
ORTHOBO: Orthogonal Bayesian Hyperparameter Optimization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Bayesian optimization is widely used for hyperparameter optimization when model
2 evaluations are expensive, but noisy acquisition estimates can lead to unstable
3 decisions. We identify acquisition estimation noise as a distinct failure mode: even
4 when the surrogate model and acquisition target are appropriate, finite-sample
5 Monte Carlo error can perturb acquisition values, flip candidate rankings, and lead
6 Bayesian optimization to evaluate suboptimal configurations. As a remedy, we
7 propose an *orthogonal acquisition estimator*, that subtracts an optimally weighted
8 score-function control variate, yielding an acquisition residual orthogonal to poste-
9 rior score directions and thereby reducing Monte Carlo variance. Building on this
10 estimator, we introduce ORTHOBO, a Bayesian optimization framework that com-
11 bines orthogonalized acquisition estimates with ensemble surrogates for structural
12 misspecification and an outer log transformation for numerical stability. Theoret-
13 ically, we prove target preservation, variance reduction, and improved pairwise
14 ranking stability. Empirically, ORTHOBO substantially reduces acquisition estima-
15 tion variance, stabilizes candidate rankings, and achieves strong performance across
16 synthetic benchmarks and downstream use cases, including vision-transformer fine-
17 tuning on an industrial wafer-map classification task. These results show that
18 stabilizing acquisition estimation can directly improve the reliability and sample
19 efficiency of Bayesian hyperparameter optimization.

20 1 Introduction

21 Hyperparameter optimization (HPO) is a fundamental component of modern machine learning
22 systems, and has been used across a wide range of domains, including biology [5, 47], chemistry [65],
23 manufacturing [1], medicine [42], and physics [55]. Because evaluating hyperparameter candidates
24 can require substantial computation, *Bayesian optimization (BO)* is frequently used to identify strong
25 hyperparameter configurations under limited evaluation budgets [e.g., 19, 51, 53, 57, 63]. BO fits
26 a probabilistic surrogate model of the objective and selects new configurations by maximizing an
27 acquisition function such as expected improvement (EI) [34].¹

28 In our work, we focus on a source of instability that is central to BO decisions but has received
29 comparatively little attention: *the estimation of the acquisition value itself*. Acquisition functions are
30 often treated as deterministic once a surrogate model has been fitted. In many practical BO pipelines,
31 however, the acquisition value is obtained by marginalizing over uncertain surrogate parameters,
32 aggregating across surrogate models, or using Monte Carlo (MC) approximations [e.g., 10, 30, 46, 53].
33 As a result, the quantity used to rank candidate configurations is itself a noisy estimate. This matters
34 because BO acts on *rankings*: even small estimation errors can flip the order of two candidates
35 with similar acquisition values and cause the algorithm to spend an expensive evaluation on a worse
36 configuration.

¹For a thorough introduction to BO see [e.g., 23, 51].

37 At a technical level, we treat the estimation of the acquisition function as a MC estimation problem
 38 [e.g., 10, 29, 30, 53] and propose a novel *orthogonalized estimator of the acquisition value* that
 39 serves as a variance-reduction technique, thus improving robustness to sampling errors and stabi-
 40 lizing acquisition rankings. Our notion of orthogonality stems from the orthogonal ML literature
 41 [e.g., 12, 20, 35, 37] in that our estimator offers robustness in terms of variance reduction guaran-
 42 tees. Note that our notion of robustness differs from that in existing ‘robust’ regression methods
 43 which aim to protect the surrogate model against corrupted or heavy-tailed observations [e.g., 38].

44 **Example.** Consider tuning a clas-
 45 sification model for quality control
 46 on a large-scale industrial manufactur-
 47 ing dataset, where each evaluation
 48 requires training on many images and
 49 assessing performance under a highly
 50 skewed class distribution. Two hyper-
 51 parameter candidate configurations
 52 may differ only slightly, e.g., a deeper
 53 network with stronger dropout versus
 54 a shallower network with weaker
 55 regularization. If acquisition values
 56 are estimated noisily, BO may spend
 57 an expensive evaluation on the worse
 58 configuration simply because it was
 59 ranked first due to sampling noise.
 60

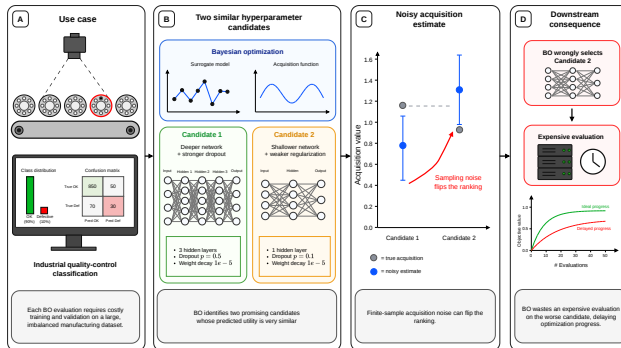


Figure 1: Motivational example of instabilities in Bayesian HPO.

61 Our contribution is complementary to existing work on robust BO. Prior methods have mainly
 62 focused on improving the surrogate model, for example by accounting for kernel misspecification,
 63 hyperparameter uncertainty, unreliable posterior uncertainty, trust-region localization, or ensemble
 64 surrogates [11, 8, 43, 16, 36, 46, 57]. Other work has improved the acquisition formula or its
 65 numerical optimization [2]. However, these approaches do **not** directly address the *variance of the*
 66 *estimated acquisition values* that drive BO decisions. In contrast, our method improves the statistical
 67 quality of the estimated acquisition values used for candidate selection through a variance-reduction
 68 techniques based on orthogonal scores.

69 Building on our orthogonal acquisition estimator, we introduce ORTHOBO, a BO framework for
 70 robust hyperparameter optimization. ORTHOBO combines three components: orthogonalized acqui-
 71 sition estimation to reduce MC acquisition noise, ensemble surrogates to mitigate structural surrogate
 72 misspecification, and an outer log transformation to improve numerical stability during acquisition
 73 maximization. The central idea is to separate these roles: surrogate ensembles improve the model
 74 class, orthogonalization improves acquisition estimation, and the outer log transformation improves
 75 numerical conditioning.

76 Our ORTHOBO offers several theoretical guarantees. The proposed orthogonalized estimator is
 77 unbiased for the marginal EI target and achieves lower variance than the corresponding naïve MC
 78 estimator (\rightarrow Theorem 4.1). Additionally, the lower acquisition-estimation variance improves the
 79 stability of pairwise acquisition rankings (\rightarrow Proposition 4.2), reducing the probability of falsely
 80 selecting a candidate based on MC noise. Empirically, we show that ORTHOBO improves robustness
 81 and sample efficiency in HPO tasks (Section 5).

82 Our **main contributions** are: (1) We propose a novel, orthogonalized estimator for hyperparameter
 83 acquisition, reducing estimation variance without changing the underlying acquisition target. (2) We
 84 develop ORTHOBO, a BO framework that explicitly targets acquisition estimation noise, structural
 85 surrogate misspecification, and numerical instabilities. (3) We provide theoretical guarantees for
 86 variance reduction and ranking stability, and demonstrate improved empirical performance on various
 87 HPO benchmarks. Importantly, ORTHOBO is also applicable to a wide range of other BO problems.

88 2 Related work²

89 **BO for hyperparameter tuning.** BO is a standard approach for HPO [e.g., 19, 53, 63], particularly
90 when evaluations of the objective function are expensive. In practice, BO relies on a surrogate model
91 to approximate the objective and guide the search. A common choice for the surrogate model are
92 Gaussian processes (GPs) [27, 34, 53], while practical alternatives include tree-based and density-
93 based surrogates such as tree-structured Parzen estimators (TPE) [6, 7, 60]. Subsequent work has
94 addressed scaling challenges through multi-fidelity methods, batch BO, and trust-region approaches
95 [e.g., 16, 17, 26, 61]. However, no method has addressed the misspecification and performance
96 impact due to acquisition estimation noise.

97 **Acquisition functions.** The surrogate models are used by acquisition functions, which are heuristics
98 to evaluate the utility of hyperparameter candidates before querying the (expensive) objective function
99 [4]. The performance of the acquisition function thus directly influences the overall performance. In
100 the literature, several types of acquisition functions or acquisition frameworks have been proposed,
101 including qLogEI [2], GIBBON [40], UCB [62], MES [59], RPR [3], and TuRBO [16]. Our work is
102 orthogonal to these works, as we focus on improving the *acquisition estimation* and not the acquisition
103 function itself.³

104 **Robustness to surrogate misspecification.** A central weakness of BO is its reliance on a surro-
105 gate model which might be misspecified. Prior work has studied BO under misspecified kernel
106 classes [11], unknown hyperparameters [8, 27], unreliable posterior uncertainty [43], and a surrogate
107 complexity-efficiency trade-off [10]. Practical strategies for mitigating such failures include trust-
108 region localization [16], robust likelihood models [3, 38], and model averaging or ensemble methods
109 [36, 46]. However, these approaches address *structural* surrogate misspecification. In contrast, our
110 method additionally targets misspecification due to *acquisition estimation* noise.

111 **Research gap.** We identify acquisition estimation errors under surrogate uncertainty as a crucial, yet
112 previously overlooked, failure mode in Bayesian optimization. We propose an orthogonalized estima-
113 tor that reduces variance and stabilizes acquisition-based decisions. To the best of our knowledge,
114 this is the first work to explicitly formulate acquisition computation in BO as a variance-reduction
115 problem and to develop an orthogonalized estimator for improving acquisition-based decisions.

116 3 Problem setting

117 **Setup.** We consider the standard setup for HPO over a search space $\Lambda \subseteq \mathbb{R}^d$ [3, 15]. For a
118 hyperparameter configuration $\lambda \in \Lambda$, let $f(\lambda)$ denote the unknown performance of interest, e.g., the
119 validation loss or validation error after training a model with hyperparameters λ . The objective f is
120 expensive to evaluate and can only be observed through noisy evaluations $y = f(\lambda) + \epsilon$, where ϵ
121 captures evaluation noise. At each iteration $t = 1, \dots, T$, the available data are $\mathcal{D}_t = \{(\lambda_i, y_i)\}_{i=1}^t$.
122 We aim to sequentially choose hyperparameter configurations $\lambda_1, \lambda_2, \dots, \lambda_T$ to identify a high-
123 performing configuration within an evaluation budget T .

124 **Recap: BO.** The central idea of BO is to maintain a probabilistic model of the unknown objective and
125 to use it to guide future evaluations [24]. At iteration $t = 1, \dots, T$, BO fits a probabilistic *surrogate*
126 *model* m to the current dataset and uses an *acquisition function* α to select the next evaluation point
127 [3]. The surrogate provides two crucial quantities for decision-making: (i) a prediction of objective
128 values in unexplored regions, and (ii) a measure of uncertainty about those predictions. BO then
129 combines these two ingredients through α , which scores candidate points according to their potential
130 utility for optimization. After selecting a new candidate point, the objective f is evaluated, and the
131 new observation is added to the dataset \mathcal{D} . The surrogate is then updated, and the process repeats. In
132 this way, BO adaptively concentrates evaluations in regions that are both informative and promising.
133 As a result, the performance of BO depends on two components: the quality of the surrogate model
134 and the stability of the acquisition values used to rank candidate configurations.

135 **Surrogate model.** Let m index a surrogate model, and θ_m denote its latent parameters. Conditional
136 on θ_m , the surrogate induces a posterior predictive distribution for the objective $f(\lambda)$ given \mathcal{D}_t . In
137 practice, the parameters θ_m are not known and must themselves be inferred from data.

²We provide an extended related work in Supplement A. Therein, we also provide a broader review of BO under different types of misspecification and orthogonal ML.

³We provide a generalization of ORTHOBO at the acquisition function level in Supplement B.

138 In our work, we focus on the *expected improvement (EI)* as our acquisition function. For a surrogate
 139 model m with parameters θ_m , we define

$$\text{EI}_m(\lambda_t; \theta_m) = \mathbb{E}[(f^* - f(\lambda_t))_+ | \mathcal{D}_t, \theta_m], \quad (1)$$

140 where $f^* = \min_{i \leq t} y_i$ is the currently best observed value up to iteration t and $(x)_+ := \max\{0, x\}$
 141 [34, 53]. In practice, the model parameters are not fixed but follow an approximate posterior
 142 $\theta_m \sim q_{m,t}(\theta)$. The marginal EI over the parameters is then given by

$$\text{EI}_m^{\text{marg}}(\lambda_t) = \mathbb{E}_{\theta_m \sim q_{m,t}}[\text{EI}_m(\lambda_t; \theta_m)], \quad (2)$$

143 which is commonly approximated through MC samples $s = 1, \dots, S$ given a specific MC sampling
 144 budget S [4, 53, 61], i.e.,

$$\widehat{\text{EI}}_m^{\text{MC}}(\lambda_t) = \frac{1}{S} \sum_{s=1}^S \text{EI}_m(\lambda_t; \theta_m^{(s)}), \quad \theta_m^{(s)} \sim q_{m,t}. \quad (3)$$

145 However, $\widehat{\text{EI}}_m^{\text{MC}}$ suffers from high variance for small budgets S and sensitivity to approximation errors
 146 in the posterior $q_{m,t}$, leading to instability in the estimated acquisition and thus the BO decisions.

147 **Failure modes in BO.** Several failure modes arise in a realistic HPO setting[e.g. 2, 8, 11, 43].
 148 Bayesian HPO performance depends on the accuracy and stability of acquisition estimates, which are
 149 computed from a surrogate posterior that is often *misspecified and approximately inferred*, leading to
 150 instability in acquisition values. Specifically, BO suffers from the following failure modes:

151 (i) *Structural surrogate misspecification* arises when the surrogate model class does not adequately
 152 represent the unknown objective. For example, the assumed kernel, likelihood, prior smoothness, or
 153 stationarity structure may be incompatible with the true response surface. As a result, the posterior
 154 mean and posterior uncertainty estimates may be systematically biased.

155 (ii) *Posterior approximation and acquisition-estimation noise.* Marginalizing over surrogate uncer-
 156 tainty mitigates plug-in overconfidence, but the marginal acquisition value must be approximated
 157 from a finite number of posterior samples. This introduces *Monte Carlo variance* into the acquisition
 158 estimate, on top of any bias from the approximate posterior $q_{m,t}$ itself differing from the true posterior.
 159 Both effects can destabilize the ranking of candidate points.

160 (iii) *Numerical instability* can arise when acquisition values are very small, highly skewed, or nearly
 161 indistinguishable across candidates. Directly optimizing such quantities can lead to unstable gradients.

162 Prior work has recognized several of these failure modes. However, the finite-sample noise induced
 163 by estimating marginal acquisition values has not been addressed so far.

164 **Our work.** In this work, we address instability in BO at the level of the acquisition estimate itself.
 165 Even with a reasonable surrogate family, BO may still perform poorly if the acquisition estimator used
 166 to rank candidate points is noisy, sensitive to approximate inference, or numerically ill-conditioned.

167 We treat the EI-scale marginal acquisition obtained by integrating over surrogate uncertainty as an
 168 MC estimation problem. We propose an orthogonalized estimator that reduces variance induced
 169 by surrogate uncertainty, while preserving the underlying EI-scale target. To address structural
 170 surrogate misspecification and overcome numerical instabilities, we combine our estimator with
 171 an ensemble of surrogate models and adaptive model weighting optimized through an outer log
 172 transformation. Overall, ORTHOBO reduces acquisition-estimation variance, stabilizes candidate
 173 rankings, and improves BO decision quality.

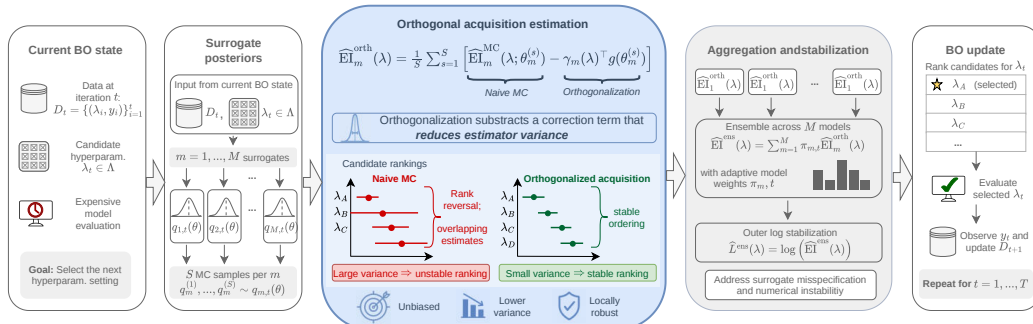


Figure 2: Our proposed ORTHOBO. It improves acquisition estimation through variance stabilization. For this, we orthogonalize the acquisition estimation without changing the acquisition target.

174 **4 ORTHOBO**

175 Below, we present ORTHOBO, our orthogonalized framework for robust HPO addressing the failure
 176 modes discussed above (an overview is shown in Fig. 2). We first present our *main contribution*, the
 177 general *orthogonalized acquisition estimation framework*. Of note, we do not change the acquisition
 178 target, but rather improve its estimation through variance stabilization. Then, we provide example
 179 instantiations of our orthogonal framework for two common surrogate classes. In the main text, we
 180 focus on EI-based acquisition functions in order to keep the notation and theoretical arguments simple.
 181 The key object throughout is the marginal EI under surrogate uncertainty and ensemble aggregation.
 182 We provide extensions to constrained and parallel acquisition functions in Supplement B.

183 **4.1 Orthogonalized acquisition estimation**

184 We now address the instability due to acquisition noise. Our key idea is to introduce an *orthogonalized*
 185 *acquisition functional* that shares the same marginal target as the original EI but has reduced variance
 186 under the posterior $q_{m,t}$. Let $g(\theta) := \nabla_{\theta} \log q_{m,t}(\theta)$ denote the posterior score function.

We define the *orthogonalized expected improvement* as

$$\text{EI}_m^{\text{orth}}(\lambda; \theta) := \text{EI}_m(\lambda; \theta) - \gamma_m(\lambda)^\top g(\theta), \quad (4)$$

where $\gamma_m(\lambda) := \Sigma_g^{-1} \text{Cov}_{q_{m,t}}(g, \text{EI}_m(\lambda; \cdot))$ and $\Sigma_g := \text{Cov}_{q_{m,t}}(g, g)$. The corresponding
 marginal acquisition is

$$\text{EI}_m^{\text{orth,marg}}(\lambda) := \mathbb{E}_{\theta \sim q_{m,t}} [\text{EI}_m^{\text{orth}}(\lambda; \theta)]. \quad (5)$$

The associated MC estimator is

$$\widehat{\text{EI}}_m^{\text{orth}}(\lambda) = \frac{1}{S} \sum_{s=1}^S \text{EI}_m^{\text{orth}}(\lambda; \theta_m^{(s)}), \quad \theta_m^{(s)} \sim q_{m,t}. \quad (6)$$

187
 188 Of note, the construction is non-trivial because $g(\theta)$ is based on the *parameter posterior*, whereas
 189 common MC pipelines sample only from the predictive posterior over function values. We therefore
 190 draw parameter samples explicitly and evaluate both EI_m and g at the same $\theta_m^{(s)}$. In practice, $\gamma_m(\lambda)$
 191 is unknown and we use the empirical plug-in $\hat{\gamma}_m(\lambda)$ estimated from the same MC samples; this
 192 introduces an in-sample bias of order $O(1/S)$ that is dominated by the $O(1/\sqrt{S})$ sampling standard
 193 deviation at the MC budgets used in our experiments.⁴

194 **Properties of the orthogonalized acquisition.** The orthogonalized EI satisfies three key properties:
 195 (i) *target preservation*: it shares the same marginal as the original EI, so the BO target is unchanged;
 196 (ii) *variance reduction*: it has smaller variance under $q_{m,t}$; and (iii) *local robustness*: it is first-order
 197 insensitive to score-tilt perturbations of $q_{m,t}$.

198 **Regularity for the score function.** Throughout our analysis, we assume that the surrogate-
 199 parameter distribution $q_{m,t}$ has a differentiable density on a support $\Theta \subseteq \mathbb{R}^{d_\theta}$ and that the cor-
 200 responding boundary term vanishes, i.e. $\int_{\Theta} \nabla_{\theta} q_{m,t}(\theta) d\theta = 0$. This holds, for example, if $\Theta = \mathbb{R}^{d_\theta}$
 201 and $q_{m,t}(\theta) \rightarrow 0$ sufficiently fast as $\|\theta\| \rightarrow \infty$.

202 **Theorem 4.1.** Assume $\mathbb{E}_{q_{m,t}}[\text{EI}_m(\lambda; \theta)^2] < \infty$ and $\mathbb{E}_{q_{m,t}}[\|g(\theta)\|_2^2] < \infty$, that Σ_g is nonsingular,
 203 and that the standard score-function regularity conditions hold for $q_{m,t}$. Then:

204 (i) **Target preservation.** $\mathbb{E}[\text{EI}_m^{\text{orth}}(\lambda; \theta)] = \text{EI}_m^{\text{marg}}(\lambda; \theta)$.

205 (ii) **Variance reduction.** Under $q_{m,t}$,

$$\text{Var}(\text{EI}_m^{\text{orth}}(\lambda; \theta)) = \text{Var}(\text{EI}_m(\lambda; \theta)) - \text{Cov}(g, \text{EI}_m)^\top \Sigma_g^{-1} \text{Cov}(g, \text{EI}_m) \leq \text{Var}(\text{EI}_m(\lambda; \theta)). \quad (7)$$

206 (iii) **Local robustness in score-tilt directions.** For any $b \in \mathbb{R}^{d_\theta}$, consider the tilted family $q_{m,t}^{(\varepsilon)}(\theta) \propto$
 207 $q_{m,t}(\theta) \exp(\varepsilon b^\top g(\theta))$ for small $|\varepsilon|$. Holding $\gamma_m(\lambda)$ fixed at its $\varepsilon = 0$ value,

$$\left. \frac{d}{d\varepsilon} \mathbb{E}_{q_{m,t}^{(\varepsilon)}} [\text{EI}_m^{\text{orth}}(\lambda; \theta)] \right|_{\varepsilon=0} = 0. \quad (8)$$

⁴Sample-splitting (cross-fitting) would restore exact finite-sample unbiasedness without changing the leading asymptotic variance; see Supplement D.

208 *Proof.* We provide a proof for Theorem 4.1 in Supplement D. \square

209 **From acquisition to estimator.** Theorem 4.1(i) shows that orthogonalization preserves the BO
 210 target. Since the MC estimator is an i.i.d. sample mean, $\text{Var}(\widehat{\text{EI}}_m^{\text{orth}}(\lambda)) = \text{Var}(\text{EI}_m^{\text{orth}}(\lambda; \theta))/S$,
 211 Theorem 4.1(ii) directly implies $\text{Var}(\widehat{\text{EI}}_m^{\text{orth}}(\lambda)) \leq \text{Var}(\widehat{\text{EI}}_m^{\text{MC}}(\lambda))$ for any MC budget S , which
 212 improves the stability of acquisition-based decisions. Theorem 4.1(iii) provides local robustness in
 213 the orthogonal ML sense [12, 35]: small score-tilt perturbations of $q_{m,t}$ do not affect the orthogonal-
 214 ized acquisition to first order. Importantly, Theorem 4.1 is a *within-model* result: orthogonalization
 215 stabilizes acquisition estimation conditional on $q_{m,t}$, but does not directly correct structural misspeci-
 216 fication of the surrogate family. We address the latter through ensemble modeling in Section 4.3.

217 **Proposition 4.2 (Pairwise ranking stability).** *Let $\lambda, \lambda' \in \Lambda$ such that $\Delta(\lambda, \lambda') := \text{EI}_m^{\text{marg}}(\lambda) -$
 218 $\text{EI}_m^{\text{marg}}(\lambda') > 0$. Let $\widehat{\Delta}_{\text{MC}}(\lambda, \lambda')$ and $\widehat{\Delta}_{\text{orth}}(\lambda, \lambda')$ denote the corresponding Monte Carlo difference
 219 estimators. Then both estimators are unbiased for $\Delta(\lambda, \lambda')$, and*

$$\mathbb{P}(\widehat{\Delta}_{\text{orth}}(\lambda, \lambda') \leq 0) \leq \frac{\text{Var}(\widehat{\Delta}_{\text{orth}}(\lambda, \lambda'))}{\text{Var}(\widehat{\Delta}_{\text{orth}}(\lambda, \lambda')) + \Delta(\lambda, \lambda')^2} \leq \frac{\text{Var}(\widehat{\Delta}_{\text{MC}}(\lambda, \lambda'))}{\text{Var}(\widehat{\Delta}_{\text{MC}}(\lambda, \lambda')) + \Delta(\lambda, \lambda')^2}. \quad (9)$$

220 *Proof.* We provide a proof for Proposition 4.2 in Supplement D. \square

221 Proposition 4.2 shows how acquisition-estimation noise propagates into BO decisions. High estimator
 222 variance can flip the estimated sign of the acquisition gap between two candidate points, causing BO
 223 to select a suboptimal point. Orthogonalization reduces the variance without changing the acquisition
 224 target, and therefore directly reduces the probability of such ranking errors.

225 4.2 Instantiations for common surrogates

226 We instantiate ORTHOBO for two surrogate classes common in HPO: **A** Gaussian process and **B**
 227 tree-based density surrogates. The first corresponds to the standard BO setting, while the second
 228 captures widely used non-Gaussian and nonparametric approaches such as TPE-style optimization.

229 **A** Gaussian process surrogate

230 For GP surrogate models, acquisition noise is often driven by sensitivity of the mean $\mu(\cdot)$ and variance
 231 $\sigma(\cdot)$ function to lengthscales and noise hyperparameters. We show that orthogonalization reduces MC
 232 fluctuations induced by uncertainty in these directions. This is especially useful early in optimization
 233 when the approximation of the posterior over the GP hyperparameters is diffuse. Let

$$f \sim \mathcal{GP}(m_\theta(\cdot), k_\theta(\cdot, \cdot)), \quad (10)$$

234 where $\theta \in \Theta \subset \mathbb{R}^{d_\theta}$ denotes surrogate hyperparameters, such as kernel lengthscales, amplitude, and
 235 noise variance. Given data $\mathcal{D}_t = \{(\lambda_i, y_i)\}_{i=1}^t$, the predictive distribution at a candidate λ is

$$f(\lambda) | \mathcal{D}_t, \theta \sim \mathcal{N}(\mu_t(\lambda; \theta), \sigma_t^2(\lambda; \theta)), \quad (11)$$

236 with

$$\mu_t(\lambda; \theta) = m_\theta(\lambda) + k_\theta(\lambda, X_t)^\top (K_\theta(X_t, X_t) + \sigma_n^2 I_t)^{-1} (y_t - m_\theta(X_t)), \quad (12)$$

$$\sigma_t^2(\lambda; \theta) = k_\theta(\lambda, \lambda) - k_\theta(\lambda, X_t)^\top (K_\theta(X_t, X_t) + \sigma_n^2 I_t)^{-1} k_\theta(\lambda, X_t), \quad (13)$$

237 where $X_t = (\lambda_1, \dots, \lambda_t)$ and $y_t = (y_1, \dots, y_t)^\top$. We aim to maximize the Expected Improvement

$$\text{EI}^{\text{GP}}(\lambda; \theta) = (f^* - \mu_t(\lambda; \theta)) \Phi(z_t(\lambda; \theta)) + \sigma_t(\lambda; \theta) \phi(z_t(\lambda; \theta)), \quad (14)$$

238 where

$$z_t(\lambda; \theta) = \frac{f^* - \mu_t(\lambda; \theta)}{\sigma_t(\lambda; \theta)}, \quad f^* = \min_{i \leq t} y_i, \quad (15)$$

239 and Φ and ϕ denote the standard normal cumulative distribution and probability density function.

240 We assume an approximate posterior $q_t(\theta)$ over GP hyperparameters, obtained for example via a
 241 Laplace approximation or variational Gaussian approximation. The marginal acquisition is then

$$\text{EI}^{\text{GP,marg}}(\lambda) = \mathbb{E}_{\theta \sim q_t} [\text{EI}^{\text{GP}}(\lambda; \theta)], \quad (16)$$

242 and the orthogonalized estimator is

$$\widehat{\text{EI}}^{\text{GP,orth}}(\lambda) = \frac{1}{S} \sum_{s=1}^S [\text{EI}^{\text{GP}}(\lambda; \theta^{(s)}) - \gamma^{\text{GP}}(\lambda)^\top g_t(\theta^{(s)})], \quad \theta^{(s)} \sim q_t, \quad (17)$$

243 where $g_t(\theta) = \nabla_{\theta} \log q_t(\theta)$ and $\gamma^{\text{GP}}(\lambda) = \text{Cov}(g_t, g_t)^{-1} \text{Cov}(g_t, \text{EI}^{\text{GP}}(\lambda; \theta))$.

244 **B** Tree-based / TPE-style surrogate

245 As a second instantiation, we consider tree-based or density-based surrogates as in Tree-structured
246 Parzen Estimators (TPE). Rather than modeling $p(y | \lambda)$ directly, TPE-style methods model

$$p(\lambda | y) = \begin{cases} \ell_{\theta}(\lambda), & y < y^*, \\ g_{\theta}(\lambda), & y \geq y^*, \end{cases} \quad (18)$$

247 where y^* is a quantile threshold and θ denotes model parameters, such as kernel density bandwidths,
248 mixture weights, or tree parameters. In this setting, the acquisition is proportional to a density ratio:

$$\text{EI}^{\text{TPE}}(\lambda; \theta) \propto \frac{\ell_{\theta}(\lambda)}{g_{\theta}(\lambda)}. \quad (19)$$

249 More generally, we denote the Expected Improvement as $h^{\text{TPE}}(\lambda; \theta) := \text{EI}^{\text{TPE}}(\lambda; \theta)$, where h^{TPE}
250 represents either the exact acquisition induced by the density model or a practical approximation.

251 Unlike the GP case, differentiability of $g_t(\theta)$ may not be possible. We therefore consider a more gen-
252 eral control-variate construction. Let $c_t(\theta) \in \mathbb{R}^{d_c}$ denote any random vector satisfying $\mathbb{E}_{q_t}[c_t(\theta)] = 0$,
253 such as centered bootstrap statistics, centered density ratio proxies, centered bandwidth or split sum-
254 maries, or score functions when they exist. We then obtain our orthogonalized estimator

$$\widehat{\text{EI}}^{\text{TPE,orth}}(\lambda) = \frac{1}{S} \sum_{s=1}^S [h^{\text{TPE}}(\lambda; \theta^{(s)}) - \gamma^{\text{TPE}}(\lambda)^\top c_t(\theta^{(s)})], \quad \theta^{(s)} \sim q_t, \quad (20)$$

255 with

$$\gamma^{\text{TPE}}(\lambda) = \text{Cov}(c_t, c_t)^{-1} \text{Cov}(c_t, h^{\text{TPE}}(\lambda; \theta)). \quad (21)$$

256 4.3 ORTHOBO optimization procedure

257 We present the full optimization in Alg. 1.

258 **Ensemble surrogates.** To mitigate structural
259 surrogate misspecification, we consider an en-
260 semble of surrogate models $m = 1, \dots, M$.
261 Such ensembles can consist of Gaussian pro-
262 cesses with different kernels, models with vary-
263 ing likelihood assumptions, various tree-based
264 or neural surrogates. For each model m , we
265 obtain the respective orthogonalized acquisition
266 estimate $\widehat{\text{EI}}_m^{\text{orth}}(\lambda)$. We then ensemble the ac-
267 quisition as a weighted average of weights $\pi_{m,t}$
268

$$\widehat{\text{EI}}^{\text{ens}}(\lambda) = \sum_{m=1}^M \pi_{m,t} \widehat{\text{EI}}_m^{\text{orth}}(\lambda). \quad (22)$$

269 **Exponentially weighted aggregation.** We up-
270 date ensemble weights through tempered ex-
271 ponentially weighted aggregation [e.g., 25], a
272 widely used approach in model aggregation
273 [18, 32, 52]. A minimum weight floor [28] pre-
274 vents premature collapse onto a single surrogate

$$\pi_{m,t} \propto \max(\delta, \pi_{m,t-1} \exp(\ell_{m,t}/\tau)), \quad (23)$$

Algorithm 1: ORTHOBO

Input: Initial data size n_0 , budget T , number of surrogate models
 M , candidate set generator \mathcal{C}_t , numerical floor $\alpha > 0$,
objective function f
Initialize $\mathcal{D}_0 = \{(\lambda_i, y_i)\}_{i=1}^{n_0}$ using random hyperparameter
configurations;
for $t = 0, 1, \dots, T - 1$ **do**
 for $m = 1, \dots, M$ **do**
 Fit surrogate model m on \mathcal{D}_t ;
 Obtain approximate posterior $q_{m,t}(\theta)$;
 end
 Update ensemble weights $\pi_{m,t}$ for $m = 1, \dots, M$;
 foreach $\lambda \in \mathcal{C}_t$ **do**
 for $m = 1, \dots, M$ **do**
 Orthogonalization: Compute orthogonalized
 marginal EI estimate $\widehat{\text{EI}}_m^{\text{orth}}(\lambda)$ via Eq. (6);
 end
 Ensemble strategy:
 $\widehat{\text{EI}}^{\text{ens}}(\lambda) = \sum_{m=1}^M \pi_{m,t} \widehat{\text{EI}}_m^{\text{orth}}(\lambda)$
 end
 Outer log transformation:
 $\lambda_{t+1} = \arg \max_{\lambda \in \mathcal{C}_t} \log(\max(\widehat{\text{EI}}^{\text{ens}}(\lambda), \alpha))$
 Evaluate $y_{t+1} = f(\lambda_{t+1})$; update
 $\mathcal{D}_{t+1} = \mathcal{D}_t \cup \{(\lambda_{t+1}, y_{t+1})\}$;
end
Output: Best observed configuration in \mathcal{D}_T

275 where $\ell_{m,t}$ is a predictive log score obtained from m_t using (λ_{t+1}, y_{t+1}) , $\tau > 0$ a temperature
 276 parameter, and $\delta > 0$ enforces a minimum weight. This results in a robust aggregation rule that
 277 balances exploitation of well-performing surrogates with persistent exploration across model classes.
 278 Note, that our overall methodology is independent and can be combined with other update strategies.

279 **Optimization based on outer log stabilization.** To improve numerical stability, we follow recent
 280 literature [2] and optimize the *outer log transformed* acquisition

$$\widehat{L}^{\text{ens}}(\lambda_t) = \log(\max(\widehat{\text{EI}}^{\text{ens}}(\lambda_t), \alpha)), \quad (24)$$

281 where $\alpha > 0$ ensures that the logarithm is well defined. We select the next hyperparameter configura-
 282 tion $\lambda_{t+1} = \arg \max_{\lambda \in \mathcal{C}_t} \widehat{L}^{\text{ens}}(\lambda)$, and evaluate it on the true objective to obtain y_{t+1} .

283 5 Empirical results

284 **Experimental setup.** We evaluate ORTHOBO against various baselines from the literature. Specifi-
 285 cally, we compare ORTHOBO, qLogEI [2] (main comparison method), UCB [62], TuRBO [16], and
 286 Sobol sampling [e.g. 3] as a random baseline. Our experimental settings follow the literature [e.g.
 287 2, 3]: unless otherwise noted, all experiments use 16 replications, $S = 512$, and we start all methods
 288 from the same $n_0 = 32$ Sobol points. In optimizing the acquisition functions, we use 512 raw
 289 samples and allow 8 restarts to prevent getting stuck in local minima [56]. To prevent confounding,
 290 we use $M = 1$ throughout the experiments in the main paper. We consider four standard synthetic
 291 regression problems with known global optima and varying dimensionality: Hartmann6, Ackley8,
 292 Michalewicz10, and Levy16 [e.g., 2, 10, 15, 41]. Details of the benchmark functions and surrogate
 293 choices are provided in Supplement E. We compare methods using best-so-far regret and use 95% CI
 294 from standard normal approximation.

295 **Evaluation.** Our goal is to demonstrate the theoretical properties of ORTHOBO and show its
 296 superiority on real-world training and fine-tuning tasks, answering four distinct research questions.
 297 We first isolate the effect of orthogonalization on acquisition estimation to RQ1) show the *variance*
 298 *reduction* and RQ2) *ranking stability* properties of ORTHOBO. Then, we evaluate the resulting
 299 RQ3) *efficiency gain* in terms of MC sampling budget. Finally, we show RQ4) the improved
 300 *downstream utility*. We present further experiments on misspecified surrogate families, weakly fitted
 301 hyperparameters or surrogates, and ensemble experiments in Supplement F.

302 **RQ1) Variance reduction.** We probe the
 303 acquisition function with random Sobol
 304 samples and compare raw and orthogonal-
 305 ized acquisition estimate variance. We
 306 use a GP surrogate. We report results
 307 for a Matérn-5/2 kernel with ARD on
 308 Michalewicz10 and Levy16 in Table 1; re-
 309 sults for a RBF kernel and a TPE surrogate are in Supplement F.1. \Rightarrow *Across all settings, orthogo-*
 310 *nalization reduces the variance, confirming that ORTHOBO substantially stabilizes the acquisition*
 311 *estimate.*

Table 1: Variance reduction on Sobol probes through orthogonalization.

S	Michalewicz10		Levy16	
	$\text{Var}(\widehat{\text{EI}}_m^{\text{MC}})$	$\text{Var}(\widehat{\text{EI}}_m^{\text{ortho}})$	$\text{Var}(\widehat{\text{EI}}_m^{\text{MC}})$	$\text{Var}(\widehat{\text{EI}}_m^{\text{ortho}})$
8	9.51e-08	4.54e-08	0.00106	0.000709
32	1.67e-08	1.07e-09	0.000208	2.27e-05

312 **RQ2) Ranking stability.**
 313 We repeatedly evaluate the
 314 acquisition function on the
 315 same fixed Sobol probe set
 316 and compare the rankings
 317 across repetitions. We re-
 318 port: (i) *probe variance*:
 319 mean variance of the acqui-
 320 sition value over the probe

Table 2: Ranking stability (\pm std), via Sobol probes. Lower variance and flip rate indicate more stable acquisition estimates; higher top-1 agreement indicates more stable candidate rankings across probes.

	Method	Probe	Top1	Flip	Regret
		variance \downarrow	agreement \uparrow	rate \downarrow	\downarrow
Michalewicz10	qLogEI	168.603 _{154.76}	0.925 _{0.16}	0.148 _{0.12}	7.140 _{0.01}
	OrthoBO (ours)	0.019 _{0.03}	0.988 _{0.06}	0.014 _{0.03}	6.961 _{0.36}
Levy16	qLogEI	233.655 _{287.67}	0.917 _{0.12}	0.069 _{0.09}	82.194 _{6.05}
	OrthoBO (ours)	0.073 _{0.21}	0.988 _{0.06}	0.048 _{0.04}	58.938 _{0.00}

321 points across repeated acquisition evaluations; (ii) *top1 agreement*: fraction of repeats that select the
 322 same top-ranked probe point, (iii) the *flip rate* of adjacent high-ranked candidate pairs, and (iv) *regret*.
 323 The flip rate is particularly close to Proposition 4.2, as it measures how often local pairwise orderings
 324 change due to acquisition-estimation noise. The main results are in Table 2; further results are in
 325 Supplement F.1. \Rightarrow *Our variance reduction yields more stable rankings.*

326 **RQ3) MC efficiency.** We use an isotropic RBF kernel, a common smooth GP prior, combining
 327 moderate surrogate mismatch with reduced acquisition-estimation budgets. We show the main results

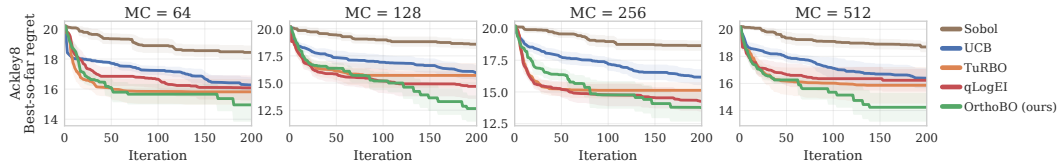


Figure 3: Best-so-far regret for different MC budgets $S \Rightarrow$ ORTHOBO achieves the lowest regrets across all budgets.

for varying S on Ackley8 in Fig. 3; further results are in Fig 12 in Supplement F. Reducing the MC budget makes the performance of acquisition-based methods more sensitive to estimation noise. Orthogonalization is most effective in earlier iterations, where MC estimation noise has a larger effect on the acquisition values. Across all benchmarks, our ORTHOBO remains competitive and frequently achieves the lowest regret.

RQ4) Downstream utility. We show how ORTHOBO improves hyperparameter optimization on two real-world tasks: (i) training a neural network and (ii) optimizing fine-tuning hyperparameters of a pre-trained vision transformer (ViT) [14] for manufacturing image classification.

Neural network. We employ the 5D CNN benchmark by Ament et al. [3] and include their RPR method as comparison. Details are in Supplement E.3. The corruption protocol prematurely stops training, providing unreliable objective values which can distort the surrogate fit, especially early in BO when only few observations are available. Results are in Fig. 4. Despite the corrupted observations, ORTHOBO improves steadily over the optimization horizon and attains the highest value.

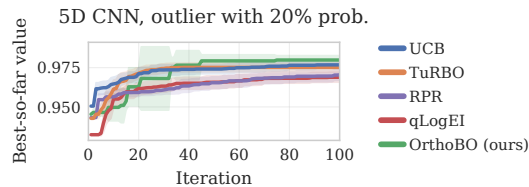


Figure 4: 5D CNN with corrupted evaluations. \Rightarrow ORTHOBO achieves the **strongest final performance**.

Vision transformer. We optimize five fine-tuning hyperparameters on the industrial WM811K wafer-map dataset [33, 64] based on the F1 score. Details are provided in Appendix E.4. We present results in Fig. 5. ORTHOBO improves rapidly within the first few evaluations and maintains the highest best-so-far validation F1 score throughout most of the optimization horizon. Compared with qLogEI, ORTHOBO improves the best observed validation F1 score by up to 20 percentage points. This demonstrates that the proposed acquisition-stabilization mechanism also improves performance in a practical, high-cost HPO setting beyond synthetic benchmarks.

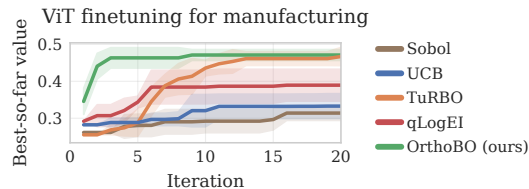


Figure 5: ViT. \Rightarrow ORTHOBO **improves faster than the baselines** and achieves the **highest best-so-far validation score**.

6 Discussion

We introduced ORTHOBO, an orthogonalized framework for stabilizing acquisition estimation in Bayesian HPO. BO decisions depend on the ranking of estimated acquisition values, so MC noise in these estimates can directly induce ranking errors and lead to suboptimal evaluations. ORTHOBO addresses this failure mode by reducing variance while preserving the underlying marginal EI target. Our results show that this variance reduction improves ranking stability and translates into better downstream BO performance. ORTHOBO has limitations. First, it introduces additional computational overhead. However, this overhead is typically small compared with expensive objective evaluations in HPO (see Supplement C). Second, we focus only on EI-based acquisition functions. Note that the same variance-reduction principle can be extended to other acquisition functions but may require acquisition-specific orthogonalization schemes (see Supplement B.1). Overall, our work identifies acquisition-estimation noise as a distinct and practically relevant failure mode in BO. By reducing this noise, ORTHOBO stabilizes acquisition rankings and improves the quality of sequential hyperparameter decisions.

374 **References**

375 [1] A. AlBahar, I. Kim, and X. Yue. A robust asymmetric kernel function for Bayesian optimization,
 376 with application to image defect detection in manufacturing systems. *IEEE Transactions on*
 377 *Automation Science and Engineering*, 19(4):3222–3233, 2021.

378 [2] S. Ament, S. Daulton, D. Eriksson, M. Balandat, and E. Bakshy. Unexpected improvements to
 379 expected improvement for Bayesian optimization. In *Neural Information Processing Systems*
 380 *(NeurIPS)*, 2023.

381 [3] S. Ament, E. Santorella, D. Eriksson, B. Letham, M. Balandat, and E. Bakshy. Robust Gaussian
 382 processes via relevance pursuit. In *Neural Information Processing Systems (NeurIPS)*, 2024.

383 [4] M. Balandat, B. Karrer, D. Jiang, S. Daulton, B. Letham, A. G. Wilson, and E. Bakshy.
 384 Botorch: A framework for efficient Monte-Carlo Bayesian optimization. In *Neural Information*
 385 *Processing Systems (NeurIPS)*, 2020.

386 [5] T. Ban, M. Ohue, and Y. Akiyama. Efficient hyperparameter optimization by using Bayesian opti-
 387 mization for drug-target interaction prediction. In *International Conference on Computational*
 388 *Advances in Bio and Medical Sciences (ICCABS)*, 2017.

389 [6] J. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl. Algorithms for hyper-parameter optimization.
 390 In *Neural Information Processing Systems (NeurIPS)*, 2011.

391 [7] J. Bergstra, D. Yamins, and D. Cox. Making a science of model search: Hyperparameter
 392 optimization in hundreds of dimensions for vision architectures. In *International Conference on*
 393 *Machine Learning (ICML)*, 2013.

394 [8] F. Berkenkamp, A. P. Schoellig, and A. Krause. No-regret Bayesian optimization with unknown
 395 hyperparameters. *Journal of Machine Learning Research*, 20(50):1–24, 2019.

396 [9] P. J. Bickel, C. A. Klaassen, P. J. Bickel, Y. Ritov, J. Klaassen, J. A. Wellner, and Y. Ritov.
 397 *Efficient and adaptive estimation for semiparametric models*, volume 4. Springer, 1993.

398 [10] E. Bodin, M. Kaiser, I. Kazlauskaitė, Z. Dai, N. Campbell, and C. H. Ek. Modulating surrogates
 399 for Bayesian optimization. In *International Conference on Machine Learning (ICML)*, 2020.

400 [11] I. Bogunovic and A. Krause. Misspecified Gaussian process bandit optimization. In *Neural*
 401 *Information Processing Systems (NeurIPS)*, 2021.

402 [12] V. Chernozhukov, D. Chetverikov, M. Demirer, E. Duflo, C. Hansen, W. Newey, and J. Robins.
 403 Double/debiased machine learning for treatment and structural parameters. *The Econometrics*
 404 *Journal*, 21:C1–C68, 2018.

405 [13] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A large-scale hierarchical
 406 image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages
 407 248–255. Ieee, 2009.

408 [14] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani,
 409 M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for
 410 image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

411 [15] D. Eriksson and M. Poloczek. Scalable constrained Bayesian optimization. In *International*
 412 *Conference on Artificial Intelligence and Statistics (AISTATS)*, 2021.

413 [16] D. Eriksson, M. Pearce, J. Gardner, R. D. Turner, and M. Poloczek. Scalable global optimization
 414 via local Bayesian optimization. In *Neural Information Processing Systems (NeurIPS)*, 2019.

415 [17] S. Falkner, A. Klein, and F. Hutter. BOHB: Robust and efficient hyperparameter optimization at
 416 scale. In *International Conference on Machine Learning (ICML)*, 2018.

417 [18] W. Fedus, B. Zoph, and N. Shazeer. Switch transformers: Scaling to trillion parameter models
 418 with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39, 2022.

419 [19] M. Feurer, A. Klein, K. Eggenberger, J. Springenberg, M. Blum, and F. Hutter. Efficient and
 420 robust automated machine learning. In *Neural Information Processing Systems (NeurIPS)*,
 421 2015.

422 [20] D. J. Foster and V. Syrgkanis. Orthogonal statistical learning. *The Annals of Statistics*, 51(3):
 423 879–908, 2023.

- 424 [21] D. Frauen, M. Schröder, K. Hess, and S. Feuerriegel. Orthogonal survival learners for estimating
425 heterogeneous treatment effects from time-to-event data. In *Neural Information Processing*
426 *Systems (NeurIPS)*, 2025.
- 427 [22] D. Frauen, A. Deviyani, M. van der Schaar, and S. Feuerriegel. Nonparametric LLM evaluation
428 from preference data. *arXiv preprint*, arXiv:2601.21816, 2026.
- 429 [23] P. I. Frazier. Bayesian optimization. In *Recent advances in optimization and modeling of*
430 *contemporary problems*, pages 255–278. 2018.
- 431 [24] P. I. Frazier. A tutorial on Bayesian optimization. *arXiv preprint*, arXiv:1807.02811, 2018.
- 432 [25] Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an
433 application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997.
- 434 [26] J. González, Z. Dai, P. Hennig, and N. Lawrence. Batch Bayesian optimization via local
435 penalization. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*,
436 2016.
- 437 [27] H. Ha, V. Nguyen, H. Tran-The, H. Zhang, X. Zhang, and A. v. d. Hengel. Provably efficient
438 Bayesian optimization with unknown Gaussian process hyperparameter estimation. *arXiv*
439 *preprint*, arXiv:2306.06844, 2023.
- 440 [28] M. Herbster and M. K. Warmuth. Tracking the best expert. *Machine learning*, 32(2):151–178,
441 1998.
- 442 [29] J. M. Hernández-Lobato, M. W. Hoffman, and Z. Ghahramani. Predictive entropy search for
443 efficient global optimization of black-box functions. In *Neural Information Processing Systems*
444 *(NeurIPS)*, 2014.
- 445 [30] J. M. Hernández-Lobato, M. Gelbart, M. Hoffman, R. Adams, and Z. Ghahramani. Predictive
446 entropy search for Bayesian optimization with unknown constraints. In *International Conference*
447 *on Machine Learning (ICML)*, 2015.
- 448 [31] K. Hess, D. Frauen, N. Kilbertus, and S. Feuerriegel. Debaised neural operators for estimating
449 functionals. *arXiv preprint*, arXiv:2604.19296, 2026.
- 450 [32] R. A. Jacobs, M. I. Jordan, S. J. Nowlan, and G. E. Hinton. Adaptive mixtures of local experts.
451 *Neural Computation*, 3(1):79–87, 1991.
- 452 [33] J.-S. R. Jang. Mir-wm811k: Dataset for wafer map failure pattern recognition. [http://](http://mir1lab.org/dataset/public/)
453 mir1lab.org/dataset/public/, 2015.
- 454 [34] D. R. Jones, M. Schonlau, and W. J. Welch. Efficient global optimization of expensive black-box
455 functions. *Journal of Global Optimization*, 13(4):455–492, 1998.
- 456 [35] E. H. Kennedy. Semiparametric doubly robust targeted double machine learning: a review.
457 *Handbook of Statistical Methods for Precision Medicine*, pages 207–236, 2024.
- 458 [36] Q. Lu, K. D. Polyzos, B. Li, and G. B. Giannakis. Surrogate modeling for Bayesian optimization
459 beyond a single Gaussian process. *IEEE Transactions on Pattern Analysis and Machine*
460 *Intelligence*, 45(9):11283–11296, 2023.
- 461 [37] L. Mackey, V. Syrgkanis, and I. Zadik. Orthogonal machine learning: Power and limitations. In
462 *International Conference on Machine Learning (ICML)*, 2018.
- 463 [38] R. Martinez-Cantin, K. Tee, and M. McCourt. Practical Bayesian optimization in the presence
464 of outliers. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*,
465 2018.
- 466 [39] V. Melnychuk and S. Feuerriegel. GDR-learners: Orthogonal learning of generative models for
467 potential outcomes. *arXiv preprint*, arXiv:2509.22953, 2025.
- 468 [40] H. B. Moss, D. S. Leslie, J. Gonzalez, and P. Rayson. GIBBON: General-purpose information-
469 based Bayesian optimisation. *Journal of Machine Learning Research*, 22(235):1–49, 2021.
- 470 [41] H. B. Moss, S. W. Ober, and V. Picheny. Inducing point allocation for sparse Gaussian processes
471 in high-throughput Bayesian optimisation. In *International Conference on Artificial Intelligence*
472 *and Statistics (AISTATS)*, 2023.
- 473 [42] V. Nath, D. Yang, A. Hatamizadeh, A. A. Abidin, A. Myronenko, H. R. Roth, and D. Xu.
474 The power of proxy data and proxy networks for hyper-parameter optimization in medical
475 image segmentation. In *International Conference on Medical Image Computing and Computer-*
476 *Assisted Intervention*, 2021.

- 477 [43] W. Neiswanger and A. Ramdas. Uncertainty quantification using martingales for misspecified
478 Gaussian processes. In *International Conference on Algorithmic Learning Theory (ALT)*, 2021.
- 479 [44] X. Nie and S. Wager. Quasi-oracle estimation of heterogeneous treatment effects. *Biometrika*,
480 108(2):299–319, 2021.
- 481 [45] M. Oprescu, V. Syrgkanis, and Z. S. Wu. Orthogonal random forest for causal inference. In
482 *International Conference on Machine Learning (ICML)*, 2019.
- 483 [46] K. D. Polyzos, Q. Lu, and G. B. Giannakis. Bayesian optimization with ensemble learning
484 models and adaptive expected improvement. In *International Conference on Acoustics, Speech
485 and Signal Processing (ICASSP)*, 2023.
- 486 [47] A. Quitadadmo, J. Johnson, and X. Shi. Bayesian hyperparameter optimization for machine
487 learning based eqtl analysis. In *International Conference on Bioinformatics, Computational
488 Biology, and Health Informatics*, 2017.
- 489 [48] J. M. Robins, A. Rotnitzky, and L. P. Zhao. Estimation of regression coefficients when some
490 regressors are not always observed. *Journal of the American Statistical Association*, 89(427):
491 846–866, 1994.
- 492 [49] J. M. Robins, M. A. Hernan, and B. Brumback. Marginal structural models and causal inference
493 in epidemiology. *Epidemiology*, 11(5):550–560, 2000.
- 494 [50] M. Schröder, V. Melnychuk, and S. Feuerriegel. Differentially private learners for heterogeneous
495 treatment effects. *International Conference on Learning Representations (ICLR)*, 2025.
- 496 [51] B. Shahriari, K. Swersky, Z. Wang, R. P. Adams, and N. De Freitas. Taking the human out of
497 the loop: A review of Bayesian optimization. *Proceedings of the IEEE*, 104(1):148–175, 2016.
- 498 [52] N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hinton, and J. Dean. Outrageously
499 large neural networks: The sparsely-gated mixture-of-experts layer. In *International Conference
500 on Learning Representations (ICLR)*, 2017.
- 501 [53] J. Snoek, H. Larochelle, and R. P. Adams. Practical Bayesian optimization of machine learning
502 algorithms. In *Neural Information Processing Systems (NeurIPS)*, 2012.
- 503 [54] N. Srinivas, A. Krause, S. M. Kakade, and M. Seeger. Gaussian process optimization in the
504 bandit setting: No regret and experimental design. In *International Conference on Machine
505 Learning (ICML)*, 2010.
- 506 [55] L. Tani and C. Veelken. Comparison of Bayesian and particle swarm algorithms for hyper-
507 parameter optimisation in machine learning applications in high energy physics. *Computer
508 Physics Communications*, 294:108955, 2024.
- 509 [56] A. Törn and A. Žilinskas. *Global optimization*, volume 350. Springer, 1989.
- 510 [57] R. Turner, D. Eriksson, M. McCourt, J. Kiili, E. Laaksonen, Z. Xu, and I. Guyon. Bayesian
511 optimization is superior to random search for machine learning hyperparameter tuning: Analysis
512 of the black-box optimization challenge 2020. In *NeurIPS Competition and Demonstration
513 Track*, 2020.
- 514 [58] M. J. Van Der Laan and D. Rubin. Targeted maximum likelihood learning. *The International
515 Journal of Biostatistics*, 2, 2006.
- 516 [59] Z. Wang and S. Jegelka. Max-value entropy search for efficient Bayesian optimization. In
517 *International conference on machine learning*, pages 3627–3635. PMLR, 2017.
- 518 [60] S. Watanabe. Tree-structured parzen estimator: Understanding its algorithm components and
519 their roles for better empirical performance. *arXiv preprint*, arXiv:2304.11127, 2023.
- 520 [61] J. Wilson, F. Hutter, and M. Deisenroth. Maximizing acquisition functions for Bayesian
521 optimization. In *Neural Information Processing Systems (NeurIPS)*, 2018.
- 522 [62] J. T. Wilson, R. Moriconi, F. Hutter, and M. P. Deisenroth. The reparameterization trick for
523 acquisition functions. *arXiv preprint*, arXiv:1712.00424, 2017.
- 524 [63] J. Wu, S. Toscano-Palmerin, P. I. Frazier, and A. G. Wilson. Practical multi-fidelity Bayesian
525 optimization for hyperparameter tuning. In *Uncertainty in Artificial Intelligence (UAI)*, 2020.
- 526 [64] M.-J. Wu, J.-S. R. Jang, and J.-L. Chen. Wafer map failure pattern recognition and similarity
527 ranking for large-scale data sets. *IEEE Transactions on Semiconductor Manufacturing*, 28(1):
528 1–12, 2015.

- 529 [65] Y. Yuan, W. Wang, and W. Pang. A systematic comparison study on hyperparameter optimisation
530 of graph neural networks for molecular property prediction. In *Proceedings of the genetic and*
531 *evolutionary computation conference*, 2021.

532 A Extended related work

533 Below, we discuss more distant related literature, completing the discussion in Section 2. We first
534 cover literature regarding the robustness of the three components of BO in our work, and finally
535 present related work on orthogonal machine learning.

536 **Marginalization over surrogate uncertainty.** A standard source of instability in BO is uncertainty
537 in the surrogate itself, especially in the hyperparameters such as kernel lengthscales. Prior work has
538 emphasized that a fully Bayesian treatment of these hyperparameters can substantially improve ro-
539 bustness relative to empirical Bayes or plug-in optimization [53]. Related work has also incorporated
540 hyperparameter uncertainty directly into acquisition design [29, 30]. While these approaches account
541 for surrogate uncertainty by integrating over posterior uncertainty in the model, they do not directly
542 address the additional MC instability that arises when the resulting acquisition value must itself be
543 estimated.

544 **Ensemble surrogates and model aggregation.** A natural strategy for mitigating surrogate mis-
545 specification is to replace a single surrogate with an ensemble of surrogate models. Recent work
546 has proposed GP ensembles and mixture-based surrogate models for BO, showing that combining
547 multiple kernels or surrogate classes can improve robustness relative to committing to a single GP
548 prior [36, 46]. These approaches address only structural misspecification at the model level and do
549 not account for instability in the individual surrogates.

550 **Acquisition functions and numerical stability.** Acquisition functions are the core decision rule
551 in BO, with classical choices including Probability of Improvement, Expected Improvement, and
552 upper-confidence-based methods [34, 54]. Recent work has highlighted that vanilla EI and its variants
553 can be numerically unstable, which can substantially degrade optimization performance [2]. This has
554 motivated numerically stable reformulations based on the logarithmic scale, which are substantially
555 easier to optimize in practice.

556 **Orthogonal machine learning.** Orthogonal machine learning is rooted in semi-parametric statistics
557 [9], and is applied in many settings where the estimation target is defined based on unknown nuisance
558 components [12, 20]. Very common examples are found in the missing data analysis [48] and causal
559 inference literature [35, 44, 45, 49, 58]. Recently, orthogonal learning has also been expanded to
560 other fields such as private machine learning [50], survival analysis [21], LLM evaluation [22], and
561 general ML optimization topics such as neural operators or generative models [31, 39].

562 **B Extension to constrained and parallel acquisitions**

563 In Section 4, we develop an orthogonalized estimator for the standard Expected Improvement setting.
 564 The same construction extends to a broader class of acquisition functions, including constrained
 565 EI, parallel or batch EI, and their combinations. Below, we first provide a general extension to our
 566 standard EI setting in the main paper and then specialize it to constrained and parallel Bayesian
 567 optimization.

568 **B.1 General acquisition-level extension**

569 Let z denote a generic acquisition input. Depending on the setting, we write z to represent (i) a single
 570 candidate point λ in standard BO, (ii) a single candidate point λ in constrained BO, or (iii) a batch
 571 $\Lambda = (\lambda_1, \dots, \lambda_q)$ in parallel BO.

572 For surrogate model m with parameter θ_m and approximate posterior $q_{m,t}(\theta)$, let $a_m(z; \theta_m)$ denote a
 573 generic acquisition function computed under parameter value θ_m . We define the marginal acquisition
 574 target

$$A_m(z) := \mathbb{E}_{\theta_m \sim q_{m,t}}[a_m(z; \theta_m)]. \quad (25)$$

575 As in Section 4, let $g_m(\theta_m) := \nabla_{\theta_m} \log q_{m,t}(\theta_m)$ denote the score function. With $\theta_m^{(s)} \stackrel{\text{i.i.d.}}{\sim} q_{m,t}$ the
 576 orthogonalized estimator is given by

$$\widehat{A}_m^{\text{orth}}(z) = \frac{1}{S} \sum_{s=1}^S \left[a_m(z; \theta_m^{(s)}) - (\text{Cov}(g_m, g_m)^{-1} \text{Cov}(g_m, a_m))^\top g_m(\theta_m^{(s)}) \right]. \quad (26)$$

577 **Proposition B.1** (Generic extension of the orthogonalized estimator). *Fix an iteration t , a surrogate*
 578 *model m , and an acquisition input z . Assume that*

$$\mathbb{E}_{q_{m,t}}[\|g_m(\theta_m)\|_2^2] < \infty, \quad \mathbb{E}_{q_{m,t}}[a_m(\theta_m; z)^2] < \infty, \quad (27)$$

579 *and that $\text{Cov}(g_m, g_m)$ is nonsingular. Then:*

- 580 1. $\widehat{A}_m^{\text{orth}}(z)$ is unbiased for $A_m(z)$
 581 2.

$$\text{Var}\left(\widehat{A}_m^{\text{orth}}(z)\right) = \frac{1}{S} \left(\text{Var}(h_m) - \text{Cov}(h_m, g_m)^\top \text{Cov}(g_m, g_m)^{-1} \text{Cov}(g_m, h_m) \right), \quad (28)$$

581 *and therefore*

$$\text{Var}\left(\widehat{A}_m^{\text{orth}}(z)\right) \leq \text{Var}\left(\widehat{A}_m^{\text{MC}}(z)\right), \quad (29)$$

582 *where $\widehat{A}_m^{\text{MC}}(z) := \frac{1}{S} \sum_{s=1}^S h_m(\theta_m^{(s)}; z)$*

- 583 3. *The residual $r_m(\theta_m; z) := a_m(\theta_m; z) - (\text{Cov}(g_m, g_m)^{-1} \text{Cov}(g_m, a_m))^\top g_m(\theta_m)$ is or-*
 584 *thogonal to the score directions.*

585 *Proof.* The proof is identical to that of the EI case in Supplement D after replacing $\text{EI}_m(\lambda; \theta_m)$ by
 586 the generic acquisition value $a_m(z; \theta_m)$. \square

587 **Remark.** Proposition B.1 shows that the theoretical properties of our orthogonalized approach are
 588 acquisition-agnostic. We only require square integrability of the acquisition value under $q_{m,t}$ and the
 589 usual score-function regularity conditions.

590 **B.2 Constrained Expected Improvement**

591 Consider constrained BO with one objective f_1 and constraint functions f_2, \dots, f_M , where, without
 592 loss of generality, the feasible set is defined by

$$f_i(\lambda) \leq 0, \quad i = 2, \dots, M. \quad (30)$$

593 A standard constrained Expected Improvement (CEI) acquisition is

$$\text{CEI}_m(\lambda; \theta_m) := \mathbb{E} \left[(f^* - f_1(\lambda))_+ \prod_{i=2}^M \mathbf{1}\{f_i(\lambda) \leq 0\} \middle| \mathcal{D}_t, \theta_m \right], \quad (31)$$

594 where f^* denotes the current best *feasible* objective value (i.e., the minimum of f_1 over feasible
595 observations).

596 The corresponding marginal target is $\text{CEI}_m^{\text{marg}}(\lambda) := \mathbb{E}_{\theta_m \sim q_{m,t}}[\text{CEI}_m(\lambda; \theta_m)]$. Defining
597 $h_m^{\text{CEI}}(\theta_m; \lambda) := \text{CEI}_m(\lambda; \theta_m)$, the orthogonalized CEI estimator is given by

$$\widehat{\text{CEI}}_m^{\text{orth}}(\lambda) = \frac{1}{S} \sum_{s=1}^S \left[\text{CEI}_m(\lambda; \theta_m^{(s)}) - (\text{Cov}(g_m, g_m)^{-1} \text{Cov}(g_m, h_m^{\text{CEI}}))^{\top} g_m(\theta_m^{(s)}) \right]. \quad (32)$$

598 **Corollary B.2** (Orthogonalized constrained EI). *Assume $\mathbb{E}_{q_{m,t}}[(\text{CEI}_m(\lambda; \theta_m))^2] < \infty$. Then
599 Proposition B.1 applies with $a_m(z; \theta_m) = \text{CEI}_m(\lambda; \theta_m)$.*

600 *Proof.* The indicator product is bounded by 1, so $0 \leq \text{CEI}_m(\lambda; \theta_m) \leq \text{EI}_{m,1}(\lambda; \theta_m) :=$
601 $\mathbb{E}[(f^* - f_1(\lambda))_+ | \mathcal{D}_t, \theta_m]$. Square integrability of $\text{EI}_{m,1}$ thus implies square integrability of CEI_m ,
602 and the corollary follows from Proposition B.1. \square

603 B.3 Parallel or batch Expected Improvement

604 In parallel BO, we select a batch $\Lambda = (\lambda_1, \dots, \lambda_q)$ of q candidate points at each iteration. The
605 standard batch expected improvement (qEI) is given by

$$\text{qEI}_m(\Lambda; \theta_m) := \mathbb{E} \left[\max_{j=1, \dots, q} (f^* - f(\lambda_j))_+ \middle| \mathcal{D}_t, \theta_m \right], \quad (33)$$

606 with corresponding marginal target is $\text{qEI}_m^{\text{marg}}(\Lambda) := \mathbb{E}_{\theta_m \sim q_{m,t}}[\text{qEI}_m(\Lambda; \theta_m)]$.

607 We define $h_m^{\text{qEI}}(\theta_m; \Lambda) := \text{qEI}_m(\Lambda; \theta_m)$. Then the orthogonalized qEI estimator is given by

$$\widehat{\text{qEI}}_m^{\text{orth}}(\Lambda) = \frac{1}{S} \sum_{s=1}^S \left[\text{qEI}_m(\Lambda; \theta_m^{(s)}) - (\text{Cov}(g_m, g_m)^{-1} \text{Cov}(g_m, h_m^{\text{qEI}}))^{\top} g_m(\theta_m^{(s)}) \right]. \quad (34)$$

608 **Corollary B.3** (Orthogonalized batch EI). *Assume $\mathbb{E}_{q_{m,t}}[(\text{qEI}_m(\Lambda; \theta_m))^2] < \infty$. Then Proposi-
609 tion B.1 applies with $a_m(z; \theta_m) = \text{qEI}_m(\Lambda; \theta_m)$.*

610 *Proof.* Observe that for every fixed batch $X = (x_1, \dots, x_q)$ and parameter value θ , $0 \leq$
611 $\text{qEI}_m(X; \theta) \leq \sum_{j=1}^q \text{EI}_m(x_j; \theta)$. Consequently, $\text{qEI}_m(X; \theta)^2 \leq \left(\sum_{j=1}^q \text{EI}_m(x_j; \theta) \right)^2$. In par-
612 ticular, if $\mathbb{E}_{\theta \sim q_{m,t}} \left[\left(\sum_{j=1}^q \text{EI}_m(x_j; \theta) \right)^2 \right] < \infty$, then as well $\mathbb{E}_{\theta \sim q_{m,t}} [\text{qEI}_m(X; \theta)^2] < \infty$. The
613 corollary now directly follows from Proposition B.1. \square

614 Unlike standard EI, qEI does not admit a closed-form expression. Therefore, we follow common
615 practice and approximate it by inner MC sampling. Let

$$\widehat{\text{qEI}}_m(\Lambda; \theta_m) = \frac{1}{N} \sum_{n=1}^N \max_{j=1, \dots, q} (f^* - \xi^{(n)}(\lambda_j))_+, \quad \xi^{(n)} \sim p(\cdot | \mathcal{D}_t, \theta_m), \quad (35)$$

616 where the inner MC estimator is assumed to be unbiased for $\text{qEI}_m(\Lambda; \theta_m)$.

617 The resulting orthogonalized MC estimator is then given by

$$\widehat{\text{qEI}}_{m, \text{nested}}^{\text{orth}}(\Lambda) = \frac{1}{S} \sum_{s=1}^S \left[\widehat{\text{qEI}}_m(\Lambda; \theta_m^{(s)}) - (\text{Cov}(g_m, g_m)^{-1} \text{Cov}(g_m, \widehat{\text{qEI}}(\Lambda; \theta_m)))^{\top} g_m(\theta_m^{(s)}) \right]. \quad (36)$$

618 Importantly, under the assumption that the inner estimator (35) is unbiased for $\text{qEI}_m(X; \theta_m)$ for
619 every fixed θ_m , and that all required second moments exist, then $\widehat{\text{qEI}}_{m,\text{nested}}^{\text{orth}}(\Lambda)$ is unbiased for
620 $\text{qEI}_m(X; \theta_m)$ as well due to the law of iterated expectations.

621 **C Theoretical properties**

622 **C.1 Computational Complexity**

623 **GP surrogate.** For an exact GP surrogate, fitting requires $O(t^3)$ time due to kernel matrix inversion.
 624 Per candidate λ , computing $\mu_t(\lambda; \theta)$ and $\sigma_t^2(\lambda; \theta)$ for a fixed θ costs $O(t^2)$ in a naive implementation,
 625 or $O(t)$ if Cholesky factors are reused appropriately. For S posterior samples, marginal EI estimation
 626 costs $O(St^2)$ or $O(St)$, depending on the implementation.

627 The orthogonalization step introduces two types of cost. First, once per BO iteration, one computes
 628 the score covariance matrix $\widehat{\Sigma}_g := \widehat{\text{Cov}}(g_t, g_t) \in \mathbb{R}^{d_\theta \times d_\theta}$, which depends only on the surrogate
 629 posterior samples and not on the candidate λ . This costs $O(Sd_\theta^2)$, followed by $O(d_\theta^3)$ to invert or
 630 factorize $\widehat{\Sigma}_g$. Second, for each candidate λ , one computes the cross-covariance vector $\widehat{c}_{gh}(\lambda) :=$
 631 $\widehat{\text{Cov}}(g_t, h(\cdot; \lambda))$, which costs $O(Sd_\theta)$, followed by $O(d_\theta^2)$ for a matrix–vector solve or multiplication
 632 with the precomputed factorization of $\widehat{\Sigma}_g^{-1}$. Hence, the additional per-candidate orthogonalization
 633 overhead is $O(Sd_\theta + d_\theta^2)$, and the total per-candidate cost is $O(St^2 + Sd_\theta + d_\theta^2)$, plus a one-time
 634 per-iteration preprocessing cost of $O(Sd_\theta^2 + d_\theta^3)$.

635 We compare the average per-step runtimes of ORTHOBO with qLogEI in Table 3. We observe that
 636 the runtimes vary with the dimensionality of the problem, but the overhead over MC-based qLogEI
 637 is often around 2 s. For our implementation, we used the BoTorch framework, which has made
 638 MC-sampling highly efficient. We thus expect further runtime reductions from a more improved
 639 orthogonalization implementation.

Table 3: Average per-step runtimes (in seconds) for ORTHOBO and qLogEI.

	Hartmann6	Ackley8	Michalewicz10	Levy16
qLogEI	0.696	0.156	0.325	0.582
ORTHOBO	1.607	0.492	2.362	2.563

640 **Tree-based / TPE-style surrogate.** Tree-based or density-based surrogates are typically cheaper
 641 to fit. A single fit often requires $O(t \log t)$ or comparable cost, depending on the underlying model.
 642 Per-candidate acquisition evaluation is usually $O(1)$ or logarithmic in t . As in the GP case, the
 643 covariance matrix of the control variate, $\widehat{\Sigma}_c := \widehat{\text{Cov}}(c_t, c_t) \in \mathbb{R}^{d_c \times d_c}$, depends only on the posterior
 644 samples and can therefore be computed once per BO iteration at cost $O(Sd_c^2)$, followed by $O(d_c^3)$
 645 for inversion or factorization. For each candidate, one then computes the cross-covariance vector
 646 $\widehat{c}_{ch}(\lambda) := \widehat{\text{Cov}}(c_t, h(\cdot; \lambda))$, at cost $O(Sd_c)$, followed by $O(d_c^2)$ for the matrix–vector solve. Hence,
 647 the additional per-candidate orthogonalization overhead is $O(Sd_c + d_c^2)$, with one-time per-iteration
 648 preprocessing cost $O(Sd_c^2 + d_c^3)$.

649 We report average per-step runtimes for orthogonalized TPE versus MC-based TPE in Table 4.
 650 Runtime differences are generally small.

Table 4: Average per-step runtimes (in seconds) for ORTHOBO and MC-based TPE.

	Hartmann6	Ackley8	Michalewicz10	Levy16
MC-TPE	0.224	0.347	0.356	0.754
ORTHOBO	0.390	0.736	2.533	0.360

651 **Ensemble overhead.** For an ensemble of M surrogate models, the total cost scales linearly as
 652 $O\left(\sum_{m=1}^M \text{cost}_m\right)$. In practice, the main computational advantage of debiasing through orthogonal-
 653 ization is not that it makes each acquisition evaluation cheaper, but that it can achieve comparable
 654 acquisition stability with a smaller MC budget S .

655 **C.2 Variance reduction and ranking stability for GP and TPE**

656 We provide versions of Theorem 4.1 and Proposition 4.2 for the two surrogate instantiations GP
657 and TPE. The proofs follow directly from the proofs of the original theorem and proposition in
658 Supplement D.

659 **Theorem C.1** (Variance reduction for the GP surrogate). *Fix an iteration t and a candidate $\lambda \in \Lambda$.
660 Assume $\mathbb{E}_{q_t} \|g_t(\theta)\|_2^2 < \infty$, $\mathbb{E}_{q_t} [(\text{EI}^{\text{GP}}(\lambda; \theta))^2] < \infty$, and that $\text{Cov}(g_t, g_t)$ is nonsingular. Define*

$$\text{EI}^{\text{GP,orth}}(\lambda; \theta) := \text{EI}^{\text{GP}}(\lambda; \theta) - \gamma^{\text{GP}}(\lambda)^\top g_t(\theta). \quad (37)$$

661 *Then:*

662 1. *Target preservation:* $\mathbb{E}_{q_t} [\text{EI}^{\text{GP,orth}}(\lambda; \theta)] = \text{EI}^{\text{GP,marg}}(\lambda)$.

663 2. *Variance reduction:*

$$\text{Var}(\text{EI}^{\text{GP,orth}}(\lambda; \theta)) = \text{Var}(\text{EI}^{\text{GP}}(\lambda; \theta)) - \text{Cov}(\text{EI}^{\text{GP}}, g_t)^\top \text{Cov}(g_t, g_t)^{-1} \text{Cov}(g_t, \text{EI}^{\text{GP}}),$$

664 *and therefore* $\text{Var}(\text{EI}^{\text{GP,orth}}(\lambda; \theta)) \leq \text{Var}(\text{EI}^{\text{GP}}(\lambda; \theta))$.

665 *By the i.i.d. MC variance identity, the corresponding estimator satisfies* $\text{Var}(\widehat{\text{EI}}^{\text{GP,orth}}(\lambda)) \leq$
666 $\text{Var}(\widehat{\text{EI}}^{\text{GP,MC}}(\lambda))$ *for any MC budget S .*

667 **Proposition C.2** (Ranking stability for the GP surrogate). *Let $\lambda, \lambda' \in \Lambda$ with $\Delta^{\text{GP}}(\lambda, \lambda') :=$
668 $\text{EI}^{\text{GP,marg}}(\lambda) - \text{EI}^{\text{GP,marg}}(\lambda') > 0$. Let $\widehat{\Delta}^{\text{GP,orth}}$ denote the MC estimator obtained by applying the
669 orthogonalization construction in Theorem C.1 to the difference functional $\text{EI}^{\text{GP}}(\lambda; \theta) - \text{EI}^{\text{GP}}(\lambda'; \theta)$.
670 Then*

$$\mathbb{P}(\widehat{\Delta}^{\text{GP,orth}} \leq 0) \leq \frac{\text{Var}(\widehat{\Delta}^{\text{GP,orth}})}{(\Delta^{\text{GP}}(\lambda, \lambda'))^2} \leq \frac{\text{Var}(\widehat{\Delta}^{\text{GP,MC}})}{(\Delta^{\text{GP}}(\lambda, \lambda'))^2}. \quad (38)$$

671 **Theorem C.3** (Variance reduction for the tree-based / TPE-style surrogate). *Fix an iteration t and a
672 candidate $\lambda \in \Lambda$. Assume $\mathbb{E}_{q_t} [\|c_t(\theta)\|_2^2] < \infty$, $\mathbb{E}_{q_t} [(h^{\text{TPE}}(\lambda; \theta))^2] < \infty$, $\mathbb{E}_{q_t} [c_t(\theta)] = 0$, and that
673 $\text{Cov}(c_t, c_t)$ is nonsingular. Define*

$$h^{\text{TPE,orth}}(\lambda; \theta) := h^{\text{TPE}}(\lambda; \theta) - \gamma^{\text{TPE}}(\lambda)^\top c_t(\theta). \quad (39)$$

674 *Then:*

675 1. *Target preservation:* $\mathbb{E}_{q_t} [h^{\text{TPE,orth}}(\lambda; \theta)] = \text{EI}^{\text{TPE,marg}}(\lambda)$.

676 2. *Variance reduction:*

$$\text{Var}(h^{\text{TPE,orth}}(\lambda; \theta)) = \text{Var}(h^{\text{TPE}}(\lambda; \theta)) - \text{Cov}(h^{\text{TPE}}, c_t)^\top \text{Cov}(c_t, c_t)^{-1} \text{Cov}(c_t, h^{\text{TPE}}),$$

677 *and therefore* $\text{Var}(h^{\text{TPE,orth}}(\lambda; \theta)) \leq \text{Var}(h^{\text{TPE}}(\lambda; \theta))$.

678 *By the i.i.d. MC variance identity, the corresponding estimator satisfies* $\text{Var}(\widehat{\text{EI}}^{\text{TPE,orth}}(\lambda)) \leq$
679 $\text{Var}(\widehat{\text{EI}}^{\text{TPE,MC}}(\lambda))$ *for any MC budget S .*

680 **Proposition C.4** (Ranking stability for the tree-based / TPE-style surrogate). *Let $\lambda, \lambda' \in \Lambda$ with
681 $\Delta^{\text{TPE}}(\lambda, \lambda') := \text{EI}^{\text{TPE,marg}}(\lambda) - \text{EI}^{\text{TPE,marg}}(\lambda') > 0$. Let $\widehat{\Delta}^{\text{TPE,orth}}$ denote the MC estimator
682 obtained by applying the orthogonalization construction in Theorem C.3 to the difference functional
683 $h^{\text{TPE}}(\lambda; \theta) - h^{\text{TPE}}(\lambda'; \theta)$. Then*

$$\mathbb{P}(\widehat{\Delta}^{\text{TPE,orth}} \leq 0) \leq \frac{\text{Var}(\widehat{\Delta}^{\text{TPE,orth}})}{(\Delta^{\text{TPE}}(\lambda, \lambda'))^2} \leq \frac{\text{Var}(\widehat{\Delta}^{\text{TPE,MC}})}{(\Delta^{\text{TPE}}(\lambda, \lambda'))^2}. \quad (40)$$

684 **D Proofs of the main theorems and propositions**

685 **D.1 Proof of Theorem 4.1**

686 **Theorem 4.1** Assume $\mathbb{E}_{q_{m,t}}[\text{EI}_m(\lambda; \theta)^2] < \infty$ and $\mathbb{E}_{q_{m,t}}[\|g(\theta)\|_2^2] < \infty$, that Σ_g is nonsingular, and
 687 that the standard score-function regularity conditions hold for $q_{m,t}$. Then:

688 **(i) Target preservation.** $\text{EI}_m^{\text{orth,marg}}(\lambda) = \text{EI}_m^{\text{marg}}(\lambda)$.

689 **(ii) Variance reduction.** Under $q_{m,t}$,

$$\text{Var}(\text{EI}_m^{\text{orth}}(\lambda; \theta)) = \text{Var}(\text{EI}_m(\lambda; \theta)) - \text{Cov}(g, \text{EI}_m)^\top \Sigma_g^{-1} \text{Cov}(g, \text{EI}_m) \leq \text{Var}(\text{EI}_m(\lambda; \theta)). \quad (41)$$

690 **(iii) Local robustness in score-tilt directions.** For any $b \in \mathbb{R}^{d_\theta}$, consider the tilted family $q_{m,t}^{(\varepsilon)}(\theta) \propto$
 691 $q_{m,t}(\theta) \exp(\varepsilon b^\top g(\theta))$ for small $|\varepsilon|$. Holding $\gamma_m(\lambda)$ fixed at its $\varepsilon = 0$ value,

$$\left. \frac{d}{d\varepsilon} \mathbb{E}_{q_{m,t}^{(\varepsilon)}}[\text{EI}_m^{\text{orth}}(\lambda; \theta)] \right|_{\varepsilon=0} = 0. \quad (42)$$

692 *Proof.* Let $h(\theta) := \text{EI}_m(\lambda; \theta)$ and $g(\theta) := \nabla_\theta \log q_{m,t}(\theta)$, and write $\gamma := \gamma_m(\lambda) = \Sigma_g^{-1} \text{Cov}(g, h)$.
 693 Then

$$\text{EI}_m^{\text{orth}}(\lambda; \theta) = h(\theta) - \gamma^\top g(\theta). \quad (43)$$

694 **(i) Target preservation.** Under the assumed regularity conditions, the score function satisfies

$$\mathbb{E}_{q_{m,t}}[g(\theta)] = \int \nabla_\theta \log q_{m,t}(\theta) q_{m,t}(\theta) d\theta = \int \nabla_\theta q_{m,t}(\theta) d\theta = 0. \quad (44)$$

695 Hence

$$\text{EI}_m^{\text{orth,marg}}(\lambda) = \mathbb{E}_{q_{m,t}}[h(\theta)] - \gamma^\top \mathbb{E}_{q_{m,t}}[g(\theta)] = \mathbb{E}_{q_{m,t}}[h(\theta)] = \text{EI}_m^{\text{marg}}(\lambda). \quad (45)$$

696 **(ii) Variance reduction.** For any $a \in \mathbb{R}^{d_\theta}$,

$$\text{Var}_{q_{m,t}}(h - a^\top g) = \text{Var}(h) - 2a^\top \text{Cov}(g, h) + a^\top \Sigma_g a. \quad (46)$$

697 Since Σ_g is positive definite by assumption, this quadratic in a is uniquely minimized at

$$a^* = \Sigma_g^{-1} \text{Cov}(g, h) = \gamma. \quad (47)$$

698 Substituting $a = \gamma$ gives

$$\text{Var}(\text{EI}_m^{\text{orth}}(\lambda; \theta)) = \text{Var}(h) - \text{Cov}(g, h)^\top \Sigma_g^{-1} \text{Cov}(g, h). \quad (48)$$

699 Because Σ_g is positive definite, the quadratic form $\text{Cov}(g, h)^\top \Sigma_g^{-1} \text{Cov}(g, h) \geq 0$, and therefore

$$\text{Var}(\text{EI}_m^{\text{orth}}(\lambda; \theta)) \leq \text{Var}(\text{EI}_m(\lambda; \theta)). \quad (49)$$

700 **(iii) Local robustness.** Throughout this part we hold γ fixed at its $\varepsilon = 0$ value. Consider the
 701 exponentially tilted family

$$q_{m,t}^{(\varepsilon)}(\theta) = \frac{q_{m,t}(\theta) \exp(\varepsilon b^\top g(\theta))}{Z(\varepsilon)}, \quad Z(\varepsilon) = \int q_{m,t}(\theta) \exp(\varepsilon b^\top g(\theta)) d\theta. \quad (50)$$

702 For any integrable test function φ , differentiation under the integral sign (justified by the square-
 703 integrability assumptions and dominated convergence in a neighborhood of $\varepsilon = 0$) yields

$$\left. \frac{d}{d\varepsilon} \mathbb{E}_{q_{m,t}^{(\varepsilon)}}[\varphi(\theta)] \right|_{\varepsilon=0} = \text{Cov}_{q_{m,t}}(\varphi(\theta), b^\top g(\theta)). \quad (51)$$

704 Apply with $\varphi(\theta) = \text{EI}_m^{\text{orth}}(\lambda; \theta) = h(\theta) - \gamma^\top g(\theta)$:

$$\text{Cov}(g, h - \gamma^\top g) = \text{Cov}(g, h) - \Sigma_g \gamma = \text{Cov}(g, h) - \Sigma_g \Sigma_g^{-1} \text{Cov}(g, h) = 0. \quad (52)$$

705 Hence

$$\left. \frac{d}{d\varepsilon} \mathbb{E}_{q_{m,t}^{(\varepsilon)}}[\text{EI}_m^{\text{orth}}(\lambda; \theta)] \right|_{\varepsilon=0} = b^\top \text{Cov}(g, h - \gamma^\top g) = 0 \quad (53)$$

706 for every $b \in \mathbb{R}^{d_\theta}$. The orthogonalized acquisition is therefore first-order insensitive to score-tilt
 707 perturbations of $q_{m,t}$. \square

708 **D.2 Proof of Proposition 4.2**

709 **Proposition 4.2** Let $\lambda, \lambda' \in \Lambda$ such that $\Delta(\lambda, \lambda') := \text{EI}_m^{\text{marg}}(\lambda) - \text{EI}_m^{\text{marg}}(\lambda') > 0$. Let $\widehat{\Delta}_{\text{MC}}(\lambda, \lambda')$
 710 and $\widehat{\Delta}_{\text{orth}}(\lambda, \lambda')$ denote the corresponding Monte Carlo difference estimators. Then both estimators
 711 are unbiased for $\Delta(\lambda, \lambda')$, and

$$\mathbb{P}\left(\widehat{\Delta}_{\text{orth}}(\lambda, \lambda') \leq 0\right) \leq \frac{\text{Var}\left(\widehat{\Delta}_{\text{orth}}(\lambda, \lambda')\right)}{\text{Var}\left(\widehat{\Delta}_{\text{orth}}(\lambda, \lambda')\right) + \Delta(\lambda, \lambda')^2} \leq \frac{\text{Var}\left(\widehat{\Delta}_{\text{MC}}(\lambda, \lambda')\right)}{\text{Var}\left(\widehat{\Delta}_{\text{MC}}(\lambda, \lambda')\right) + \Delta(\lambda, \lambda')^2}. \quad (54)$$

712 *Proof.* By unbiasedness from Theorem 4.1,

$$\mathbb{E}\left[\widehat{\Delta}_{\text{orth}}(\lambda, \lambda')\right] = \Delta(\lambda, \lambda'). \quad (55)$$

713 Hence

$$\mathbb{P}\left(\widehat{\Delta}_{\text{orth}}(\lambda, \lambda') \leq 0\right) = \mathbb{P}\left(\widehat{\Delta}_{\text{orth}}(\lambda, \lambda') - \Delta(\lambda, \lambda') \leq -\Delta(\lambda, \lambda')\right). \quad (56)$$

714 Applying Cantelli's inequality to the centered random variable

$$Z := \widehat{\Delta}_{\text{orth}}(\lambda, \lambda') - \Delta(\lambda, \lambda') \quad (57)$$

715 gives

$$\mathbb{P}(Z \leq -a) \leq \frac{\text{Var}(Z)}{\text{Var}(Z) + a^2}, \quad a > 0. \quad (58)$$

716 With $a = \Delta(\lambda, \lambda')$, this yields

$$\mathbb{P}\left(\widehat{\Delta}_{\text{orth}}(\lambda, \lambda') \leq 0\right) \leq \frac{\text{Var}\left(\widehat{\Delta}_{\text{orth}}(\lambda, \lambda')\right)}{\text{Var}\left(\widehat{\Delta}_{\text{orth}}(\lambda, \lambda')\right) + \Delta(\lambda, \lambda')^2}. \quad (59)$$

717 The second inequality follows from Theorem 4.1, which gives

$$\text{Var}\left(\widehat{\Delta}_{\text{orth}}(\lambda, \lambda')\right) \leq \text{Var}\left(\widehat{\Delta}_{\text{MC}}(\lambda, \lambda')\right). \quad (60)$$

718 Since the map $v \mapsto v/(v + \Delta^2)$ is increasing for $v \geq 0$, the result follows. \square

719 **E Implementation details**

720 **General implementation details** All experiments were conducted on a proprietary compute cluster.
721 Our experiments utilized 4 CPU cores, 16 GB of main memory, and a A100 MIG slice with 7GB
722 of VRAM. We used BoTorch for all code implementations. In total, we use ~ 1000 GPU hours
723 including initial experimentation and bug fixing.

724 **License details** We used BoTorch, which is MIT licensed. We implemented all methods and
725 benchmarks in this framework. The WM811K dataset [33, 64] is available at <http://mirmlab.org/dataSet/public/> and allows *usage* under the following conditions (copied verbatim from the
726 license document):
727

728 2. Redistribution and use in any form must be
729 accompanied by the following two citations:

730
731 [1] Ming-Ju Wu, Jyh-Shing Roger Jang, and Jui-Long Chen,
732 "Wafer Map Failure Pattern Recognition and
733 Similarity Ranking for Large-Scale Data Sets,"
734 in IEEE Transactions on Semiconductor Manufacturing,
735 vol. 28, no. 1, pp. 1-12, Feb. 2015, doi: 10.1109/TSM.2014.2364237.

736
737 [2] MIR-WM811K: Dataset for wafer map failure pattern recognition,
738 2015 <http://mirmlab.org/dataset/public/>

739 We included these citations at the appropriate places. Further, we received written permission by the
740 authors to use the dataset for our research.

741 **E.1 TPE implementation**

742 **TPE implementation.** For the TPE-based experiments, we approximate the parameter distribution
743 $q_t(\theta)$ by a nonparametric bootstrap over the current BO history. Concretely, for each sample
744 $s = 1, \dots, S$, we resample the observed data with replacement, fit a TPE surrogate $\theta^{(s)}$, and evaluate
745 the corresponding TPE acquisition

$$h_{\text{TPE}}(\lambda; \theta^{(s)}) \propto \frac{\ell_{\theta^{(s)}}(\lambda)}{g_{\theta^{(s)}}(\lambda)}.$$

746 We use the empirical γ -quantile split with $\gamma = 0.2$ to define the good and bad sets.

747 As control variate, we use centered bootstrap summary statistics. For each bootstrap fit $\theta^{(s)}$, we form

$$\phi(\theta^{(s)}) = [\log b_{\text{good},1}^{(s)}, \dots, \log b_{\text{good},d}^{(s)}, \log b_{\text{bad},1}^{(s)}, \dots, \log b_{\text{bad},d}^{(s)}, y^{*,(s)}]^\top,$$

748 where $b_{\text{good},j}^{(s)}$ and $b_{\text{bad},j}^{(s)}$ are the fitted KDE bandwidths in dimension j , and $y^{*,(s)}$ is the bootstrap
749 split threshold. The implemented control variate is then

$$c_t(\theta^{(s)}) = \phi(\theta^{(s)}) - \frac{1}{S} \sum_{r=1}^S \phi(\theta^{(r)}),$$

750 which is zero-mean by construction. The orthogonalized TPE estimator is

$$\widehat{\text{El}}_{\text{TPE,orth}}^c(\lambda) = \frac{1}{S} \sum_{s=1}^S \left(h_{\text{TPE}}(\lambda; \theta^{(s)}) - \gamma_{\text{TPE}}(\lambda)^\top c_t(\theta^{(s)}) \right).$$

751 **E.2 Misspecified kernel family experiments**

752 **E.3 CNN outlier experiment**

753 We include the CNN outlier experiment used by Ament et al. [3]. Specifically, the goal is to optimize
754 the test-set accuracy of a CNN trained on MNIST images. To model corrupted evaluations, we

Function	Qualitative structure	Kernel family		Reason for misspecification
		Mildly misspecified	Strongly misspecified	
Hartmann6	Smooth, anisotropic, interaction-heavy, multimodal	Isotropic RBF or isotropic Matérn	Linear	Isotropic kernel variants miss dimension-specific lengthscales. Linear kernels cannot capture the nonlinear multimodal structure.
Ackley8	Highly multimodal, oscillatory, with many local minima	Smooth RBF	Linear	Very smooth kernels tend to smear out Ackley’s rugged local structure. Linear kernels are too simple for the landscape.
Michalewicz10	Sharp valleys, strong multimodality, narrow optima	Isotropic RBF or Matérn	Linear	Isotropic kernels blur the narrow valleys and local optima. Linear kernels are strongly misspecified.
Levy16	Rugged, oscillatory, with nonlinear interactions	Overly smooth isotropic RBF	Linear	Isotropic smooth kernels underfit sharper local variation and ignore anisotropy. Linear kernels cannot represent the nonlinear oscillatory structure.

Table 5: Synthetic benchmark functions, their qualitative structure, and examples of mildly and strongly misspecified kernel families.

755 prematurely stop training with probability 20 % after observing between 100 and 1000 training
756 samples. We optimize the same parameters as Ament et al. [3]. For comparison, we also include their
757 RPR method. To not implicitly mitigate the corrupted evaluations, we only use $n_0 = 2$ initial points.
758 In initial experiments, we found that using larger n_0 made the problem *too easy for all methods*; all
759 reached $\sim 100\%$ test accuracy after few iterations.

760 E.4 Case study experiment, additional details

761 **Dataset.** For our real-world case study, we use the large-scale industrial WM811K dataset [33, 64].
762 It contains 811 457 real-world wafer map images and annotations of common failure types. The
763 images have pixel values 0, 1, 2, where 0 is background, 1 are OK regions, and 2 are faulty regions.
764 Depending on the clustering of the failures, one of eight class labels is assigned. Exemplary failure
765 patterns are shown in Figure 6. Images without any failures are labelled with an additional class label.
766 We resize all wafer maps to a common (48×48) resolution, rescale the data to $(0, 1)$ by dividing by
767 $/2$, and then normalize using ImageNet statistics.

768 The dataset is strongly imbalanced, and we randomly split the data into a 80 % train and 20 % test set,
769 stratifying on the labels to keep class distributions equal across the two splits. During fine-tuning, we
770 use an (inverse) class-weighted cross-entropy loss.

771 **Model.** As a model, we use an ImageNet [13] pretrained ViT B16 [14] model. We freeze the
772 backbone model and train only the penultimate output layers.

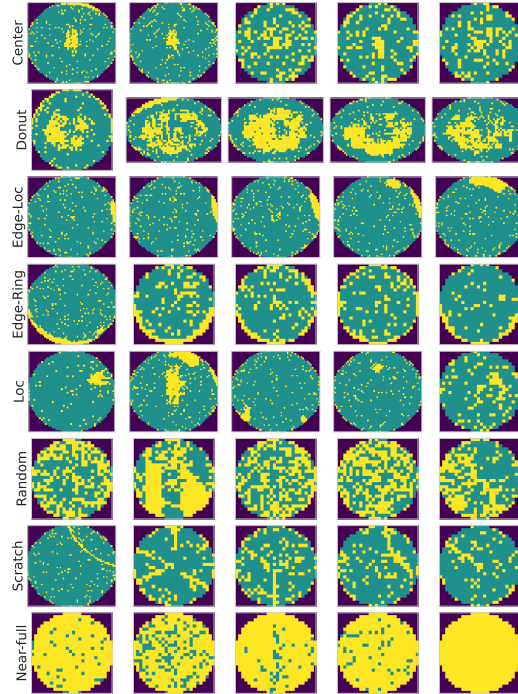


Figure 6: Examples of wafer-map manufacturing failures.

773 **Experimental settings.** For this experiment, due to computational constraints, we deviate from the
 774 settings used in the main paper. Specifically, we here use $S = 512$, $n_0 = 4$, and run four trials of 20
 775 iterations. We use a RBF kernel.

776 **Problem.** The goal is to optimize the fine-tuning test accuracy of a pretrained ViT B16 model on the
 777 industrial WM811K dataset. The tunable parameters are: learning rate in $[1e-4, 1e-1]$, momentum
 778 in $[0, 1]$, weight decay in $[0, 1]$, step size (for the learning rate scheduling) in $[1, 100]$, and $\gamma \in [0, 1]$.⁵

⁵This experiment conceptually builds on Ament et al.’s experiment on CNN training.

779 **F Further results**

780 **F.1 Variance reduction and ranking stability**

781 We estimate the empirical macro variance of the acquisition estimator by fixing a BO state, drawing
 782 a fixed set of $n = 64$ Sobol probe samples, rebuilding the acquisition independently 16 times, and
 783 computing the sample variance across repeated estimator values at each probe point. We report the
 784 mean pointwise macro variance for the raw and orthogonalized estimators. (Note: this does *not* affect
 785 the fitted surrogate model, which stays untouched). We repeat this at fixed points throughout training.
 786 Due to space reasons, we only report abridged results in the main text. We report results for (i) a
 787 Matérn-5/2 kernel in section F.1.1, (ii) a RBF kernel in section F.1.2, and (iii) a TPE surrogate in
 788 section F.1.3.

789 **F.1.1 Matérn-5/2 kernel**

790 In Table 6, we give the variance reduction of ORTHOBO for all benchmark functions. For higher-
 791 dimensional Michalewicz10 and Levy16, we found that larger n_0 were required. Across all bench-
 792 marks, we observe a strong reduction in variance, confirming that our proposed ORTHOBO stabilizes
 793 the acquisition estimate. Further, we report the ranking stabilities in Table 7. We observe that
 794 ORTHOBO again reduces the variance and consistently stabilizes the probe ranking.

Table 6: Variance reduction by orthogonalization.

n_0	S	Hartmann6		Ackley8		Michalewicz10		Levy16	
		$\text{Var}(\widehat{\text{EI}}_m^{\text{MC}})$	$\text{Var}(\widehat{\text{EI}}_m^{\text{ortho}})$	$\text{Var}(\widehat{\text{EI}}_m^{\text{MC}})$	$\text{Var}(\widehat{\text{EI}}_m^{\text{ortho}})$	$\text{Var}(\widehat{\text{EI}}_m^{\text{MC}})$	$\text{Var}(\widehat{\text{EI}}_m^{\text{ortho}})$	$\text{Var}(\widehat{\text{EI}}_m^{\text{MC}})$	$\text{Var}(\widehat{\text{EI}}_m^{\text{ortho}})$
32	8	2e-08	1.61e-07	9.68e-09	5.24e-09	9.51e-08	4.54e-08	0.00106	0.000709
32	32	2.59e-09	3.02e-10	6.3e-09	4.17e-10	1.67e-08	1.07e-09	0.000208	2.27e-05

Table 7: Ranking stability. Lower probe variance and flip rate indicate more stable acquisition estimates; higher ranking correlation and top-1 agreement indicate more stable candidate rankings across probes.

Benchmark	Method	Probe variance	Top1 agr.	Flip rate	Regret
Hartmann6	qLogEI	8.416 _{15.43}	0.925 _{0.16}	0.114 _{0.12}	1.411 _{0.42}
	OrthoBO (ours)	0.038 _{0.04}	0.950 _{0.10}	0.051 _{0.04}	1.232 _{0.67}
Ackley8	qLogEI	19.471 _{39.29}	0.958 _{0.11}	0.102 _{0.12}	17.864 _{2.08}
	OrthoBO (ours)	0.017 _{0.04}	0.986 _{0.06}	0.021 _{0.05}	18.136 _{0.78}
Michalewicz10	qLogEI	168.603 _{154.76}	0.925 _{0.16}	0.148 _{0.12}	7.140 _{0.01}
	OrthoBO (ours)	0.019 _{0.03}	0.988 _{0.06}	0.014 _{0.03}	6.961 _{0.36}
Levy16	qLogEI	233.655 _{287.67}	0.917 _{0.12}	0.069 _{0.09}	82.194 _{6.05}
	OrthoBO (ours)	0.073 _{0.21}	0.988 _{0.06}	0.048 _{0.04}	58.938 _{0.00}

795 **F.1.2 RBF kernel**

796 We repeat our previous analysis, but use RBF kernel for the GP surrogate. We report the raw and
 797 orthogonalized variances in Table 8. Our observations are consistent: orthogonalization substantially
 798 reduces estimation variance.

Table 8: Variance reduction through orthogonalization, using a RBF kernel.

n_0	S	Hartmann6		Ackley8		Michalewicz10		Levy16	
		$\text{Var}(\widehat{\text{EI}}_m^{\text{MC}})$	$\text{Var}(\widehat{\text{EI}}_m^{\text{ortho}})$	$\text{Var}(\widehat{\text{EI}}_m^{\text{MC}})$	$\text{Var}(\widehat{\text{EI}}_m^{\text{ortho}})$	$\text{Var}(\widehat{\text{EI}}_m^{\text{MC}})$	$\text{Var}(\widehat{\text{EI}}_m^{\text{ortho}})$	$\text{Var}(\widehat{\text{EI}}_m^{\text{MC}})$	$\text{Var}(\widehat{\text{EI}}_m^{\text{ortho}})$
32	8	5.91e-08	1.91e-07	1.28e-08	7.19e-09	1.15e-07	4.56e-08	0.00129	0.000758
32	32	8.62e-09	7.3e-10	3.71e-09	4e-10	2.61e-08	1.76e-09	0.000301	2.85e-05

Table 9: Ranking stability, **RBF kernel**. Lower probe variance and flip rate indicate more stable acquisition estimates; higher ranking correlation and top-1 agreement indicate more stable candidate rankings across probes.

Benchmark	Method	Probe variance	Top1 agr.	Flip rate	Regret
Hartmann6	qLogEI	5.373 _{11.64}	0.944 _{0.15}	0.085 _{0.10}	0.126 _{0.00}
	OrthoBO (ours)	0.031 _{0.04}	0.986 _{0.08}	0.031 _{0.04}	0.125 _{0.00}
Ackley8	qLogEI	10.051 _{31.34}	0.986 _{0.06}	0.052 _{0.10}	12.476 _{2.06}
	OrthoBO (ours)	0.013 _{0.04}	0.993 _{0.04}	0.019 _{0.04}	12.581 _{4.03}
Michalewicz10	qLogEI	55.138 _{100.54}	0.931 _{0.15}	0.085 _{0.09}	4.517 _{0.20}
	OrthoBO (ours)	0.002 _{0.01}	1.000 _{0.00}	0.031 _{0.04}	4.004 _{0.52}
Levy16	qLogEI	82.913 _{154.71}	0.889 _{0.18}	0.201 _{0.10}	3.109 _{0.57}
	OrthoBO (ours)	0.034 _{0.17}	1.000 _{0.00}	0.049 _{0.04}	2.819 _{0.57}

799 F.1.3 TPE-based surrogate

800 In Table 10, we show the variance reduction from ORTHOBO for a TPE surrogate, using a MC budget
 801 of $S = 32$. Our observations are in line with the GP surrogate: ORTHOBO can substantially reduce
 802 the variance.

Table 10: Variance reduction by orthogonalization, for a TPE surrogate.

Hartmann6		Ackley8		Michalewicz10		Levy16	
$\text{Var}(\text{EI}^{\text{TPE}})$	$\text{Var}(\widehat{\text{EI}}^{\text{TPE,orth}})$	$\text{Var}(\text{EI}^{\text{TPE}})$	$\text{Var}(\widehat{\text{EI}}^{\text{TPE,orth}})$	$\text{Var}(\text{EI}^{\text{TPE}})$	$\text{Var}(\widehat{\text{EI}}^{\text{TPE,orth}})$	$\text{Var}(\text{EI}^{\text{TPE}})$	$\text{Var}(\widehat{\text{EI}}^{\text{TPE,orth}})$
11.222 _{16.405}	1.214 _{2.620}	13.182 _{18.233}	1.284 _{2.601}	21.712 _{25.982}	2.197 _{3.358}	20.944 _{26.200}	1.399 _{2.259}

803 F.2 Misspecified surrogate families

804 We first study robustness to surrogate misspecification. In practical BO, the surrogate family is
 805 only an approximation to the unknown objective, and kernel hyperparameters are often estimated
 806 from limited data. As a result, the posterior quantities used inside MC acquisition estimates can
 807 be noisy or biased, which may destabilize candidate rankings. This experiment evaluates whether
 808 the variance reduction induced by ORTHOBO improves BO decisions under such challenging, but
 809 realistic, surrogate conditions.

810 We consider four standard synthetic regression problems with known global optima and varying
 811 dimensionality: Hartmann6, Ackley8, Michalewicz10, and Levy16 [e.g., 2, 10, 15, 41]. Details of
 812 the benchmark functions and surrogate choices are provided in Table 5. We compare methods using
 813 best-so-far regret.

814 **Mild misspecification: isotropic RBF kernel.** We first use an isotropic RBF kernel, a common
 815 smooth GP prior. For anisotropic and multimodal objectives it is mildly misspecified because it uses
 816 a single lengthscale across all dimensions. (The dimension of the benchmark functions varies from 6
 817 to 16.) The results are in Fig. 7.

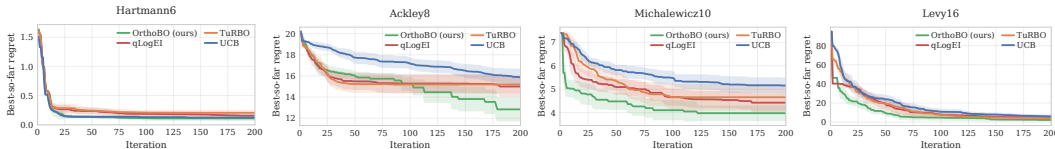


Figure 7: Best-so-far regret as a function of iterations on four widely used regression problems, using a **mildly misspecified isotropic RBF kernel**. Our proposed ORTHOBO is competitive, and in higher-dimensional functions (Michalewicz10, Levy16) outperforms the baselines.

818 **Results.** Under mild surrogate misspecification, our proposed ORTHOBO remains competitive across
 819 all benchmarks. On Hartmann6 and Ackley8, it improves over qLogEI and reaches the lowest regret.
 820 On the higher-dimensional Michalewicz10 and Levy16 functions, ORTHOBO achieves the lowest

821 regret and improves substantially earlier than the baselines. For example, ORTHOBO reaches a regret
 822 of ~ 5 after 5 iterations, whereas the baseline methods require at least 50 to 60 *additional iterations*.
 823 Together, these results indicate that variance reduction at the acquisition-estimation level improves
 824 the reliability of BO decisions under realistic surrogate mismatch, with the largest gains appearing in
 825 higher-dimensional settings where candidate rankings are more sensitive to estimation noise.

826 **Strong misspecification: linear kernel.** We next consider a more challenging setting by using a
 827 linear kernel. This surrogate is unable to represent the nonlinear and multimodal structure of the
 828 benchmark functions, which amplifies the effect of posterior and acquisition-estimation errors. We
 829 therefore test whether stabilizing the acquisition estimate can still improve BO decisions when the
 830 surrogate family is substantially imperfect. The results are shown in Fig. 8.

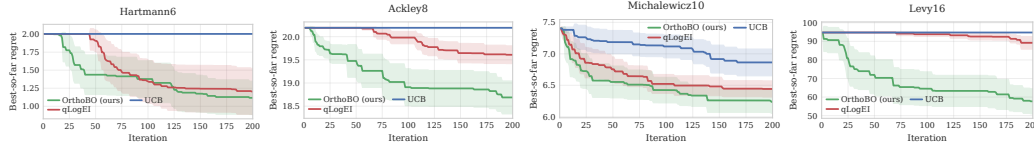


Figure 8: Best-so-far regret as a function of iterations on four widely used regression problems, using a **strongly misspecified linear kernel**. Our proposed ORTHOBO substantially outperforms other methods. The performance difference increases as the problem dimensionality grows. Notably, on Hartmann6, qLogEI takes five times more iterations than ORTHOBO to reach the same level of regret.

831 **Results.** We observe: (i) On all regression problems, our ORTHOBO outperforms the baselines. (ii)
 832 on Hartmann6, ORTHOBO reaches the same regret as qLogEI in half the time. (iii) On Levy16,
 833 ORTHOBO is the only method to substantially improve over time. Together, the results demonstrate
 834 the orthogonalization can successfully overcome strong misspecifications of the surrogate family.

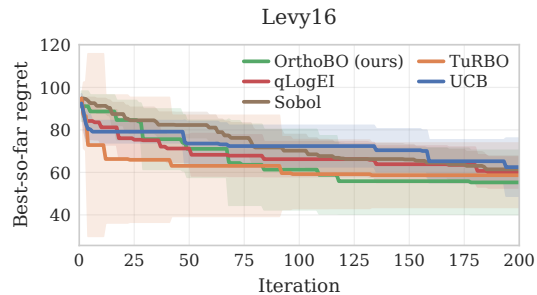
835 With the strongly misspecified linear kernel, ORTHOBO outperforms the baselines on all four
 836 functions. On Hartmann6, it reaches the same regret level as qLogEI with substantially fewer
 837 evaluations. On Levy16, ORTHOBO is the only method that shows a clear improvement over
 838 the optimization horizon. Together, these results indicate that acquisition-estimation noise can be
 839 a significant failure mode when the surrogate is imperfect, and that orthogonalization helps by
 840 stabilizing the acquisition values used to rank candidates. Thus, ORTHOBO does not require a
 841 perfectly specified surrogate to improve BO performance; rather, it provides a variance-reduction
 842 mechanism that is especially useful when MC acquisition estimates become unreliable.

843 F.3 Weakly fitted hyperparameters

844 We next study a data-scarce regime in which the surrogate family is expressive, but its hyperparameters
 845 are weakly identified. This situation is common in expensive HPO problems: only a small initial
 846 design is available, yet the GP must already estimate lengthscales, output scale, and noise parameters
 847 that determine the posterior and hence the acquisition values. Uncertainty or instability in these
 848 hyperparameters can therefore propagate to the acquisition function and distort candidate rankings.

849 To isolate this effect, we use a GP with a Matérn-5/2 kernel and ARD. Unlike the isotropic
 850 kernels studied above, ARD assigns a separate lengthscale to each input dimension. This
 851 makes the surrogate more flexible, but also harder to fit when the initial design is small.
 852 We therefore initialize BO with only $n_0 = 4$ Sobol points and evaluate performance on
 853 Levy16, which is particularly challenging in this setting because the surrogate must estimate
 854 16 lengthscales and a noise parameter from very limited data. The results are shown in Fig. 9.
 855

856 **Results.** All methods initially exhibit high
 857 regret, reflecting the difficulty of fitting an
 858 ARD surrogate from only four initial observa-
 859 tions. Nevertheless, ORTHOBO decreases re-
 860 gret steadily over the optimization horizon and
 861 reaches one of the lowest final regrets among
 862 the compared methods. These results show that
 863 orthogonalization remains useful even when the
 864 kernel class is flexible but its hyperparameters



27 Figure 9: Best-so-far regret on Levy16 with a Matérn-5/2 ARD kernel and only $n_0 = 4$ initial Sobol points. ORTHOBO remains among the strongest methods and improves steadily despite weakly identified surrogate hyperparameters.

865 are poorly estimated. The gains are consistent
 866 with the proposed mechanism: reducing sensi-
 867 tivity of the acquisition estimate to surrogate-
 868 posterior perturbations yields more reliable candi-
 869 date rankings in early, data-limited BO.

870 F.4 Ensembling

871 In the main experiments, we used $M = 1$ to
 872 avoid confounding the effect of orthogonaliza-
 873 tion with mere performance gains from a larger ensemble. We now use an ensemble of $M = 3$
 874 models, each with a different kernel: linear, RBF, and Matérn-5/2. We increased the number of
 875 iterations by 25% to account for additional fitting iterations for the models. We compute $\ell_{m,t}$ from
 876 the GP surrogate posteriors, set $\tau = 1$, and use $\delta = 0.001$ as a minimum weight floor.

877 The results for ensembling are in Figure 10.

878 Further, we compute (i) the variance reduction, (ii) ranking stability metrics, and (iii) weighting
 879 metrics. These additional measures are computed on a Sobol probe. We present the results in Table 11
 880 and Figure 11.

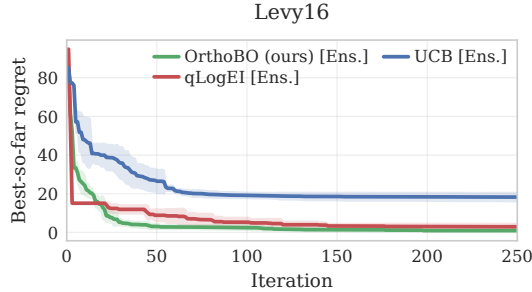


Figure 10: Results for ensembling.

Table 11: Variance reduction for ensembling on Levy16.

Method	Probe variance	Top1 agr.	Flip rate	Regret	Entropy	Weight 1	Weight 2	Weight 3
UCB [Ens.]	39.255 _{48.64}	1.000 _{0.00}	0.026 _{0.04}	18.267 _{2.40}	0.163 _{0.18}	0.008 _{0.06}	0.062 _{0.14}	0.933 _{0.14}
qLogEI [Ens.]	390.598 _{1033.83}	0.875 _{0.20}	0.167 _{0.10}	2.926 _{1.66}	0.145 _{0.16}	0.014 _{0.07}	0.073 _{0.21}	0.917 _{0.21}
ORTHOBO (Ours) [Ens.]	0.039 _{0.12}	0.958 _{0.14}	0.081 _{0.08}	0.941 _{0.54}	0.169 _{0.23}	0.008 _{0.06}	0.069 _{0.14}	0.927 _{0.14}

881 **Results.** We observe: (1) ORTHOBO achieves the lowest regret. (ii) ORTHOBO has substantially
 882 lower Sobol probing variance. (iii) Weight are highest for the Matérn-5/2 kernel (“Weight 3”)
 883 demonstrating that all methods primarily focus the “most correct” kernel. We plot the evolution of the
 884 entropy over time in Figure 11, which shows that the weight aggregation scheme balances exploration
 885 and exploitation.

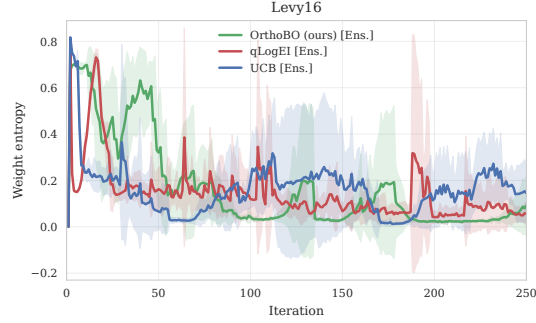


Figure 11: Evolution of the entropy across ensembling weights. For all methods, the used tempered exponentially weighted aggregation scheme is effective in balancing exploration and exploitation at the *ensemble* level.

886 F.5 Monte Carlo sampling ablation experiment

887 We here study the effect of the MC budget used to estimate the acquisition function. Here, finite-
 888 sample acquisition noise can change candidate rankings, with the effect becoming more pronounced
 889 for small sampling budgets S [cf. 2]. This setting is practically relevant because acquisition functions
 890 are evaluated many times during optimization, so increasing S can be computationally expensive.
 891 We therefore test whether the variance reduction of ORTHOBO improves BO performance when
 892 acquisition estimates must be computed from limited MC samples.

893 We use an isotropic RBF kernel, thereby combining moderate surrogate mismatch with reduced
 894 acquisition-estimation budgets. We vary $S \in \{64, 128, 256, 512\}$. The main results for Ackley8 and
 895 Michalewicz10 are shown in Fig. 3; additional results are provided in Fig. 12.

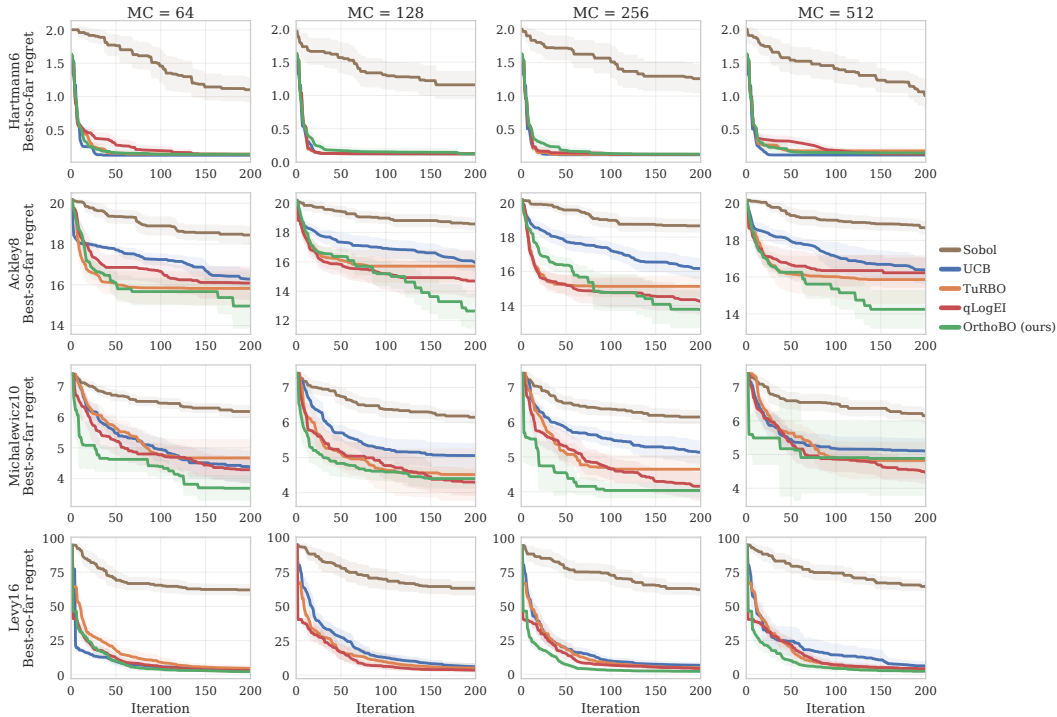


Figure 12: Different mc samples

896 **Results** Reducing the MC budget makes the performance of acquisition-based methods more sensitive
 897 to estimation noise. Across the tested budgets, ORTHOBO remains competitive on Ackley8 and
 898 consistently achieves low regret on Michalewicz10. The effect of orthogonalization are most visible

899 at earlier iterations, where MC estimation noise has a larger effect on the acquisition values and
900 hence on the induced candidate rankings. Overall, these results support the central mechanism of
901 ORTHOBO: orthogonalization improves the reliability of acquisition estimates when MC budgets are
902 limited, leading to more stable BO decisions and thus lower regret.