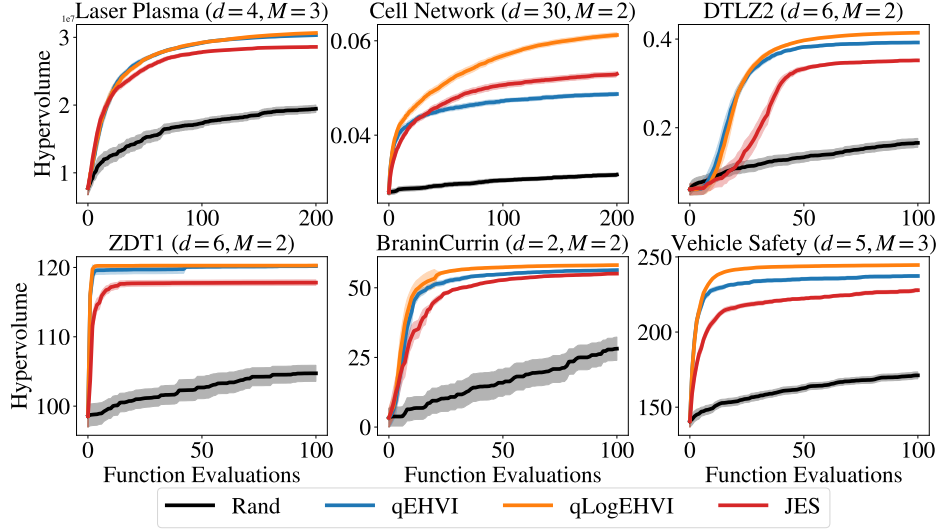


Figure 15: Sequential ( $q = 1$ ) optimization performance on multi-objective problems, as measured by the hypervolume of the Pareto frontier across observed points. This plot includes JES [55]. Similar to the single-objective case, qLogEHVI significantly outperforms all baselines on all test problems.

## 863 D.5 Details on Multi-Objective Problems

864 We consider a variety of multi-objective benchmark problems. We evaluate performance on three  
 865 synthetic biobjective problems Branin-Currin ( $d = 2$ ) [4], ZDT1 ( $d = 6$ ) [66], and DTLZ2 ( $d = 6$ )  
 866 [10]. As described in 5, we also evaluated performance on three real world inspired problems. For  
 867 the laser plasma acceleration problem, we used the public data available at Irshad et al. [27] to fit  
 868 an independent GP surrogate model to each objective. We only queried the surrogate at the highest  
 869 fidelity to create a single fidelity benchmark.

## 870 D.6 Effect of Temperature Parameter

871 In Figure 16, we examine the effect of fixed  $\tau$  for the softplus operator on optimization performance.  
 872 We find that smaller values typically work better.

## 873 D.7 Effect of the initialization strategy

874 Packages and frameworks commonly utilize smart initialization heuristics to improve acquisition  
 875 function optimization performance. In Figure 17, we compare simple random restart optimization,  
 876 where initial points are selected uniformly at random, with BoTorch’s default initialization strategy,  
 877 which evaluates the acquisition function on a large number of points selected from a scrambled Sobol  
 878 sequence, and selects  $n$  points at random via Boltzman sampling (e.g., sampling using probabilities  
 879 computed by taking a softmax over the acquisition values [2]. Here we consider 1024 initial  
 880 candidates. We find that the BoTorch initialization strategy improves regret for all cases, and that  
 881 qLogEI, followed by UCB show less sensitivity to the choice of initializations strategy. Figure 18  
 882 examines the sensitivity of qEI to the number of initial starting points when performing standard  
 883 random restart optimization and jointly optimizing the  $q$  points in the batch. We find that, consistent  
 884 with our empirical and theoretical results in the main text, qEI often gets stuck in local minima for  
 885 the Ackley test function, and additional random restarts often improve results but do not compensate  
 886 for the fundamental optimality gap. The performance of qLogEI also improves as the number of  
 887 starting points increases.

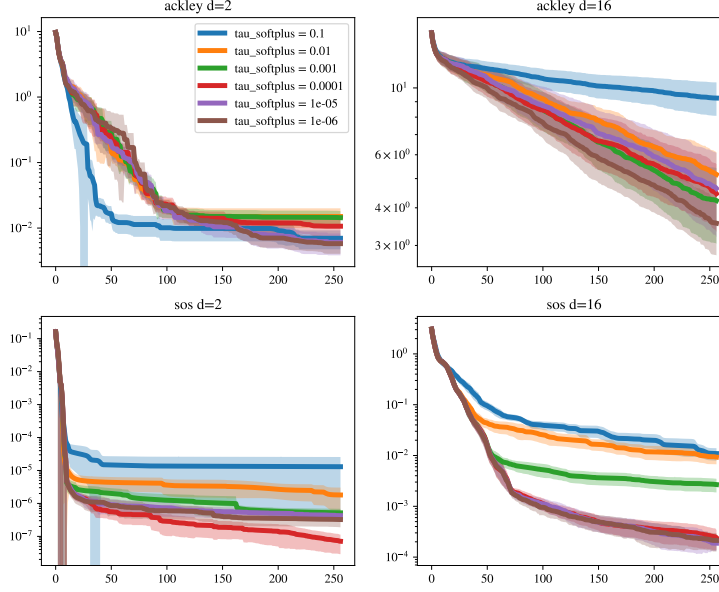


Figure 16: Ablation study on the convergence characteristics of LogEI on Ackley and sum of squares (SOS) problems in 2 and 16 dimensions. The study shows that it is important to choose a small  $\tau_0$  for the best convergence properties, which results in a very tight approximation to the original ReLU non-linearity in the integrand. Critically, setting  $\tau_0$  as low as  $10^{-6}$  is only possible due to the transformation of all computations into log-space. Otherwise, the smoothed acquisition utility would exhibit similarly numerically vanishing gradients as the original ReLU non-linearity.

	CELL NETWORK	BRANIN-CURRIN	DTLZ2	LASER PLASMA	ZDT1	VEHICLE SAFETY
JES	21.6 (+/- 1.1)	89.6 (+/- 3.3)	33.6 (+/- 1.0)	57.3 (+/- 0.7)	72.7 (+/- 1.0)	47.0 (+/- 1.6)
QEHVI	0.6 (+/- 0.0)	0.7 (+/- 0.0)	1.0 (+/- 0.0)	3.0 (+/- 0.1)	0.6 (+/- 0.0)	0.6 (+/- 0.0)
QLOGEHVI	9.2 (+/- 0.8)	10.0 (+/- 0.4)	5.8 (+/- 0.2)	31.6 (+/- 1.7)	7.2 (+/- 0.7)	2.1 (+/- 0.1)
RAND	0.2 (+/- 0.0)	0.2 (+/- 0.0)	0.2 (+/- 0.0)	0.3 (+/- 0.0)	0.3 (+/- 0.0)	0.3 (+/- 0.0)

Table 1: Acquisition function optimization wall time in seconds on CPU (2x Intel Xeon E5-2680 v4 @ 2.40GHz) . We report the mean and  $\pm 2$  standard errors.

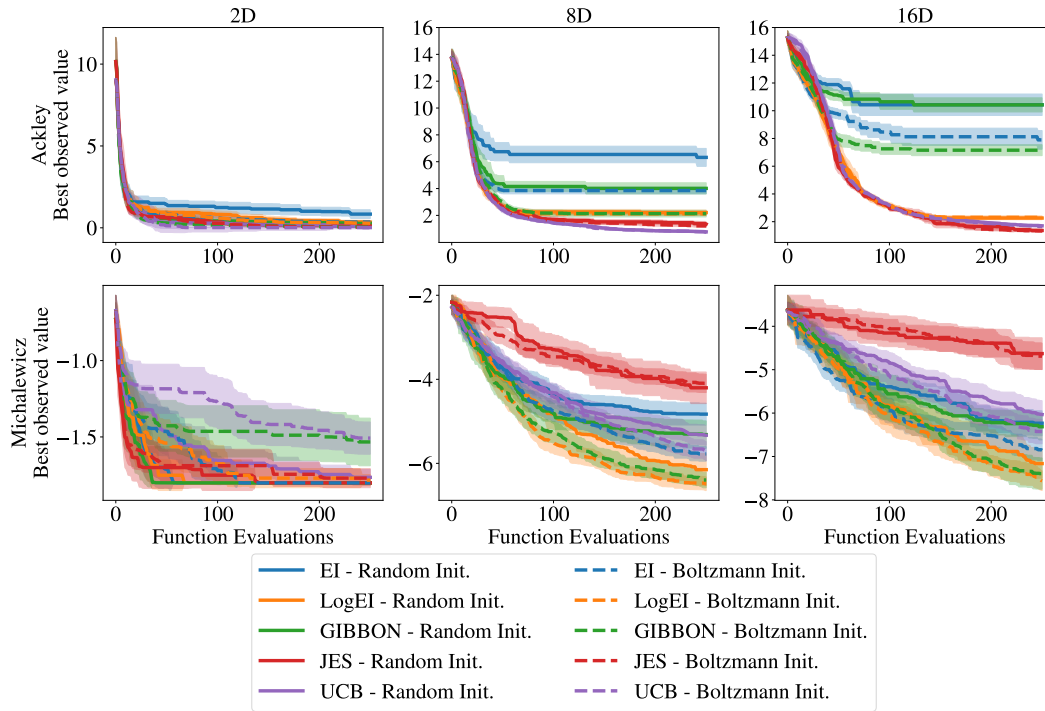


Figure 17: Sensitivity to the initialization strategy. Random selects random restart points from the design space uniformly at random, whereas Boltzmann initialization is the default BoTorch initialization strategy which selects points with higher acquisition function values with a higher probability via Boltzmann sampling.

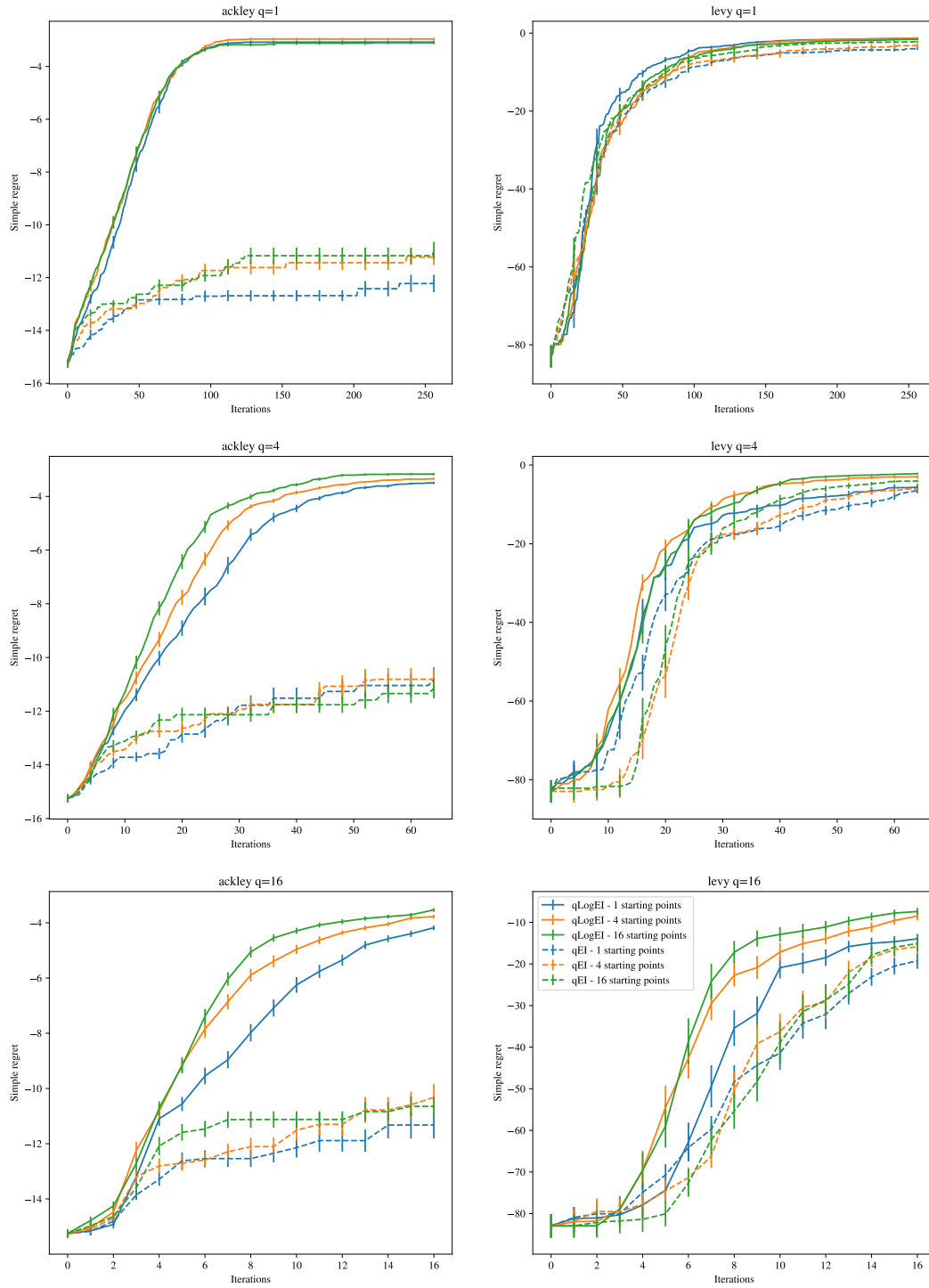


Figure 18: Sensitivity to number of starting points with multi-start optimization for the 16D Ackley and Levy test problems. Note: We plot negative regret, so higher is better.