
BOLT: A Benchmark to Democratize Black-box Optimization Research for Expensive LLM Tasks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Optimization of LLM training and inference configurations, such as hyperpa-
2 rameters, data mixtures, and prompts, is critical to performance, but it is often
3 approached heuristically in practice, leading to potentially suboptimal outcomes.
4 By framing them as *noisy*, *expensive*, and *derivative-free* optimization problems,
5 Bayesian optimization (BO) and other black-box optimization (BBO) methods
6 offer a promising yet underexplored direction for principled, sample-efficient meth-
7 ods. However, LLM training and inference costs are prohibitively high for most
8 of the BBO research community, and new methods are often only evaluated on
9 synthetic test functions and small-scale datasets that fail to capture the challenges
10 of modern LLM optimization problems. This impedes the development of BBO
11 methods and makes it difficult to assess their effectiveness on modern LLM tasks.
12 We introduce BOLT, the first LLM-centric benchmark that democratizes LLM
13 research for the BBO community. BOLT covers broad and well-motivated LLM
14 optimization problems, involving multi-fidelity, multi-objective, heteroscedastic
15 noise, and high-dimensional search spaces. Each problem in BOLT is grounded
16 in real experimental data and made fully reproducible and accessible through
17 lightweight surrogate models fitted to the results of thousands of real LLM experi-
18 ments. We benchmark BOLT against an extensive range of BO and BBO methods,
19 showing that selected BO methods consistently outperform others across tasks and
20 highlighting gaps in existing BBO methods on LLM tasks, underscoring the need
21 to modernize benchmarks for the BBO community. An anonymized version of
22 BOLT is released at <https://github.com/anonom799/bolt>.

23 1 Introduction

24 Optimizing hyperparameters, data mixtures, and prompts for LLMs are critical decisions in any LLM
25 pipeline and often have substantial computational cost. By framing them as black-box optimization
26 (BBO) problems in the *noisy*, *expensive*, and *derivative-free* settings, Bayesian optimization (BO) and
27 other BBO methods emerge as promising solutions to tackle these problems, capable of handling noisy
28 observations and often requiring no gradients [24, 26]. Compared to heuristic approaches that are
29 often inefficient and suboptimal [2], BO and BBO are principled and sample-efficient approaches for
30 optimizing LLM configurations. However, researchers in the BBO community often lack sufficient
31 computational resources to validate new methods (and ideas) on LLM-centric problems. As a
32 result, many research endeavors in BBO are only validated by numerical or classical optimization
33 benchmarks (e.g., BBOB [28], HPO-B [4]), making it difficult to reliably assess their effectiveness
34 on modern LLM tasks. This impedes the development of BBO methods and leaves their potential
35 untapped for LLM tasks.

36 We aim to democratize black-box optimization research for modern LLM problems by releasing
37 BOLT, a **Benchmark for optimization of expensive LLM Tasks**, the first black-box optimization

38 benchmark grounded in real LLM objectives. The key goal of BOLT is to modernize the research
39 efforts of the BBO community by allowing researchers to validate their methods on LLMs without
40 requiring access to expensive computational resources. BOLT comprises three families of LLM
41 optimization tasks: hyperparameter optimization (HPO), data mixture optimization (DMO), and
42 prompt optimization (PO), all of which are active areas of research in LLMs [13, 70, 75]. Each task
43 family consists of multiple problem instances, framed as derivative-free optimization problems that
44 span key challenges in black-box optimization research, including multi-fidelity, multi-objective,
45 heteroscedastic noise, and high-dimensional settings. We hope our benchmark serves as a modern,
46 practical, and accessible validation process for future BBO research ideas.

47 BOLT lowers the barrier for BBO community to adopt LLM tasks as standard, reproducible bench-
48 marks via emulators and tabular datasets of precomputed evaluations. In particular, BOLT consist of
49 a suite of emulators, which are surrogate models trained on over 20k real LLM evaluations across
50 all tasks, and require little compute and time to query. Each emulator is validated to ensure that
51 conclusions drawn are consistent with the true objectives.

52 Our experiments show that BO methods consistently outperform standard HPO and evolutionary
53 baselines across BOLT’s tasks, while exposing structural and practical challenges that emerge under
54 our LLM-based benchmarks. These results suggest that BO is a promising but underexplored
55 direction for LLM optimization, and further underscores the importance of principled, accessible,
56 and reproducible benchmarking to realize this potential.

57 In summary, our contributions are:

- 58 • **A suite of LLM tasks as BBO problems.** We introduce three task families grounded in real
59 LLM objectives spanning hyperparameter (Sec. 4.3), data mixture (Sec. 4.4), and prompt
60 optimization (Sec. 4.5). Optimization problems within each task family are designed to cover
61 key challenges in black-box optimization research including multi-fidelity optimization,
62 multi-objective trade-offs, and high-dimensional search spaces.
- 63 • **Emulators to democratize reproducible evaluation on LLM tasks.** Emulators trained on
64 extensive real LLM experiments replace expensive LLM training and inference with cheap
65 surrogate queries, providing access to validated LLM-grounded objectives and enabling
66 reproducible comparisons without large-scale compute (Sec. 4.2).
- 67 • **Extensive empirical evaluation of BO and black-box baselines.** We benchmark a range
68 of BO methods alongside HPO and evolutionary baselines, and provide insights into method
69 selection and the structural and practical challenges of BO-based LLM optimization, high-
70 lighting their potential for LLM tasks (Sec. 5).

71 BOLT will be released as a Python library. The anonymized version of our code is available at
72 <https://github.com/anonom799/bolt>, and emulator models and tabular data are available on
73 the Huggingface Hub at <https://huggingface.co/collections/anonom799/bolt-models>.
74 Sec. A provides usage examples.

75 2 Related Works

76 **Black-box optimization benchmarks.** BBO algorithms are most commonly evaluated on syn-
77 thetic functions, such as Ackley, Branin-Hoo, and Hartmann-6 [9], with COCO [29] providing a
78 broader suite including 24 noiseless BBOB functions [28] and noisy, large-scale, and bi-objective
79 variants. While synthetic functions are convenient benchmarks as they are inexpensive to evaluate,
80 analytically tractable, with known global optima, they may not adequately represent real-world BBO
81 challenges [58].

82 Several real-world BBO benchmarks exist, each targeting specific problem settings. Liang et al. [49]
83 introduced five datasets in the material science domain and Kumar et al. [42] introduced a suite of
84 constrained optimization problems drawn from real-world applications, though closed-form equations
85 are used rather than simulation or experimental results. EXPObench [10] provides four expensive
86 real-world problems, such as wind farm layout optimization, with mixed and conditional variables.
87 Dreczkowski et al. [17] addresses mixed-variable and combinatorial BO including RNA inverse
88 folding. Our benchmark complements these efforts with a focus on optimization problems arising in
89 LLM development, a setting not addressed by existing BBO benchmarks.

Table 1: Summary of BOLT optimization problems. **MF**: multi-fidelity. **MO**: multi-objective.

Family	Problem	Dim	Optimization challenges
Hyperparameters	HPO	7	mixed variables
	HPO-MF-Cont	8	mixed variables, MF (continuous)
	HPO-MF-Disc	8	mixed variables, MF (discrete)
Data mixtures	DMO	6	simplex constraint
	DMO-MO	6	simplex constraint, MO
	DMO-Het	6	simplex constraint, heteroscedastic noise
Prompts	PO-128	128	high-dim.
	PO-256	256	high-dim.
	PO-512	512	high-dim.
	PO-768	768	high-dim.

90 **Optimization benchmarks for deep learning tasks.** Closer to our setting are benchmark suites
 91 targeting HPO in machine learning. HPOBench [18] and YAHPO Gym [61] provide surrogate-
 92 and tabular-based HPO problems over classical ML models, with support for multi-fidelity and
 93 multi-objective settings. HPO-B [4] is another large-scale black-box HPO benchmark derived from
 94 OpenML. These benchmarks have driven meaningful progress in HPO research. However, these
 95 benchmarks target classical or moderately-sized ML models and do not capture the optimization
 96 challenges that arise in LLM development, a gap our benchmark is designed to address.

97 **BBO for LLM.** Concurrently, a growing body of work has framed LLM optimization problems as
 98 BBO. In prompt optimization, Zhou et al. [81] and Yang et al. [77] first formulated it as a black-box
 99 problem and used the LLM itself as the optimizer. InstructZero [13] optimized prompts for black-box
 100 LLMs via BO with Gaussian process (GP) surrogates and instruction-coupled kernels. ZOPO [34]
 101 proposed NTK-based GPs targeting local optima, INSTINCT [51] used a neural surrogate over
 102 embeddings, and EASE [74] jointly optimized instructions and exemplars via neural bandits.

103 HPO for LLM fine-tuning presents another natural application for BO, as evaluation cost makes
 104 sample efficiency crucial. BO is well-established as an effective approach to HPO in machine learning
 105 [64, 66], and recent work extends this to LLM settings. Jang et al. [37] applied multi-objective BO
 106 to jointly optimize hyperparameters and fusion weights, while Seong-Eun et al. [63] injected LoRA
 107 domain knowledge into the search space to achieve strong configurations in as few as 30 iterations.

108 Beyond prompting and HPO, BO has been applied to data mixture optimization. Yen et al. [78]
 109 proposed a multi-fidelity framework that models uncertainty across mixtures and scales, while Chen
 110 et al. [14] used BO to re-weight post-training mixtures, achieving performance gains. These works
 111 indicate that BO is effective for a range of LLM optimization problems, with its sample efficiency
 112 and principled uncertainty quantification. Our benchmark complements these efforts by enabling
 113 accessible and reproducible comparison of BBO methods across these emerging problems.

114 3 Preliminaries

115 Black-box optimization seeks to optimize a black-box objective function $f : \mathcal{X} \rightarrow \mathbb{R}^d$, where d is
 116 the output dimension, $d = 1$ for single-objective problems and $d > 1$ if there are multiple objectives.
 117 Observed values are often noisy: $y = f(x) + \varepsilon$, where $\varepsilon \sim \mathcal{N}(0, \sigma_N^2)$. For instance, finding the
 118 hyperparameter that produces the best LLM performance is one such black-box optimization problem.

119 4 BOLT

120 BOLT provides a suite of LLM problems for reproducible and accessible benchmarking of BO
 121 and other blackbox optimization methods, via emulators and tabular lookups backed by real LLM
 122 experiments. An overview of task families and optimization problems is in Table 1. There are 3 task
 123 families: **hyperparameter optimization** (HPO), **data mixture optimization** (DMO), and **prompt**
 124 **optimization** (PO). We provide detailed descriptions, problem formulation, and problem variants of
 125 HPO in Sec. 4.3, DMO in Sec. 4.4, and PO in Sec. 4.5.

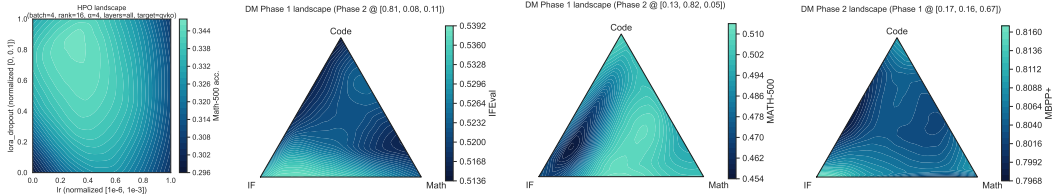


Figure 1: 2D slices of emulator landscapes

Table 2: Validation of objective function emulators using Spearman’s rank correlation ρ_s between emulator predictions and true evaluations. HF: high fidelity. LF: low fidelity.

Emulator	Problems	Objective	$\rho_{\text{train}} / \rho_{\text{test}}$	$n_{\text{train}} / n_{\text{test}}$
HPO-8B	HPO, HPO-MF-Cont, HPO-MF-Disc (HF)	MATH-500	0.967 / 0.938	7404 / 785
HPO-4B	HPO-MF-Disc (LF)		0.916 / 0.905	7091 / 819
DMO-4B	DMO, DMO-MO, DMO-Het	IFEval	0.759 / 0.721	3014 / 160
		MATH-500	0.780 / 0.779	
		MBPP+	0.758 / 0.721	
DMO-Noise	DMO-Het	MATH-500 σ	0.945 / -	50 / -

126 4.1 Design Principles

127 BOLT is designed around three principles that distinguish it from existing optimization benchmarks.

128 **Grounded in real LLM development workflows.** Each task family in BOLT corresponds to critical
 129 decisions in the modern LLM pipeline. Objective landscapes are obtained via real LLM experiments,
 130 ensuring that algorithm comparisons on BOLT will be meaningful for practical LLM research.

131 **Structured coverage of BO challenges.** BOLT provides controlled variation of challenges most
 132 relevant to modern BO research: mixed continuous-discrete-categorical search spaces, multi-fidelity
 133 observations with both continuous and discrete fidelity parameters, multi-objective trade-offs, het-
 134 eroscedastic noise, and high dimensions.

135 **Reproducibility through emulators.** Running BBO experiments with LLM training or inference in
 136 the loop is prohibitively expensive and often irreproducible due to under-specified configurations,
 137 making it hard to reliably compare methods. BOLT addresses both issues through deterministic
 138 *emulators*, which are surrogate models fitted on real experimental runs to reproduce the objective
 139 landscape at negligible query cost. Every blackbox function in BOLT is backed by an emulator or
 140 tabular data, so BBO methods can be evaluated without access to LLM infrastructure.

141 4.2 Emulators

142 The key purpose of our emulators is to *democratize* BBO research on modern LLM problems.
 143 Currently, validating BBO methods on LLM tasks requires large-scale compute that is inaccessible to
 144 most researchers. Our emulators remove this compute barrier, enabling BBO researchers to engage
 145 with modern LLM problems against validated surrogate objective landscapes at negligible cost. Using
 146 emulators in place of expensive oracles is also well-established in optimization benchmarking [30];
 147 HPOBench and YAHPO Gym both offer surrogate models, and YAHPO Gym shows that surrogate-
 148 based benchmarks match real functions more closely than tabular-based ones.

149 We fit 2-layer MLPs for the optimization objective of HPO and DMO tasks, which map a given LLM
 150 training configuration (e.g., data mixture) to the observed performance. Data for the emulators were
 151 generated via random sampling and adaptive sampling strategies: one maximizing emulator gradient
 152 magnitudes (to focus samples in high-variation regions) and one maximizing predicted objectives
 153 (to ensure the emulator is anchored by observed high-performing configurations). Random samples
 154 provide broad input space coverage, while adaptive strategies concentrate effort where it matters.

155 Emulators are validated on a Sobol-sampled held-out test set to ensure uniform space-filling coverage.
 156 We report Spearman’s rank correlation ρ as the validation metric, since candidate ranking matters

157 more than absolute calibration in BO. Train and test ρ are close across all emulators (Table 2),
158 indicating no overfitting. The more complex landscape of the DMO emulator in Fig. 1 and lower ρ
159 compared to HPO, suggests that the DMO objective is intrinsically harder to optimize than HPO’s.

160 For the heteroscedastic task in DMO, the noise model is a non-parametric kernel regression model
161 fitted on 50 configurations, each with score standard deviation estimated from 5 repeated evaluations,
162 with k -fold cross-validation. Additional details on emulator training, further analysis of emulators,
163 including for multi-fidelity settings and noise emulators, are in Sec. B.

164 4.3 Hyperparameter optimization (HPO)

165 The choice of hyperparameters (e.g., learning rate, batch size) plays a significant role in determining
166 the effectiveness of LLM training [46]. Our HPO benchmark consists of a moderate-dimensional
167 search space that describes how the LLM performance varies w.r.t. the choice of LoRA fine-tuning
168 hyperparameters (rank, alpha, target modules), learning rate, and batch size.

169 Our benchmark data in BOLT consists of datapoints from real LLM training runs with different
170 training hyperparameters. Each training run fine-tunes a model on a fixed 100k-sample subset of
171 OpenMathInstruct-2 [69] using LoRA under a given hyperparameter configuration. We then evaluated
172 the fine-tuned model on MATH-500 [31, 50] using the `lm-eval` harness framework [25]. We have
173 multiple observations for each training run via checkpoint evaluations throughout training, and used
174 early stopping to increase coverage of the hyperparameter space within a fixed compute budget.
175 In total, this yielded 7k+ observations per emulator across all runs and checkpoints. We trained 2
176 emulators, one for each model size (Qwen3-4B and Qwen3-8B [76]). More details are in Sec. C.1.

177 **Problem formulation.** The base problem search space is a 7 dimensional *mixed*-variable space,
178 comprising continuous variables (learning rate, LoRA dropout), integer variables (batch size, LoRA
179 rank, LoRA alpha, LoRA layers), and categorical variables (LoRA target modules). Full search space
180 specification in Table 4. After fitting, the emulator is a 2-layer MLP and is deterministic.

181 **Problem variants.** There are three variants in this task family. *HPO* base problem is single-fidelity,
182 single-objective at the full training budget, with a 7-dimensional search space. The other two are
183 multi-fidelity variants with an additional fidelity parameter to indicate the observed LLM performance
184 at a given scale (i.e., training tokens or model size), making them well-suited for multi-fidelity BO
185 research [23]. *HPO with continuous token fidelity* (HPO-MF-Cont) adds a fidelity parameter $s \in [0, 1]$,
186 normalized from the raw range $[10^5, 9.1 \times 10^6]$ that controls the number of training tokens seen.
187 *HPO with discrete model-size fidelity* (HPO-MF-Disc) uses a 4B low-fidelity and 8B high-fidelity
188 model. As a result, both multi-fidelity variants are 8-dimensional problems. Evaluating the function
189 at a higher fidelity incurs more optimization budget, but yields more accurate observations.

190 4.4 Data mixture optimization (DMO)

191 Optimization over data mixtures can lead to performance gains for both pretraining [54, 75] and
192 supervised fine-tuning (SFT)[14, 48]. Data curriculum adjusts mixture proportions across training
193 stages, which further improves performance but increases search space complexity [57, 79]. There is
194 also growing research interest in examining the optimization of LLM data mixtures over multiple,
195 possibly competing, objectives. Our DMO problem settings span these research areas.

196 Datapoints were collected from real LLM training runs across different training data mixtures in
197 both single and multi-objective settings. Each training data mixture contains 10k samples randomly
198 drawn from three TULU-3 SFT data domains (instruction following, math, and code) [44] at a
199 specified mixing proportion. We perform SFT on a pretrained Qwen3-4B model and evaluate on
200 three objectives: IFEval [80], MATH-500 (4-shot), and MBPP+ [5, 53]. The DMO-4B emulator is
201 trained on 3k+ datapoints, spanning diverse mixture configurations, and predicts all three objective
202 scores given a data mixture curriculum. See Sec. C.2 for more training details.

203 **Problem formulation.** The search space is $\mathcal{X}_{\text{DMO}} \triangleq \left\{ x \in \mathbb{R}_{\geq 0}^6 \mid \sum_{i=1}^3 x_i = 1, \sum_{j=4}^6 x_j = 1 \right\}$,
204 where $x_{1:3}$ and $x_{4:6}$ are the mixture proportions over the three data sources at stage 1 (tokens 0–5M)
205 and stage 2 (tokens 5M–10M), respectively. Full simplex-constrained space is specified in Table 5.
206 Evaluations comprise IFEval, MATH-500, and MBPP+ scores and these are combined into different
207 objectives for each of the three problem variants.

Table 3: Example prompt candidates and corresponding scores

Prompt	Score
Write the solution in a concise, technical style typical of advanced mathematics.	0.810
Break down the problem into smaller subproblems and solve each one sequentially.	0.716
Provide only the final answer without any explanation or intermediate steps.	0.374

208 **Problem variants.** All three variants share \mathcal{X}_{DMO} as the search space and a single multi-output
 209 emulator, differing only in objective definition and noise structure. *DMO* base problem has a single
 210 objective and defines the objective as the unweighted mean of the three scores. This serves as the
 211 entry point for BO on a simplex-constrained space. *Multi-objective DMO* (DMO-MO) defines the
 212 objective function $f_{\text{DMO}} : \mathcal{X}_{\text{DMO}} \rightarrow [0, 1]^3$, returning the three scores directly and targeting Pareto-
 213 front methods such as multi-objective BO [15]. *DMO with heteroscedastic noise* (DMO-Het) has
 214 a single MATH-500 objective and adds $\varepsilon \sim \mathcal{N}(0, \sigma_{\text{DMO}}^2(x))$ to each output, where $\sigma_{\text{DMO}}(x)$ is the
 215 prediction from the noise model. This variant is a natural testbed for noise-aware methods [39, 45].

216 4.5 Prompt optimization (PO)

217 The right prompts can substantially improve performance at inference time [81]. However, the text
 218 space is discrete and unstructured, and its embedding space high-dimensional, making undirected
 219 search inefficient. Prompt selection at test time is also time-sensitive, and long-context inference is
 220 computationally costly. Our PO dataset and emulator makes it assessible for researchers to evaluate
 221 their methods on prompt optimization problems, without needing repeated, expensive inference.

222 We pre-generated a corpus of 5,014 candidate prompts for mathematical reasoning and embedded
 223 each prompt using EmbeddingGemma [62], which supports Matryoshka Representation Learning
 224 (MRL) [43] and can produce meaningful representations at varying dimensionalities. Each prompt
 225 candidate was evaluated on MATH-500 (0-shot), to maximize the influence of the prompt rather than
 226 the few-shot demonstrations. Some candidate prompts and their scores are in Table 3. Tabular data is
 227 used for this task rather than a learned emulator as interpolated embeddings cannot be easily decoded
 228 back to natural language prompts, and selecting from a pre-generated candidate pool is standard
 229 practice in prompt optimization [34, 51, 65, 81]. Further details on corpus generation and evaluation
 230 are in Sec. C.3.

231 **Problem formulation.** The search space is $\mathcal{X}_{\text{PO}} \triangleq \{z_i\}_{i=1}^{5014} \subset \mathbb{R}^d$, where each point z_i is the d -dim
 232 embedding of a candidate prompt i . The objective $f_{\text{PO}} : \mathcal{X}_{\text{PO}} \rightarrow [0, 1]$ is the MATH-500 accuracy
 233 obtained under prompt i from Qwen3-14B model. Acquisition functions are optimized over the
 234 discrete candidate set.

235 **Problem variants.** *PO-128*, *PO-256*, *PO-512*, and *PO-768* correspond to MRL truncations at
 236 $d \in \{128, 256, 512\}$ and the full embedding dimension $d = 768$, respectively.

237 5 Experiments

238 We benchmark a wide range of BO methods on BOLT, as well as some common black-box opti-
 239 mization methods. All experiments have noisy observations with σ_N set to 0.001, except the
 240 heteroscedastic noise problem DMO-Het, where noise is input-dependent and thus requires a per-
 241 point noise model $\sigma_{\text{DMO}}(x)$. For single-objective problems, we plot the logarithm of simple regret,
 242 $\log(f^* - \max_{x_{t'} \leq t} f(x_{t'}))$, where $x_{t'}$ is the observed point at iteration t' . For multi-objective prob-
 243 lems, we plot the logarithm of hypervolume difference, $\log(\text{HV}(\text{PF}_t) - \text{HV}^*)$, where PF_t is the
 244 noiseless Pareto frontier across all observations at iteration t . These are common evaluation metrics
 245 found in optimization literature [15, 67]. The optimal values f^* and HV^* are found empirically via
 246 exhaustive evaluation for tabular data, or gradient ascent of 200 best points found via grid evaluation
 247 for MLP emulators. Each method has 10 initial observations. All GP-based methods are implemented
 248 in BoTorch [6], using an RBF kernel with a dimensionality-scaled lengthscale prior [36]. Results are
 249 averaged over 5 runs and 95% confidence intervals are shown as shaded areas, with initial observa-
 250 tions excluded in the plots. Additional experimental details and discussion on evaluation metrics are
 251 in Sec. D.

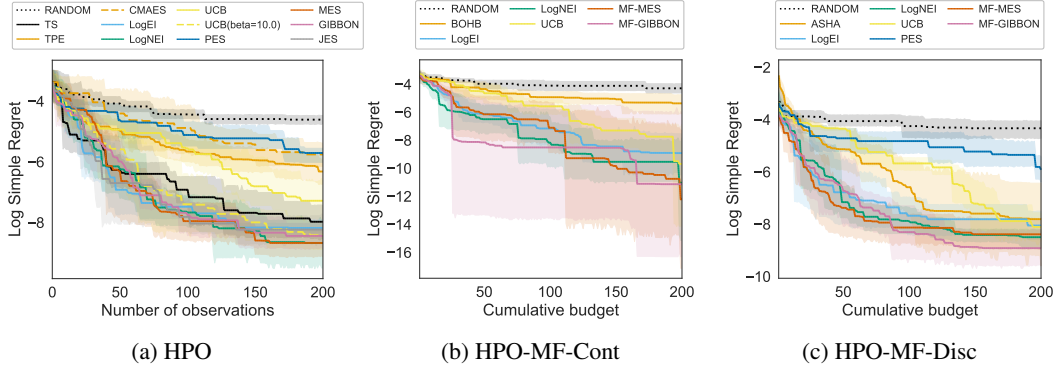


Figure 2: Results on BOLT’s HPO problems. MF variants are plotted against the cumulative budget, where each observation cost range from 0.1 to 1, depending on the fidelity chosen.

252 5.1 HPO experiments and analysis

253 We compare a range of acquisition functions including Thompson sampling (TS), expected im-
 254 provement (EI, NEI), upper confidence bound (UCB), knowledge-gradients (KG), entropy-search
 255 methods (PES, MES, GIBBON), as well as common black-box HPO baselines (TPE, CMA-ES
 256 BOHB, ASHA). For multi-fidelity problems, we adapt acquisition functions to be cost-aware [23]
 257 and show the results w.r.t. the cumulative budget. The cost for each multi-fidelity evaluation was set
 258 to $0.9 \cdot \text{fid}(x) + 0.1$. More implementation details are in Sec. D.1 and results are shown in Fig. 2.

259 **GP-based BO methods generally outperform non-GP-based HPO baselines.** LogNEI, (MF-
 260)MES, and (MF-)GIBBON are the strongest performers across all three settings over 200 iterations.
 261 We note that non-GP baselines have substantially faster wall-clock times (Table 9), although this
 262 overhead is negligible compared to the much more costly, repeated LLM trainings with different
 263 hyperparameters at every iteration.

264 **Evaluation budget is an important consideration for acquisition function selection.** In the
 265 standard HPO setting (Fig. 2a), TS outperforms other methods under 25 iterations but is subsequently
 266 overtaken. UCB faces a similar budget dependency through its β parameter that balances exploration
 267 and exploitation, which we analyze in Sec. D.2.

268 **Cost-scale sensitivity is a practical liability of MF methods.** Multi-fidelity acquisition functions
 269 incorporate a cost utility to discount evaluations at lower fidelities, but there is no principled way to set
 270 this. We select cost scales by comparing early-stopped runs, though this choice can significantly affect
 271 performance and fidelity allocation across methods, affecting the transferability to new problems.
 272 Analysis of acquisitions functions at different cost scales on both MF problems are in Sec. D.3.

273 We include a breakdown of fidelity queries per method for both multi-fidelity problems in Sec. D.4.

274 5.2 DMO experiments and analysis

275 Since LLM finetuning and evaluation are often conducted in parallel, we evaluate both single-
 276 observation and batched acquisition strategies, plotting performance against the number of optimiza-
 277 tion iterations. Batched candidates are selected via sequential greedy conditioning [6]. Batch sizes
 278 and population (for genetic algorithms) are noted where they differ from a default of one.

279 We compare against a range of acquisition functions, including multi-objective variants and an
 280 evolutionary method (NSGA2). Since entropy-based methods (MES, PES, JES) rely on a candidate
 281 set to approximate the optimal-value/optimal-point distribution, we replace the standard box-uniform
 282 grid with Dirichlet-sampled points to respect the simplex constraint on each mixture group. For
 283 DMO-Het, we additionally include noise-aware MLHGP that learns a noise model with expectation
 284 maximization (EM) [39] and a GP with known ground-truth noise to serve as an oracle upper bound.
 285 Additional implementation details are in Sec. D.5 and results are shown in Fig. 3.

286 **NEI-based methods dominate overall.** LogNEI performs best on single-objective problems, while
 287 NEHVI($q = 5$) leads on multi-objective DMO-MO problem. Notably, NEHVI with a single observa-

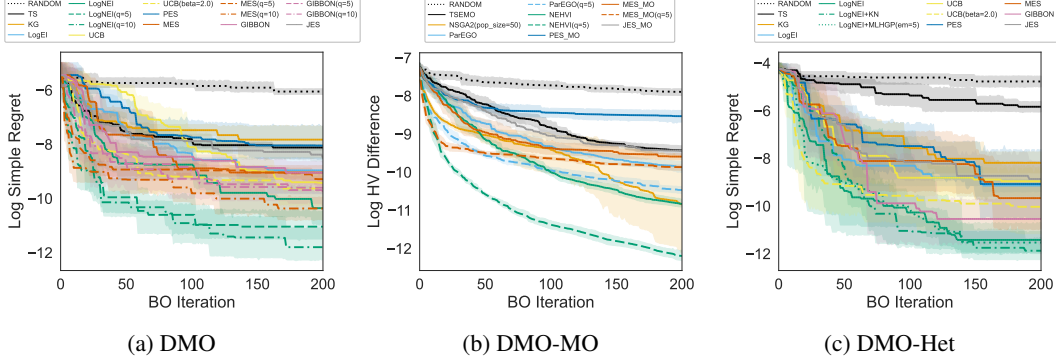


Figure 3: Results on BOLT’s DMO problems. ($q = 5$) indicates a batch size of 5, and NSGA2 has 50 observations in each iteration. All other methods use only 1 observation per iteration.

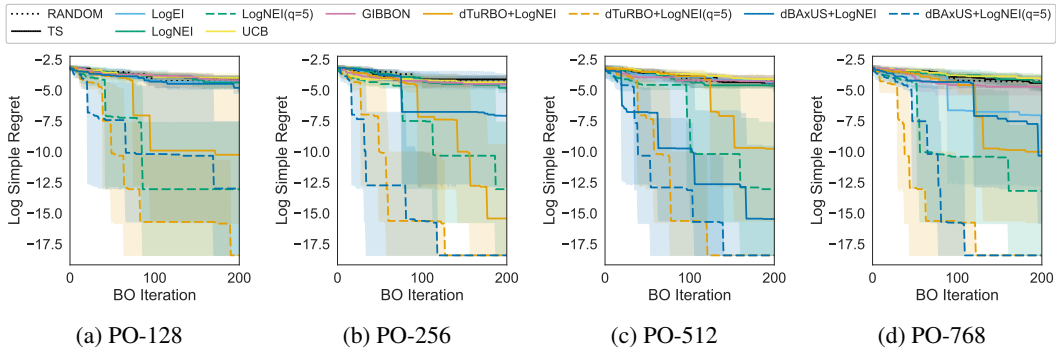


Figure 4: Results on BOLT’s PO problems. ($q = 5$) indicates a batch size of 5. Note that dTURBO and dBaxUS are our adapted methods for a discretized search space.

288 tion per iteration matches NSGA2(pop_size=50) despite using $50\times$ fewer function evaluations per
 289 iteration, and with substantially tighter confidence intervals, indicating more consistent performance.

290 **Batch BO accelerates convergence.** Since batch methods collect more observations per iteration, it
 291 is expected and observed that they outperform single-observation BO over the same iterations. Larger
 292 batch sizes generally accelerate early convergence, though gains can diminish over time, for instance,
 293 MES($q = 5$) plateaus after 50 iterations. The effect of batch size is acquisition-function dependent:
 294 GIBBON shows little difference between $q = 5$ and $q = 10$, whereas LogNEI and MES exhibit more
 295 pronounced gains, possibly reflecting differences in the acquisition landscape across methods.

296 **Noise-aware methods yield modest improvements under heteroscedastic noise.** When noise is
 297 known (LogNEI+KN), performance surpasses standard homoscedastic LogNEI in 200 iterations,
 298 confirming that correct noise modelling is beneficial. LogNEI+MLHGP performed marginally better
 299 than LogNEI, suggesting the EM-learned noise model captures sufficient noise structure.

300 5.3 PO experiments and analysis

301 Beyond benchmarking acquisition functions, we evaluate trust-region and subspace-reduction meth-
 302 ods, specifically TuRBO [20] and BaxUS [59]. Since PO problems operate over a discrete candidate
 303 set, we introduce discretized variants of these methods, which we term *dTuRBO* and *dBaxUS*. In these
 304 variants, the trust-region hyperrectangle is replaced by the k nearest neighbors of the current best
 305 solution in embedding space, where k is doubled upon success and halved upon failure. Consequently,
 306 Sobol candidate generation and sparse perturbation are no longer required. We omit SAASBO [19]
 307 due to its prohibitive computational cost, and MSR [60] as it is inapplicable to discretized search
 308 spaces. Finally, for acquisition functions that support batch selection, we include results with batch
 309 size 5, reflecting the parallel evaluation setting common in practice. More implementation details,
 310 including on dTuRBO and dBaxUS, are provided in Sec. D.6. Results are shown in Fig. 4.

311 **Standard acquisition functions fail to scale.** LogEI, LogNEI, UCB, TS, and GIBBON largely
 312 plateau between -3 and -5 log simple regret within the first 50 iterations, with negligible improve-
 313 ment thereafter, making them indistinguishable from random search. A notable exception is LogEI at
 314 PO-768, which descends to approximately -7 , suggesting it benefits from the richer expressivity of
 315 the full embedding dimensionality, though it remains far behind trust-region methods.

316 **Batch evaluation partially compensates for the absence of trust-region structure.** Log simple
 317 regret of $\text{LogNEI}(q = 5)$ reaches approximately -12.5 for all PO problems without any local search
 318 constraints, far exceeding its sequential counterpart and showing that parallel acquisition provides
 319 meaningful gains independent of trust-region structures.

320 **Local search methods lead to substantial performance gains.** Across all problem dimensions,
 321 dTuRBO+LogNEI and dBaxUS+LogNEI reach log simple regret of approximately -10 to -15 ,
 322 demonstrating that the k -NN trust region adaptation preserves the effectiveness of local search in
 323 discretized search spaces. Combining local search with batch evaluation yields further gains, with
 324 dTuRBO+LogNEI($q = 5$) and dBaxUS+LogNEI($q = 5$) achieving the lowest regret and approaching
 325 the optimal solution within 200 iterations, across almost all PO problems.

326 **Subspace reduction may be unnecessary for LLM embedding spaces.**

327 At PO-128, dBaxUS performs notably worse than dTuRBO, as PO-
 328 128 embeddings require 74.2% of principal components to explain 95%
 329 of variance (Fig. 5), violating the low-dimensional active subspace as-
 330 sumption underlying BaxUS [59]. At higher dimensionalities, where
 331 variance concentrates in fewer principal components, dBaxUS recovers
 332 but results are mixed for single-observation runs. With batch size 5,
 333 dBaxUS matches dTuRBO but converges more slowly, likely due to
 334 wasted evaluations in very low-dimensional subspaces before expansion.
 335 This suggests subspace reduction is most effective when the embedding
 336 space has genuine dimensional redundancy, a property that may not hold
 337 for LLM embeddings, as they tend to be broadly informative across dimensions.

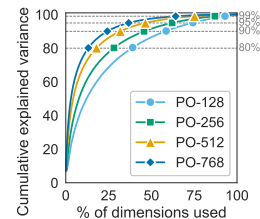


Figure 5: PCA explained variance for PO

338 6 Limitations

339 To collect sufficient data to fit reliable emulators for BOLT, we required practical measures such as
 340 smaller training token budgets, reduced maximum token limits, and, for HPO tasks, the exclusion of
 341 chain-of-thought reasoning. This may not represent the best achievable performance of the evaluated
 342 models under unconstrained inference. Nevertheless, the optimization challenges BOLT is designed to
 343 benchmark (mixed-variable search spaces, simplex-constraints, and more) are intrinsic to the
 344 structure of LLM tasks themselves and remain consistent across inference and training configurations.
 345 We also plan to update the benchmark with new LLM architectures in the future, in an effort to
 346 continue modernizing optimization benchmarks.

347 7 Conclusion and Future Work

348 Black-box optimization methods hold significant promise for automating and accelerating LLM
 349 pipeline optimization, but progress has been hampered by the lack of accessible and reproducible
 350 evaluations grounded in real LLM experiments. BOLT is critical to unlocking this progress, by
 351 providing emulators and tabular datasets built on real LLM experiments, enabling fast, reproducible
 352 benchmarking without large-scale compute.

353 Our evaluation reveals that mature BBO methods, in particular, NEI- and GIBBON-based approaches,
 354 transfer effectively to LLM tasks, but strong performance required adaptation to LLM-specific
 355 structural challenges, including simplex-constrained and discretized search spaces. These challenges
 356 reflect the structure of real LLM pipeline optimization problems and highlight gaps for the BBO
 357 community to tackle.

358 BOLT is extensible to more optimization problems, for example contextual and constrained optimiza-
 359 tion, as new optimization settings can be constructed using our publicly available emulators and data.
 360 There are also many other LLM-centric tasks that fit the black-box optimization framework, such as
 361 decoding configuration and RLHF reward weighting, that can be future extensions of BOLT.

362 **References**

- 363 [1] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama. Optuna: A next-generation hyperparameter
364 optimization framework. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge
365 discovery & data mining*, pages 2623–2631, 2019.
- 366 [2] A. Albalak, Y. Elazar, S. M. Xie, S. Longpre, N. Lambert, X. Wang, N. Muennighoff, B. Hou, L. Pan,
367 H. Jeong, et al. A survey on data selection for language models. *arXiv preprint arXiv:2402.16827*, 2024.
- 368 [3] S. Ament, S. Daulton, D. Eriksson, M. Balandat, and E. Bakshy. Unexpected improvements to expected
369 improvement for Bayesian optimization. *Advances in Neural Information Processing Systems*, 36:20577–
370 20612, 2023.
- 371 [4] S. P. Arango, H. S. Jomaa, M. Wistuba, and J. Grabocka. Hpo-b: A large-scale reproducible benchmark
372 for black-box hpo based on openml. In *Thirty-fifth Conference on Neural Information Processing Systems
373 Track on Datasets and Benchmarks*, 2021.
- 374 [5] J. Austin, A. Odena, M. Nye, M. Bosma, H. Michalewski, D. Dohan, E. Jiang, C. Cai, M. Terry, Q. Le,
375 et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.
- 376 [6] M. Balandat, B. Karrer, D. Jiang, S. Daulton, B. Letham, A. G. Wilson, and E. Bakshy. BoTorch: A
377 framework for efficient monte-carlo bayesian optimization. *Advances in Neural Information Processing
378 Systems*, 33:21524–21538, 2020.
- 379 [7] S. Belakaria, A. Deshwal, and J. R. Doppa. Max-value entropy search for multi-objective Bayesian
380 optimization. *Advances in Neural Information Processing Systems*, 32, 2019.
- 381 [8] J. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl. Algorithms for hyper-parameter optimization. *Advances
382 in Neural Information Processing Systems*, 24, 2011.
- 383 [9] D. Bingham and S. Surjanovic. Virtual library of simulation experiments: Test functions and datasets,
384 2013. URL <https://www.sfu.ca/~ssurjano/optimization.html>.
- 385 [10] L. Bliet, A. Guijt, R. Karlsson, S. Verwer, and M. De Weerd. Benchmarking surrogate-based optimisation
386 algorithms on expensive black-box functions. *Applied Soft Computing*, 147:110744, 2023.
- 387 [11] E. Bradford, A. M. Schweidtmann, and A. Lapkin. Efficient multiobjective optimization employing
388 Gaussian processes, spectral sampling and a genetic algorithm. *Journal of global optimization*, 71(2):
389 407–438, 2018.
- 390 [12] O. Chapelle and L. Li. An empirical evaluation of Thompson sampling. *Advances in Neural Information
391 Processing Systems*, 24, 2011.
- 392 [13] L. Chen, J. Chen, T. Goldstein, H. Huang, and T. Zhou. Instructzero: Efficient instruction optimization for
393 black-box large language models. In *International Conference on Machine Learning*, 2024.
- 394 [14] Z. Chen, G. K. R. Lau, C.-S. Foo, and B. K. H. Low. DUET: Optimizing training data mixtures via
395 feedback from unseen evaluation tasks. *arXiv:2502.00270*, 2025.
- 396 [15] S. Daulton, M. Balandat, and E. Bakshy. Parallel bayesian optimization of multiple noisy objectives with
397 expected hypervolume improvement. *Advances in Neural Information Processing Systems*, 34:2187–2200,
398 2021.
- 399 [16] K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan. A fast and elitist multiobjective genetic algorithm:
400 NSGA-II. *IEEE transactions on evolutionary computation*, 6(2):182–197, 2002.
- 401 [17] K. Drezkowski, A. Grosnit, and H. Bou Ammar. Framework and benchmarks for combinatorial and
402 mixed-variable Bayesian optimization. *Advances in Neural Information Processing Systems Track on
403 Datasets and Benchmarks*, 36:69464–69489, 2023.
- 404 [18] K. Eggenberger, P. Müller, N. Mallik, M. Feurer, R. Sass, A. Klein, N. Awad, M. Lindauer, and F. Hutter.
405 Hpobench: A collection of reproducible multi-fidelity benchmark problems for HPO. In *Neural Information
406 Processing Systems Track on Datasets and Benchmarks*, 2021.
- 407 [19] D. Eriksson and M. Jankowiak. High-dimensional Bayesian optimization with sparse axis-aligned sub-
408 spaces. In *Uncertainty in artificial intelligence*, pages 493–503. PMLR, 2021.
- 409 [20] D. Eriksson, M. Pearce, J. Gardner, R. D. Turner, and M. Poloczek. Scalable global optimization via local
410 Bayesian optimization. *Advances in Neural Information Processing Systems*, 32, 2019.
- 411 [21] S. Falkner, A. Klein, and F. Hutter. Bohb: Robust and efficient hyperparameter optimization at scale. In
412 *International Conference on Machine Learning*, pages 1437–1446. PMLR, 2018.
- 413 [22] C. Fernando, D. S. Banarse, H. Michalewski, S. Osindero, and T. Rocktäschel. Promptbreeder: Self-
414 referential self-improvement via prompt evolution. In *International Conference on Machine Learning*,
415 pages 13481–13544. PMLR, 2024.
- 416 [23] P. Frazier, W. Powell, and S. Dayanik. The knowledge-gradient policy for correlated normal beliefs.
417 *INFORMS journal on Computing*, 21(4):599–613, 2009.

- 418 [24] P. I. Frazier. A tutorial on Bayesian optimization. *arXiv preprint arXiv:1807.02811*, 2018.
- 419 [25] L. Gao, J. Tow, B. Abbasi, S. Biderman, S. Black, A. DiPofi, C. Foster, L. Golding, J. Hsu, A. Le Noac’h,
420 H. Li, K. McDonell, N. Muennighoff, C. Ociepa, J. Phang, L. Reynolds, H. Schoelkopf, A. Skowron,
421 L. Sutawika, E. Tang, A. Thite, B. Wang, K. Wang, and A. Zou. The language model evaluation harness,
422 07 2024. URL <https://zenodo.org/records/12608602>.
- 423 [26] R. Garnett. *Bayesian optimization*. Cambridge University Press, 2023.
- 424 [27] N. Hansen. The CMA evolution strategy: a comparing review. *Towards a new evolutionary computation: Advances in the estimation of distribution algorithms*, pages 75–102, 2006.
- 425 [28] N. Hansen, S. Finck, R. Ros, and A. Auger. *Real-parameter black-box optimization benchmarking 2009: Noiseless functions definitions*. PhD thesis, INRIA, 2009.
- 426 [29] N. Hansen, A. Auger, R. Ros, O. Mersmann, T. Tušar, and D. Brockhoff. COCO: A platform for comparing
429 continuous optimizers in a black-box setting. *Optimization Methods and Software*, 36(1):114–144, 2021.
- 430 [30] F. Häse, M. Aldeghi, R. J. Hickman, L. M. Roch, M. Christensen, E. Liles, J. E. Hein, and A. Aspuru-Guzik.
431 Olympos: a benchmarking framework for noisy optimization and experiment planning. *Machine Learning: Science and Technology*, 2(3):035021, 2021.
- 432 [31] D. Hendrycks, C. Burns, S. Kadavath, A. Arora, S. Basart, E. Tang, D. Song, and J. Steinhardt. Measuring
433 mathematical problem solving with the MATH dataset. In *Thirty-fifth Conference on Neural Information Processing Systems Track on Datasets and Benchmarks*, 2021.
- 434 [32] D. Hernández-Lobato, J. Hernandez-Lobato, A. Shah, and R. Adams. Predictive entropy search for multi-
437 objective Bayesian optimization. In *International Conference on Machine Learning*, pages 1492–1501.
438 PMLR, 2016.
- 439 [33] J. M. Hernández-Lobato, M. W. Hoffman, and Z. Ghahramani. Predictive entropy search for efficient
440 global optimization of black-box functions. *Advances in Neural Information Processing Systems*, 27, 2014.
- 441 [34] W. Hu, Y. Shu, Z. Yu, Z. Wu, X. Lin, Z. Dai, S.-K. Ng, and B. K. H. Low. Localized zeroth-order prompt
442 optimization. *Advances in Neural Information Processing Systems*, 37:86309–86345, 2024.
- 443 [35] C. Hvarfner, F. Hutter, and L. Nardi. Joint entropy search for maximally-informed Bayesian optimization.
444 *Advances in Neural Information Processing Systems*, 35:11494–11506, 2022.
- 445 [36] C. Hvarfner, E. O. Hellsten, and L. Nardi. Vanilla Bayesian optimization performs great in high dimensions.
446 In *International Conference on Machine Learning*, pages 20793–20817. PMLR, 2024.
- 447 [37] C. Jang, H. Lee, J. Kim, and J. Lee. Model fusion through Bayesian optimization in language model
448 fine-tuning. *Advances in Neural Information Processing Systems*, 37:29878–29912, 2024.
- 449 [38] D. R. Jones, M. Schonlau, and W. J. Welch. Efficient global optimization of expensive black-box functions.
450 *Journal of Global optimization*, 13(4):455–492, 1998.
- 451 [39] K. Kersting, C. Plagemann, P. Pfaff, and W. Burgard. Most likely heteroscedastic gaussian process
452 regression. In *International Conference on Machine Learning*, pages 393–400, 2007.
- 453 [40] J. Knowles. Parego: A hybrid algorithm with on-line landscape approximation for expensive multiobjective
454 optimization problems. *IEEE transactions on evolutionary computation*, 10(1):50–66, 2006.
- 455 [41] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa. Large language models are zero-shot reasoners.
456 *Advances in Neural Information Processing Systems*, 35:22199–22213, 2022.
- 457 [42] A. Kumar, G. Wu, M. Z. Ali, R. Mallipeddi, P. N. Suganthan, and S. Das. A test-suite of non-convex
458 constrained optimization problems from the real-world and some baseline results. *Swarm and evolutionary
459 computation*, 56:100693, 2020.
- 460 [43] A. Kusupati, G. Bhatt, A. Rege, M. Wallingford, A. Sinha, V. Ramanujan, W. Howard-Snyder, K. Chen,
461 S. Kakade, P. Jain, et al. Matryoshka representation learning. *Advances in Neural Information Processing
462 Systems*, 35:30233–30249, 2022.
- 463 [44] N. Lambert, J. Morrison, V. Pyatkin, S. Huang, H. Ivison, F. Brahman, L. J. V. Miranda, A. Liu, N. Dziri,
464 X. Lyu, et al. Tulu 3: Pushing frontiers in open language model post-training. In *Second Conference on
465 Language Modeling*, 2025.
- 466 [45] M. Lázaro-Gredilla and M. K. Titsias. Variational heteroscedastic Gaussian process regression. In
467 *International Conference on Machine Learning*, pages 841–848, 2011.
- 468 [46] H. Li, W. Zheng, Q. Wang, H. Zhang, Z. Wang, S. Xuyang, Y. Fan, Z. Ding, H. Wang, N. Ding, S. Zhou,
469 X. Zhang, and D. Jiang. Predictable scale: Part i, step law – optimal hyperparameter scaling law in large
470 language model pretraining. *arXiv:2503.04715*, 2025.
- 471 [47] L. Li, K. Jamieson, A. Rostamizadeh, E. Gonina, J. Ben-Tzur, M. Hardt, B. Recht, and A. Talwalkar. A
472 system for massively parallel hyperparameter tuning. *Proceedings of machine learning and systems*, 2:
473 230–246, 2020.

- 474 [48] Y. Li, Z. Liu, and E. Xing. Data mixing optimization for supervised fine-tuning of large language models.
475 In *International Conference on Machine Learning*, pages 35419–35437. PMLR, 2025.
- 476 [49] Q. Liang, A. E. Gongora, Z. Ren, A. Tiihonen, Z. Liu, S. Sun, J. R. Deneault, D. Bash, F. Mekki-Berrada,
477 S. A. Khan, et al. Benchmarking the performance of Bayesian optimization across multiple experimental
478 materials science domains. *npj Computational Materials*, 7(1):188, 2021.
- 479 [50] H. Lightman, V. Kosaraju, Y. Burda, H. Edwards, B. Baker, T. Lee, J. Leike, J. Schulman, I. Sutskever, and
480 K. Cobbe. Let’s verify step by step. In *The Twelfth International Conference on Learning Representations*,
481 2023.
- 482 [51] X. Lin, Z. Wu, Z. Dai, W. Hu, Y. Shu, S.-K. Ng, P. Jaillet, and B. K. H. Low. Use your INSTINCT: Instruc-
483 tion optimization for llms using neural bandits coupled with transformers. In *International Conference on*
484 *Machine Learning*, pages 30317–30345. PMLR, 2024.
- 485 [52] M. Lindauer, K. Eggensperger, M. Feurer, A. Biedenkapp, D. Deng, C. Benjamins, T. Ruhkopf, R. Sass,
486 and F. Hutter. Smac3: A versatile Bayesian optimization package for hyperparameter optimization. *Journal*
487 *of Machine Learning Research*, 23(54):1–9, 2022.
- 488 [53] J. Liu, C. S. Xia, Y. Wang, and L. Zhang. Is your code generated by chatgpt really correct? rigorous
489 evaluation of large language models for code generation. *Advances in Neural Information Processing*
490 *Systems*, 36:21558–21572, 2023.
- 491 [54] Q. Liu, X. Zheng, N. Muennighoff, G. Zeng, L. Dou, T. Pang, J. Jiang, and M. Lin. Regmix: Data mixture
492 as regression for language model pre-training. In *The Thirteenth International Conference on Learning*
493 *Representations*, 2025.
- 494 [55] H. B. Moss, D. S. Leslie, J. Gonzalez, and P. Rayson. Gibbon: General-purpose information-based
495 Bayesian optimisation. *Journal of Machine Learning Research*, 22(235):1–49, 2021.
- 496 [56] M. Nomura, S. Watanabe, Y. Akimoto, Y. Ozaki, and M. Onishi. Warm starting CMA-ES for hyperpa-
497 rameter optimization. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages
498 9188–9196, 2021.
- 499 [57] T. Olmo, A. Ettinger, A. Bertsch, B. Kuehl, D. Graham, D. Heineman, D. Groeneveld, F. Brahman,
500 F. Timbers, H. Ivison, et al. Olmo 3. *arXiv preprint arXiv:2512.13961*, 2025.
- 501 [58] P. S. Palar, R. P. Liem, L. R. Zuhail, and K. Shimoyama. On the use of surrogate models in engineering
502 design optimization and exploration: The key issues. In *Proceedings of the genetic and evolutionary*
503 *computation conference companion*, pages 1592–1602, 2019.
- 504 [59] L. Papenmeier, L. Nardi, and M. Poloczek. Increasing the scope as you learn: Adaptive Bayesian
505 optimization in nested subspaces. *Advances in Neural Information Processing Systems*, 35:11586–11601,
506 2022.
- 507 [60] L. Papenmeier, M. Poloczek, and L. Nardi. Understanding high-dimensional Bayesian optimization. In
508 *International Conference on Machine Learning*, pages 47902–47923. PMLR, 2025.
- 509 [61] F. Pfisterer, L. Schneider, J. Moosbauer, M. Binder, and B. Bischl. Yahpo gym-an efficient multi-objective
510 multi-fidelity benchmark for hyperparameter optimization. In *International Conference on Automated*
511 *Machine Learning*, pages 3–1. PMLR, 2022.
- 512 [62] H. Schechter Vera, S. Dua, B. Zhang, D. Salz, R. Mullins, S. Raghuram Panyam, S. Smoot, I. Naim, J. Zou,
513 F. Chen, D. Cer, A. Lisak, M. Choi, L. Gonzalez, O. Sanseviero, G. Cameron, I. Ballantyne, K. Black,
514 K. Chen, W. Wang, Z. Li, G. Martins, J. Lee, M. Sherwood, J. Ji, R. Wu, J. Zheng, J. Singh, A. Sharma,
515 D. Sreepat, A. Jain, A. Elarabawy, A. Co, A. Doumanoglou, B. Samari, B. Hora, B. Potetz, D. Kim,
516 E. Alfonseca, F. Moiseev, F. Han, F. Palma Gomez, G. Hernández Ábrego, H. Zhang, H. Hui, J. Han,
517 K. Gill, K. Chen, K. Chen, M. Shanbhogue, M. Boratko, P. Suganthan, S. M. K. Duddu, S. Mariserla,
518 S. Ariafar, S. Zhang, S. Zhang, S. Baumgartner, S. Goenka, S. Qiu, T. Dabral, T. Walker, V. Rao,
519 W. Khawaja, W. Zhou, X. Ren, Y. Xia, Y. Chen, Y.-T. Chen, Z. Dong, Z. Ding, F. Visin, G. Liu, J. Zhang,
520 K. Kenealy, M. Casbon, R. Kumar, T. Mesnard, Z. Gleicher, C. Brick, O. Lacombe, A. Roberts, Y. Sung,
521 R. Hoffmann, T. Warkentin, A. Joulain, T. Duerig, and M. Seyedhosseini. Embeddingemma: Powerful and
522 lightweight text representations. *arXiv preprint arXiv:2509.20354*, 2025.
- 523 [63] B. Seong-Eun, L. Jung-Mok, K. Sung-Bin, and T.-H. Oh. Efficient hyper-parameter search for LoRA via
524 language-aided Bayesian optimization. *arXiv preprint arXiv:2602.11171*, 2026.
- 525 [64] B. Shahriari, K. Swersky, Z. Wang, R. P. Adams, and N. De Freitas. Taking the human out of the loop: A
526 review of Bayesian optimization. *Proceedings of the IEEE*, 104(1):148–175, 2015.
- 527 [65] C. Shi, K. Yang, Z. Chen, J. Li, J. Yang, and C. Shen. Efficient prompt optimization through the lens of
528 best arm identification. *Advances in Neural Information Processing Systems*, 37:99646–99685, 2024.
- 529 [66] J. Snoek, H. Larochelle, and R. P. Adams. Practical Bayesian optimization of machine learning algorithms.
530 *Advances in Neural Information Processing Systems*, 25, 2012.

- 531 [67] N. Srinivas, A. Krause, S. M. Kakade, and M. Seeger. Gaussian process optimization in the bandit setting:
532 No regret and experimental design. In *International Conference on Machine Learning*, 2010.
- 533 [68] S. Takeno, H. Fukuoka, Y. Tsukada, T. Koyama, M. Shiga, I. Takeuchi, and M. Karasuyama. Multi-fidelity
534 Bayesian optimization with max-value entropy search and its parallelization. In *International Conference*
535 *on Machine Learning*, pages 9334–9345. PMLR, 2020.
- 536 [69] S. Toshniwal, W. Du, I. Moshkov, B. Kisacanin, A. Ayrapetyan, and I. Gitman. Openmathinstruct-2:
537 Accelerating ai for math with massive open-source instruction data. In *The Thirteenth International*
538 *Conference on Learning Representations*, 2023.
- 539 [70] C. Tribes, S. Benarroch-Lelong, P. Lu, and I. Kobyzev. Hyperparameter optimization for large language
540 model instruction-tuning. In *AAAI Conference on Artificial Intelligence*, 2024.
- 541 [71] B. Tu, A. Gandy, N. Kantas, and B. Shafei. Joint entropy search for multi-objective Bayesian optimization.
542 *Advances in Neural Information Processing Systems*, 35:9922–9938, 2022.
- 543 [72] L. Wang, W. Xu, Y. Lan, Z. Hu, Y. Lan, R. K.-W. Lee, and E.-P. Lim. Plan-and-solve prompting: Improving
544 zero-shot chain-of-thought reasoning by large language models. In *Proceedings of the 61st annual meeting*
545 *of the association for computational linguistics (volume 1: long papers)*, pages 2609–2634, 2023.
- 546 [73] Z. Wang and S. Jegelka. Max-value entropy search for efficient Bayesian optimization. In *International*
547 *Conference on Machine Learning*, pages 3627–3635. PMLR, 2017.
- 548 [74] Z. Wu, X. Lin, Z. Dai, W. Hu, Y. Shu, S.-K. Ng, P. Jaillet, and B. K. H. Low. Prompt optimization
549 with ease? efficient ordering-aware automated selection of exemplars. *Advances in Neural Information*
550 *Processing Systems*, 37:122706–122740, 2024.
- 551 [75] S. M. Xie, H. Pham, X. Dong, N. Du, H. Liu, Y. Lu, P. S. Liang, Q. V. Le, T. Ma, and A. W. Yu.
552 Doremi: Optimizing data mixtures speeds up language model pretraining. *Advances in Neural Information*
553 *Processing Systems*, 36:69798–69818, 2023.
- 554 [76] A. Yang, A. Li, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Gao, C. Huang, C. Lv, et al. Qwen3
555 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- 556 [77] C. Yang, X. Wang, Y. Lu, H. Liu, Q. V. Le, D. Zhou, and X. Chen. Large language models as optimizers.
557 In *The Twelfth International Conference on Learning Representations*, 2023.
- 558 [78] T. Yen, A. W. T. Siah, H. Chen, C. D. Guetta, T. Peng, and H. Namkoong. Data mixture optimization: A
559 multi-fidelity multi-scale bayesian framework. In *Advances in Neural Information Processing Systems*,
560 2025.
- 561 [79] Y. Zhang, A. Mohamed, H. Abdine, G. Shang, and M. Vazirgiannis. Beyond random sampling: Efficient
562 language model pretraining via curriculum learning. In *Proceedings of the 19th Conference of the European*
563 *Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5776–5794,
564 2026.
- 565 [80] J. Zhou, T. Lu, S. Mishra, S. Brahma, S. Basu, Y. Luan, D. Zhou, and L. Hou. Instruction-following
566 evaluation for large language models. *arXiv preprint arXiv:2311.07911*, 2023.
- 567 [81] Y. Zhou, A. I. Muresanu, Z. Han, K. Paster, S. Pitis, H. Chan, and J. Ba. Large language models are
568 human-level prompt engineers. In *The Eleventh International Conference on Learning Representations*,
569 2022.

570 **A BOLT usage examples**

```

1 from bolt import HPO
2 import torch
3
4 # Instantiate a problem
5 problem = HPO(noise_std=0.001)
6
7 # BoTorch-compatible metadata
8 problem.bounds          # (2, 7) tensor of lower/upper bounds
9 problem.dim             # search-space dimensionality
10 problem.continuous_inds # indices of continuous variables
11 problem.discrete_inds  # indices of integer variables
12 problem.categorical_inds # indices of categorical variables
13
14 # Sample random candidates and round integer/categorical variables
15 lb, ub = problem.bounds
16 X = lb + torch.rand(10, problem.dim) * (ub - lb)
17 X[:, problem.discrete_inds] = X[:, problem.discrete_inds].round()
18 X[:, problem.categorical_inds] = X[:, problem.categorical_inds].round()
19 y = problem(X) # (10, 1)

```

Listing 1: BOLT usage example.

571 **B Additional emulator details and analysis**

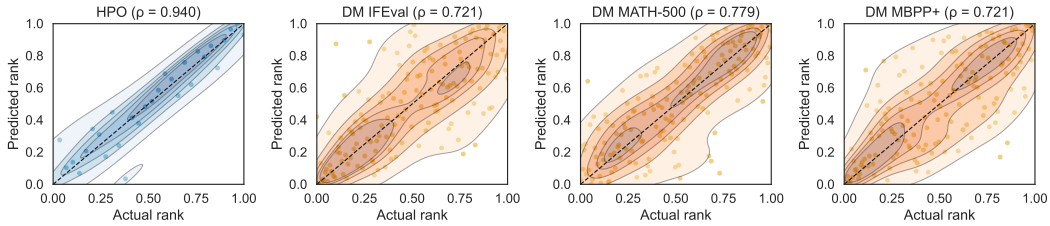


Figure 6: Real and predicted ranks on test set.

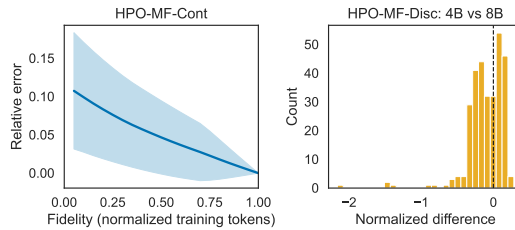


Figure 7: Differences in emulated objective at different fidelities

572 **Emulator training.** We chose 2-layer MLP for objective emulators to balance the trade off between
573 model complexity and amount of training data needed, to prevent overfitting. Dropout and layer norm
574 was also used during training for regularization.

575 **Actual and predicted rank comparison** of data points in the test set is shown Fig. 6, where predicted
576 ranks largely align with actual ranks. For the DMO emulator, lower- and higher-ranked points are
577 generally better aligned than intermediate ones, which could be due to mixture compositions in the
578 intermediate range being more similar and harder to discriminate.

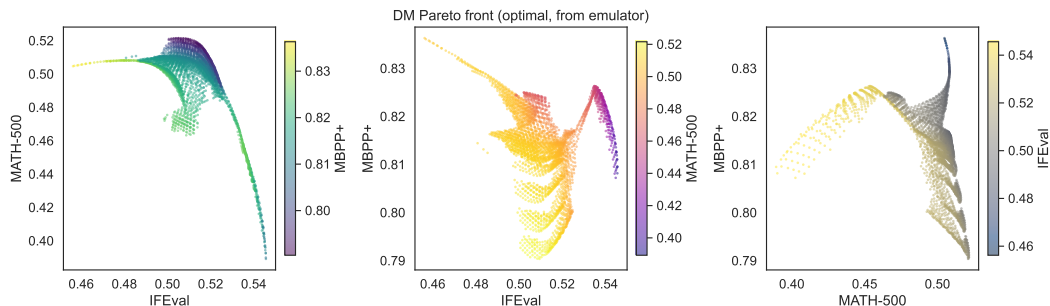


Figure 8: Estimated optimal Pareto front, as 2D slices of a 3D space

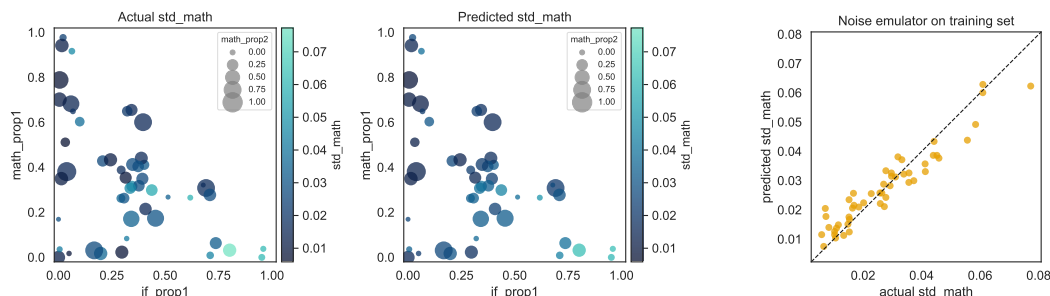


Figure 9: Noise emulator predictions for MATH-500 standard deviation. `if_prop1` indicates mixture proportion of information following data at stage 1 of the curriculum, while `math_prop1` and `math_prop2` represent math data proportions at stage 1 and 2 respectively.

579 **Multi-fidelity settings.** We analyze the HPO emulators under multi-fidelity settings in Fig. 7. For
 580 continuous fidelity, as expected, error relative to the highest fidelity is largest at lower fidelities. For
 581 discrete fidelity, differences largely cluster around -0.6 to 0.2 , indicating that the 4B model serves as
 582 a reliable proxy for the 8B model’s performance.

583 **Multi-objective setting.** We visualize the optimal Pareto frontier of the DMO-4B emulator in Fig. 8,
 584 shown as three 2D slices of the joint IFEval–MATH-500–MBPP+ objective space. A clear trade-off
 585 is visible between IFEval and MATH-500, forming a concave frontier in that projection. MBPP+
 586 varies more smoothly across the frontier, suggesting code performance is less in tension with the
 587 other two objectives. Overall, the frontier spans a diverse set of Pareto-optimal mixtures, reflecting
 588 meaningful flexibility across objective priorities.

589 **For the noise model,** hyperparameters such as kernel choice were selected via k -fold cross-validation
 590 with $k = 5$. The final model was trained on all 50 data samples using the selected hyperparameters: a
 591 Laplacian kernel with $\gamma = 0.1$. This model choice was motivated by the limited number of samples,
 592 as each data point requires multiple training runs and evaluations, making large-scale collection
 593 computationally expensive.

594 **Visualization of the noise model predictions** are shown in Fig. 9. The predicted noise distribution
 595 closely mirrors the actual, with higher noise concentrated at low `if_prop1` (stage 1 information
 596 following proportion) and high `math_prop1` (stage 1 math proportion) values. The predicted vs.
 597 actual plot on the training set also shows points tracking closely along the diagonal, confirming that
 598 the noise emulator fits the training data well. Some underestimation is visible at higher `std_math`
 599 values, which is expected given the limited number of high-noise samples in the 50-point training set.

600 C Details on data generation for emulators

601 C.1 Hyperparameter optimization

602 OpenMathInstruct-2 [69] was selected as the training data. We used only the answer for training,
 603 due to time and compute constraint. Long reasoning chains increase the training and evaluation time

604 substantially. Evaluation was similarly done by prompting for only an answer without reasoning, on
 605 MATH-500.

606 We fine-tune from Qwen3 Base models rather than the Instruct models for this task, for more variance
 607 in the evaluation scores. 2 H200 was used for each training run, and checkpoints are saved every 200
 608 steps for evaluation.

609 HPO’s mixed-variable search space is specified in Table 4.

Table 4: Mixed-variable search space for HPO.

Parameter	Type	Range
Learning rate	Cont.	$[10^{-6}, 10^{-3}]$ normalized to $[0, 1]$
Per device batch size (2^n)	Int.	$[2, 4]$ (effective batch size is $2 \times$)
LoRA rank (2^n)	Int.	$[2, 5]$
LoRA alpha (2^n)	Int.	$[2, 5]$
LoRA dropout (2^n)	Cont.	$[0, 0.1]$ normalized to $[0, 1]$
LoRA layers	Int.	$[1, 30]$
LoRA target modules	Cat.	0: (q_proj, v_proj) 1: (q_proj, v_proj, k_proj, o_proj) 2: (gate_proj, up_proj, down_proj) 3: (all-linear)
[HPO-MF-Cont] Token fidelity	Cont.	$[10^5, 9.1 \times 10^6]$ normalized to $[0, 1]$
[HPO-MF-Disc] Model fidelity	Cat.	0: HPO-4B emulator 1: HPO-8B emulator

610 C.2 Data mixture optimization

611 We construct training mixtures from three datasets spanning instruction following, math, and
 612 coding: `tulu-3-sft-personas-instruction-following`, `tulu-3-sft-personas-math`, and
 613 `tulu-3-sft-personas-code` [44]. Each training run has two stages: stage 1 (tokens 0-5M) and
 614 stage 2 (tokens 5M-10M). At each stage, 10k samples are drawn from the three datasets following
 615 a mixing proportion. This roughly corresponds to 5M tokens or less for each 1 epoch. Training
 616 hyperparameters are: learning rate 10^{-4} , per device batch size 16 (effective batch size 32), lora rank
 617 8, lora alpha 8, lora dropout 0.05, lora layers 30, and lora target modules `q_proj, v_proj`.

618 Evaluations were done with `lm-eval` [25]. We use the `inst_level_strict_acc` metric for IFEval,
 619 `math_verify` metric for MATH-500 (4-shot, minerva format), and `pass@1` for MBPP+.

Table 5: Simplex-constrained search space for DMO.

Parameter	Type	Range	Constraint
x_1	IF proportion (stage 1)	$[0, 1]$	$x_1 \geq 0, \sum_{i=1}^3 x_i = 1$
x_2	Math proportion (stage 1)	$[0, 1]$	$x_2 \geq 0, \sum_{i=1}^3 x_i = 1$
x_3	Code proportion (stage 1)	$[0, 1]$	$x_3 \geq 0, \sum_{i=1}^3 x_i = 1$
x_4	IF proportion (stage 2)	$[0, 1]$	$x_4 \geq 0, \sum_{j=4}^6 x_j = 1$
x_5	Math proportion (stage 2)	$[0, 1]$	$x_5 \geq 0, \sum_{j=4}^6 x_j = 1$
x_6	Code proportion (stage 2)	$[0, 1]$	$x_6 \geq 0, \sum_{j=4}^6 x_j = 1$

$\mathcal{X}_{\text{DMO}} \triangleq \{x \in \mathbb{R}_{\geq 0}^6 \mid \sum_{i=1}^3 x_i = 1, \sum_{j=4}^6 x_j = 1\}$. Effective free parameters: 4.

620 C.3 Prompt optimization

621 C.3.1 Prompt generation

622 Qwen3-14B was used to generate a set of prompts for the candidate set, using multiple meta-prompts
 623 to ensure a diverse set of generated prompts. Prompts are embedding using `EmbeddingGemma` [62].

624 Duplicate prompts were removed via exact text match and thresholding on embeddings distance. In
 625 addition to generated prompts, common prompts selected by prompt optimization papers in Table 6
 626 were also added to the candidate set. Meta-prompts used to generate the rest of the prompts are
 627 shown in Table 7.

Table 6: Prompts from prior prompt optimization methods.

Method	Prompt
Zero-Shot CoT [41]	“Let’s think step by step.”
APE [81]	“Let’s work this out in a step by step way to be sure we have the right answer.”
OPRO [77]	“Take a deep breath and work through this problem step-by-step.” “Break this down.” “A little bit of arithmetic and a logical approach will help us quickly arrive at the solution to this problem.” “Let’s combine our numerical command and clear thinking to quickly and accurately decipher the answer.” “Let’s work through this problem step-by-step.”
PromptBreeder [22]	“SOLUTION:”
ZOPO [34]	“Let’s find the solution by using the given information.”
INSTINCT [51]	“Let’s think about it.”
InstructZero [13]	“Let’s use the instruction to solve the problem.”
Plan-and-Solve+ [72]	“Let’s first understand the problem, extract relevant variables and their corresponding numerals, and make a complete plan. Then, let’s carry out the plan, calculate intermediate variables (pay attention to correct numerical calculation and commonsense), solve the problem step by step, and show the answer.”

628 C.3.2 Prompt evaluation

629 We evaluate candidate instruction prompts on the MATH-500 benchmark [31] using the `lm-eval`
 630 library [25], on Qwen3-14B in non-thinking mode. The candidate instruction is placed in the system
 631 instructions. All evaluations use greedy decoding (temperature $T = 0$), 4-shot prompting, and
 632 `math_verify` as the scoring metric. We cap generation at 1024 tokens.

633 The best instruction prompt recovered by our BO search is: “Write the solution in a concise, technical
 634 style typical of advanced mathematics.”, achieving 81.0% accuracy on MATH-500. Note that the
 635 optimal instruction is specific to our generation configuration, and may differ especially under a
 636 different token budget. Our cap of 1024 tokens is substantially smaller than those used in most
 637 frontier model papers, for example, the evaluation results in Qwen3 use a max output length of
 638 32,768 [76]. We chose a smaller budget to keep evaluation cost tractable across 5014 prompt
 639 candidates evaluations. Qwen3-14B in non-thinking mode is also able to produce concise, direct
 640 answers that fit comfortably within 1024 tokens.

641 D More experiment details

642 Monte Carlo (MC)-based implementations were used for LogNEI, PES, MES, GIBBON, JES, and
 643 KG acquisition functions, with 128 quasi-MC samples using Sobol sequences per iteration.

644 D.1 HPO experiments

645 For the single-fidelity HPO problem, we compare random sampling, Thompson sampling (TS) [12],
 646 log expected improvement (LogEI) [3, 38], log noisy expected improvement (LogNEI) [3], upper
 647 confidence bounds (UCB) [67], predictive entropy search (PES) [33], max-value entropy search
 648 (MES) [73], GIBBON [55], and joint entropy search (JES) [35]. Knowledge gradient (KG) [23]
 649 was excluded as it was prohibitively slow for mixed-variable search spaces. These methods were
 650 implemented using the BoTorch library [6].

651 To handle the mixed-variable space, we use a sum-and-product composite kernel comprising a
 652 Hamming-distance kernel for categorical variables and an RBF kernel for continuous and discrete

Table 7: Meta-prompt categories and generation prompts. All meta-prompts follow the template: “Generate $[N]$ diverse and distinct short instructions (1–4 sentences each) that $[DESCRIPTION]$. Each instruction should be phrased as a direct directive to the solver. Output one instruction per line with no numbering or extra formatting.” The table shows only the $[DESCRIPTION]$ component and count N for each category.

Category	Description	Count
Reasoning Style	tell a math solver to use a specific <i>reasoning style</i> (e.g., step-by-step breakdown, chain-of-thought, working backwards, proof by contradiction, elimination of wrong choices, estimation first, drawing a diagram, considering special cases)	1000
Verbosity	tell a math solver to operate at a specific <i>level of detail or verbosity</i> (e.g., give only the final numerical answer, briefly sketch the key steps, show every algebraic manipulation, explain each line as if teaching a beginner, include intermediate checks)	1000
Persona	adopt a specific <i>voice or role</i> for solving math problems (e.g., address the reader as a student being guided, write as a peer collaborator, explain as a strict grader, narrate as if thinking aloud, use the tone of a competition preparation coach)	1000
Format	specify a <i>notation or output format</i> for math solutions (e.g., write all equations in LaTeX, box the final answer, use numbered lines for each step, present results as a table when possible, state assumptions explicitly before solving)	1000
Subject-Specific	guide a math solver to <i>identify and leverage the relevant mathematical structure</i> of a problem before solving (e.g., identify whether the problem is algebraic, geometric, or combinatorial, reframe the problem in the most natural mathematical language, check whether the problem reduces to a known result)	1000
Common Heuristics	are common, general-purpose <i>problem-solving heuristics</i> applicable to any math problem (e.g., think step-by-step, verify your answer at the end, re-read the problem carefully before starting, check your solution against special cases, estimate the answer before computing)	1000
Adversarial	are <i>unhelpful or counterproductive</i> for solving math problems (e.g., vague directives with no clear method, instructions to skip showing work, contradictory guidance, focus on irrelevant topics, overly restrictive rules that prevent flexible reasoning)	1000
CoT Triggers	serve as <i>zero-shot chain-of-thought triggers</i> — brief directives that prompt a solver to reason carefully before giving the final answer, as alternatives to “Let’s think step by step” (e.g., “Take a deep breath and work through this carefully.”, “Reason slowly and methodically.”, “Walk through your logic before concluding.”)	200
CoT Structured	impose an <i>explicit reasoning structure</i> on the solving process (e.g., plan your approach before computing, write out a scratchpad of intermediate values, number each logical sub-step, annotate each line with the rule or property used, decompose the problem into clearly labelled stages)	200
CoT Verification	ask the solver to <i>verify or sanity-check</i> during or after solving (e.g., substitute the answer back into the original equation, estimate the order of magnitude and check against it, consider a simple special case to test your formula, check each step for arithmetic or sign errors)	200
CoT Decomposition	guide the solver to <i>decompose the problem</i> into manageable parts before solving (e.g., identify all given quantities and unknowns, break into sub-problems and solve each in turn, reduce to a simpler analogous case first, separate into cases and handle each independently)	200
CoT Reflection	encourage <i>metacognitive self-monitoring</i> during solving — prompting the solver to reflect on their own reasoning process (e.g., pause after each step and ask whether it follows logically, flag uncertain steps and revisit them, backtrack if a contradiction is reached, ask “does this answer make intuitive sense?”)	200

653 ordinal variables. Acquisition function optimization alternates between local search over discrete
654 and categorical dimensions, and gradient-based optimization over continuous dimensions, with
655 continuous relaxation and rounding applied to high-cardinality discrete variables (e.g., number of
656 LoRA layers, with 30 discrete values) [6].

657 For multi-fidelity problems, we compare MF-MES [68], MF-GIBBON, as well as cost-scaled variants
658 of single-fidelity acquisition functions (LogEI, LogNEI, UCB, and PES). In cost-scaled acquisition
659 functions, acquisition values are divided by an affine cost model $c(x) = \alpha \cdot \text{fid}(x) + 1$, where $\text{fid}(x)$
660 is the fidelity parameter and α controls the cost sensitivity. For log-valued acquisition functions,
661 $\log(c(x))$ is subtracted from the acquisition values instead. We set $\alpha = 0.005$ for UCB and PES,
662 $\alpha = 0.01$ for EI, and $\alpha = 0.05$ for NEI, MF-MES, and MF-GIBBON. A sensitivity analysis over α
663 is in Sec. D.3.

664 Finally, we include dedicated HPO baselines from popular HPO libraries: tree-structured Parzen
 665 estimator (TPE) [8] and covariance matrix adaptation evolution strategy (CMA-ES) [27, 56] from
 666 Optuna [1] for the single-fidelity setting problem, BOHB [21] from SMAC3 [52] for the continuous-
 667 fidelity setting, and asynchronous successive halving algorithm (ASHA) [47] from Optuna for the
 668 discrete-fidelity setting. Note that CMA-ES cannot handle categorical variables and will fall back to
 669 random sampling for the categorical LoRA target module variable.

670 D.2 Discussion on UCB’s exploration parameter

671 Results comparing different β in UCB are shown in Fig. 10. While β_t following Srinivas et al.
 672 [67] often overexplores in practice, a fixed β allows direct control over the exploration-exploitation
 673 tradeoff and thus the convergence behavior. Smaller β values reduce exploration, enabling faster
 674 initial convergence but potentially at the cost of solution optimality, as seen for $\beta = 1.0$ in both
 675 HPO and DMO plots. Conversely, larger β values converge more slowly but tend to yield better
 676 final solutions. This suggests that β should be chosen with the compute budget in mind. To balance
 677 exploration and exploitation effectively, we selected $\beta = 10.0$ for HPO and $\beta = 2.0$ for DMO for
 678 comparison in main regret plots in Fig. 2a and Fig. 3a, where these values enable rapid convergence
 679 to a good solution.

680 These results also suggest that an adaptive β schedule could be beneficial in practice, starting with a
 681 smaller β to exploit promising regions early, then increasing it over iterations to ensure sufficient
 682 exploration for final solution quality. The schedule of Srinivas et al. [67] is designed for theoretical
 683 no-regret guarantees in the infinite-horizon setting and grows too aggressively for finite budgets. A
 684 schedule calibrated to the evaluation budget could offer a principled middle ground between rapid
 685 early convergence and high final solution quality, which we leave for future work.

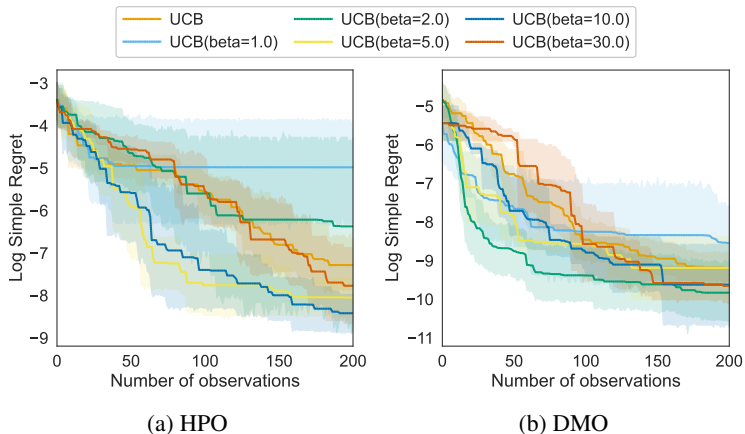


Figure 10: Comparison of β for UCB on HPO and DMO base problems. "UCB" refers to using scheduled β_t from Srinivas et al. [67]

686 D.3 Cost scale sensitivity

687 In cost-scaled acquisition functions, acquisition values are divided by an affine cost model $c(x) =$
 688 $\alpha \cdot \text{fid}(x) + 1$, where $\text{fid}(x)$ is the fidelity parameter and α controls the cost sensitivity. Larger α
 689 penalizes high-fidelity queries more heavily in their acquisition scores, encouraging more low-fidelity
 690 exploration. We evaluate sensitivity to the cost scale parameter α , comparing $\alpha \in \{0.005, 0.01, 0.05\}$
 691 on both multi-fidelity HPO variants (continuous token-based fidelity and discrete model-size fidelity)
 692 in Fig. 11. The degree of sensitivity varies across acquisition functions and problems. On the
 693 discrete-fidelity problem (HPO-MF-Disc), methods are generally less sensitive to cost scale, likely
 694 because the binary fidelity choice limits the ways in which cost weighting can reshape query patterns.
 695 On the continuous-fidelity problem (HPO-MF-Cont), LogNEI, MF-MES, and MF-GIBBON exhibit
 696 the greatest sensitivity, with certain cost scales yielding substantially better performance over all
 697 50 iterations. MF-GIBBON in particular shows large spread across α values. While LogEI, UCB,
 698 and PES are comparatively robust across cost scales on both problems, they also achieve worse final

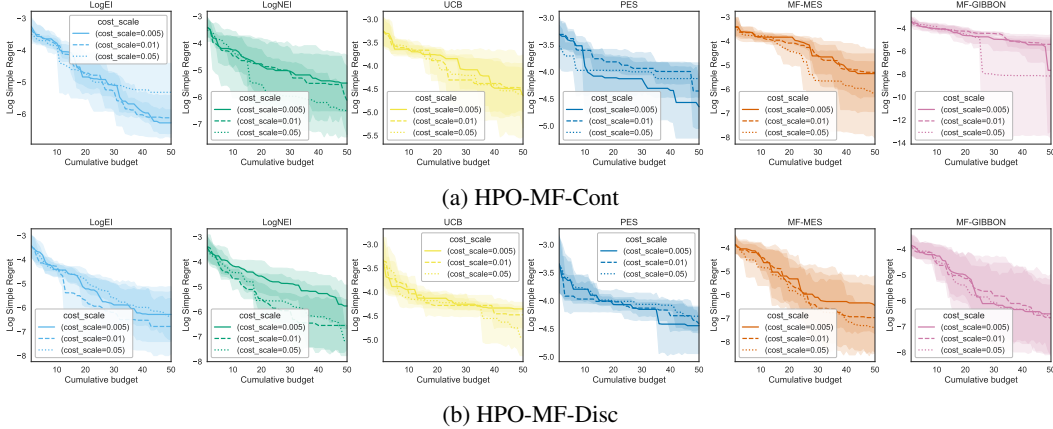


Figure 11: Performance on MF HPO problems on different cost scale.

699 performance. This suggests that cost scale sensitivity and strong performance may go hand-in-hand:
 700 methods that respond to cost scale can exploit it to achieve better results, making cost scale an
 701 important parameter that requires tuning.

702 D.4 Fidelity proportions over iterations

703 Figure 12 shows how each method allocates queries across fidelity levels over the course of optimization.
 704 Recall the acquisition functions are all cost scaled, with $\alpha = 0.005$ for UCB and PES, $\alpha = 0.01$
 705 for EI, and $\alpha = 0.05$ for NEI, MF-MES, and MF-GIBBON.

706 For the continuous-fidelity setting (HPO-MF-Cont), BOHB concentrates mass at a small number of
 707 discrete fidelity levels determined by its bracket schedule, as expected from its Hyperband-based
 708 design. All GP-based methods, including MF-MES and MF-GIBBON, rapidly saturate near fidelity=1
 709 after a brief initial exploration phase, with only a sparse cloud of low-fidelity evaluations in the early
 710 queries. Interestingly, LogEI, LogNEI, and UCB each exhibit a secondary cluster of observations
 711 near the lowest fidelity, indicating a bimodal allocation pattern in which intermediate fidelities are
 712 largely avoided. This does not appear to harm performance, as LogNEI achieves competitive regret
 713 with MF-MES on this task (Figure 2b).

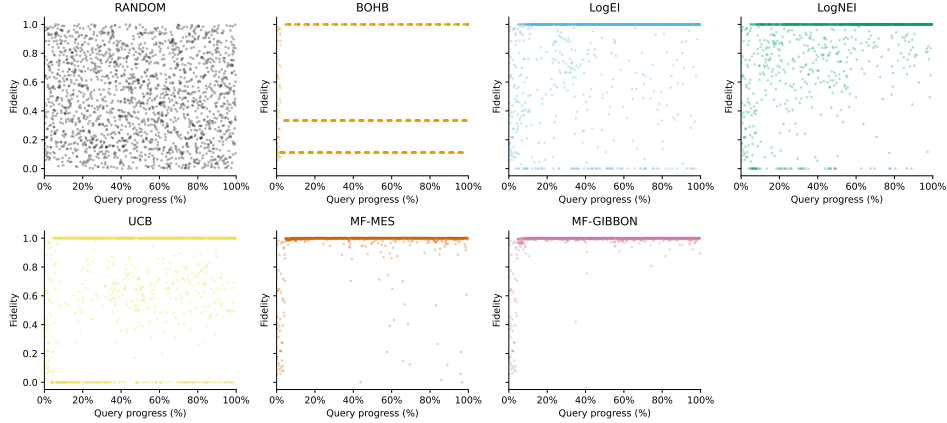
714 For the discrete-fidelity setting (HPO-MF-Disc), acquisition-based methods generally shift toward
 715 high-fidelity evaluations as optimization progresses, suggesting a tendency to exploit the high-fidelity
 716 region despite its greater cost. The notable exception is MF-MES, which begins with nearly all
 717 high-fidelity queries and progressively reallocates toward the cheap surrogate ($\text{fid}=0$) as the budget is
 718 consumed. The empirical performance of MF-MES on this task, achieving the third-lowest regret
 719 behind MF-GIBBON and LogNEI (Figure 2c), suggests that this allocation is not detrimental.

720 Given that methods with vastly different fidelity allocation strategies can perform well, these results
 721 suggest a complex relationship between fidelity allocation and final optimization quality, one that
 722 likely depends on the specific configurations evaluated at each fidelity level rather than the allocation
 723 pattern alone.

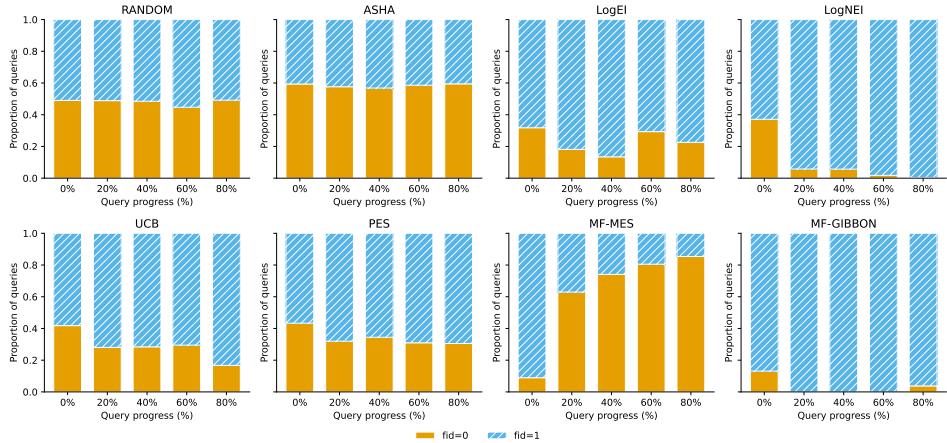
724 D.5 DMO experiments

725 For the single-objective problems, we compare random sampling, TS, LogEI, LogNEI, KG, UCB,
 726 PES, MES, GIBBON, and JES. For the multi-objective problem, we benchmarked TSEMO [11],
 727 NSGA2 [16], ParEGO [40], NEHVI [15], PESMO [32] MESMO [7], and JESMO [71].

728 Information-theoretic acquisition functions (PES, MES, GIBBON, JES, and their multi-objective
 729 variants) require approximations of the global optimum or Pareto-optimal set. Since standard box-
 730 constrained sampling cannot satisfy the simplex equality constraints of the data-mixture search space,
 731 we instead draw feasible candidates via symmetric Dirichlet sampling and select approximate optima
 732 by Thompson sampling and taking the argmax of each posterior draw over the candidate set. For
 733 JES, optimal inputs and values are drawn from the same posterior sample to ensure consistency.



(a) HPO-MF-Cont



(b) HPO-MF-Disc

Figure 12: Fidelity proportions of observations on MF HPO problems.

734 For multi-objective entropy-search methods, we draw multiple joint posterior samples over feasible
 735 candidates, extract the non-dominated set from each sample as an approximate Pareto front.

736 The reference HV point HV^* was set per objective as 10% below the observed minimum, with slack
 737 scaled to each objective’s range. Minima were computed from surrogate predictions over a dense
 738 regular grid on the mixture-ratio simplex.

739 D.6 PO experiments

740 For PO experiments, we compare random sampling, TS, LogEI, LogNEI, UCB, and GIBBON, as well
 741 as batched, dTuRBO, and dBaxUS variants with LogNEI. dTuRBO and dBaxUS are our discretized
 742 adaptations of TuRBO and BaxUS.

743 TuRBO searches within trust-regions for scalability, expanding and contracting trust-region hy-
 744 perrectangles upon successive improvements and failures respectively, while BaxUS searches in
 745 reduced-dimensional subspaces via random projection and applies trust-region filtering within that
 746 subspace. TuRBO and BaxUS are originally designed for continuous spaces, particularly in their
 747 trust-region methodology.

748 We adapted TuRBO and BaxUS to discretized search spaces by representing the trust regions as
 749 the k nearest neighbors of the current best solution, where k is doubled after τ_{succ} consecutive
 750 improvements and halved after τ_{fail} consecutive non-improvements. The trust region is initialized
 751 with $k_{\text{init}} = \lfloor N/2 \rfloor$ candidates, where $N = 5014$ is the total prompt pool size. Both methods use
 752 $\tau_{\text{succ}} = 3$, dTuRBO uses a fixed $\tau_{\text{fail}} = 10$, while for dBaxUS, τ_{fail} is set dynamically per subspace

753 expansion stage, calibrated so that the number of failures required to fully contract the trust region
 754 matches the per-stage evaluation budget, following the continuous-space contraction rate in BAXUS.
 755 We show how selected enhancements (batch size, dTuRBO, dBAXUS) affect the simple regret at
 756 iteration 200 in Table 8.

Table 8: Delta in log simple regret at iteration 200 vs LogNEI baseline (mean \pm std over 5 seeds). Negative is better as it indicates lower log regret.

Enhancement	PO128	PO256	PO512	PO768
q=5	-8.67 ± 7.75	-8.23 ± 7.61	-8.40 ± 7.54	-8.85 ± 7.67
TuRBO	-5.87 ± 7.55	-10.61 ± 7.10	-5.12 ± 7.78	-5.67 ± 7.67
TuRBO+q=5	$-14.05 \pm 0.63^\dagger$	$-13.61 \pm 0.46^\dagger$	$-13.78 \pm 0.41^\dagger$	$-14.09 \pm 0.87^\dagger$
BAXUS	-0.43 ± 1.09	-2.27 ± 6.57	-10.82 ± 6.72	-5.99 ± 7.38
BAXUS+q=5	-8.61 ± 6.88	$-13.61 \pm 0.46^\dagger$	$-13.78 \pm 0.41^\dagger$	$-14.09 \pm 0.87^\dagger$

† All 5 seeds hit the log-regret floor (≈ -18.42); deltas coincide when methods fully floor out.

757 D.7 On evaluation metric

758 In Sec. 5, for single-objective problems, we evaluate performance using simple regret, which is
 759 defined as $f^* - \max_{t' \leq t} f(x_{t'})$. While simple regret provides a clean picture of optimization progress
 760 over the budget t , it requires access to the noiseless f in their max term, which is computable on
 761 benchmarks but unavailable in practice.

762 Another popular evaluation metric is inference regret, defined as $f^* - f(\hat{x}_t)$, where $\hat{x}_t \triangleq$
 763 $\arg \max_{x \in \mathcal{X}} \mu_t(x)$ is the point that maximizes the posterior mean at iteration t . Inference regret
 764 is commonly used in entropy search-based methods [73], and it utilizes the surrogate model in
 765 estimating the optimal point. We additionally report the *per-step* variant of simple regret $f^* - f(\tilde{x}_t)$,
 766 where $\tilde{x}_t \triangleq \arg \max_{t' \leq t} y_{t'}$ is the best point by noisy observations. Both of these metrics are not
 767 guaranteed to be monotonically non-increasing, but matches practical evaluation conditions more
 768 closely.

769 Inference regret plots are Fig. 13 and per-step simple regret plots are in Fig. 14. With inference
 770 regret, methods are less differentiable in some problems like HPO, suggesting that they may have
 771 similar posterior mean despite different observations. Nevertheless, insights drawn from both plots
 772 are generally similar to the main simple regret plots in Sec. 5: LogNEI, MF-MES, and MF-GIBBON
 773 performs well for multi-fidelity HPO-MF-Disc and DMO experiments, methods using noise models
 774 and information perform better in heteroscedastic noise settings, and batched and dTuRBO methods
 775 are important for high-dimensional PO problems.

776 D.8 Wall-clock time for experiments

777 Wall-clock runtimes are reported in Table 9 and Table 10. Note that these runtimes include emulator
 778 query time, which is negligible. Most experiments ran on L40 GPUs, with a subset using H100 or
 779 H200 GPUs to accommodate the higher VRAM demands of the GPU-accelerated implementation,
 780 which scales with the number of observations and input dimensionality (e.g., PO-768 has $d = 768$
 781 which is memory intensive).

782 E Compute usage

783 We detail the compute used for generation of BOLT’s data, and BO experiments in Table 11.

784 F Broader Impact

785 Our work introduces BOLT, a benchmark to spur development of BBO methods for LLM config-
 786 uration optimization, and does not directly affect LLM model outputs where societal impacts are
 787 more pronounced. By enabling more efficient LLM pipeline optimization, BOLT could indirectly

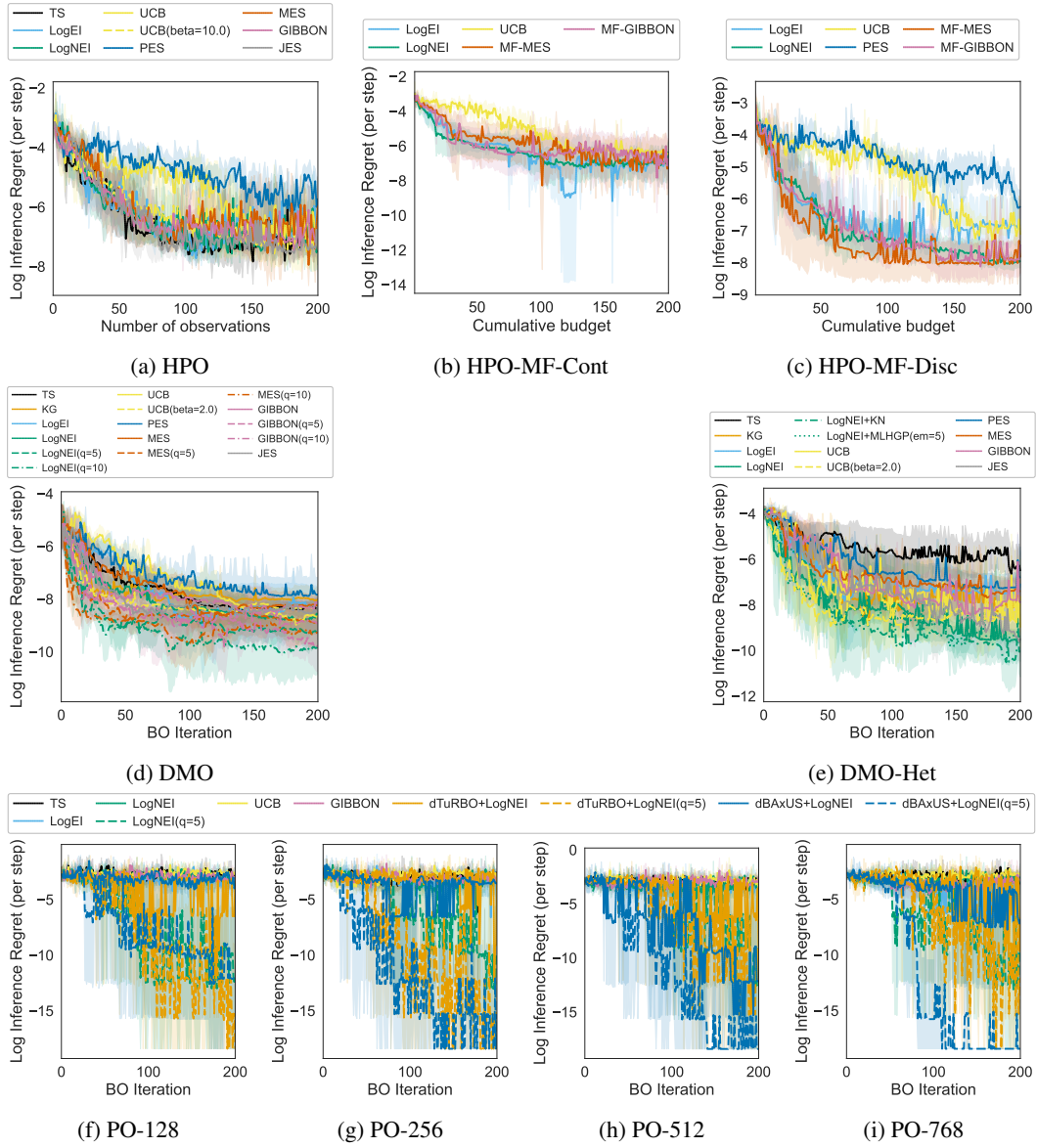


Figure 13: Results on BoLT's HPO problems using inference regret. Random is not included as random sampling does not maintain a posterior mean.

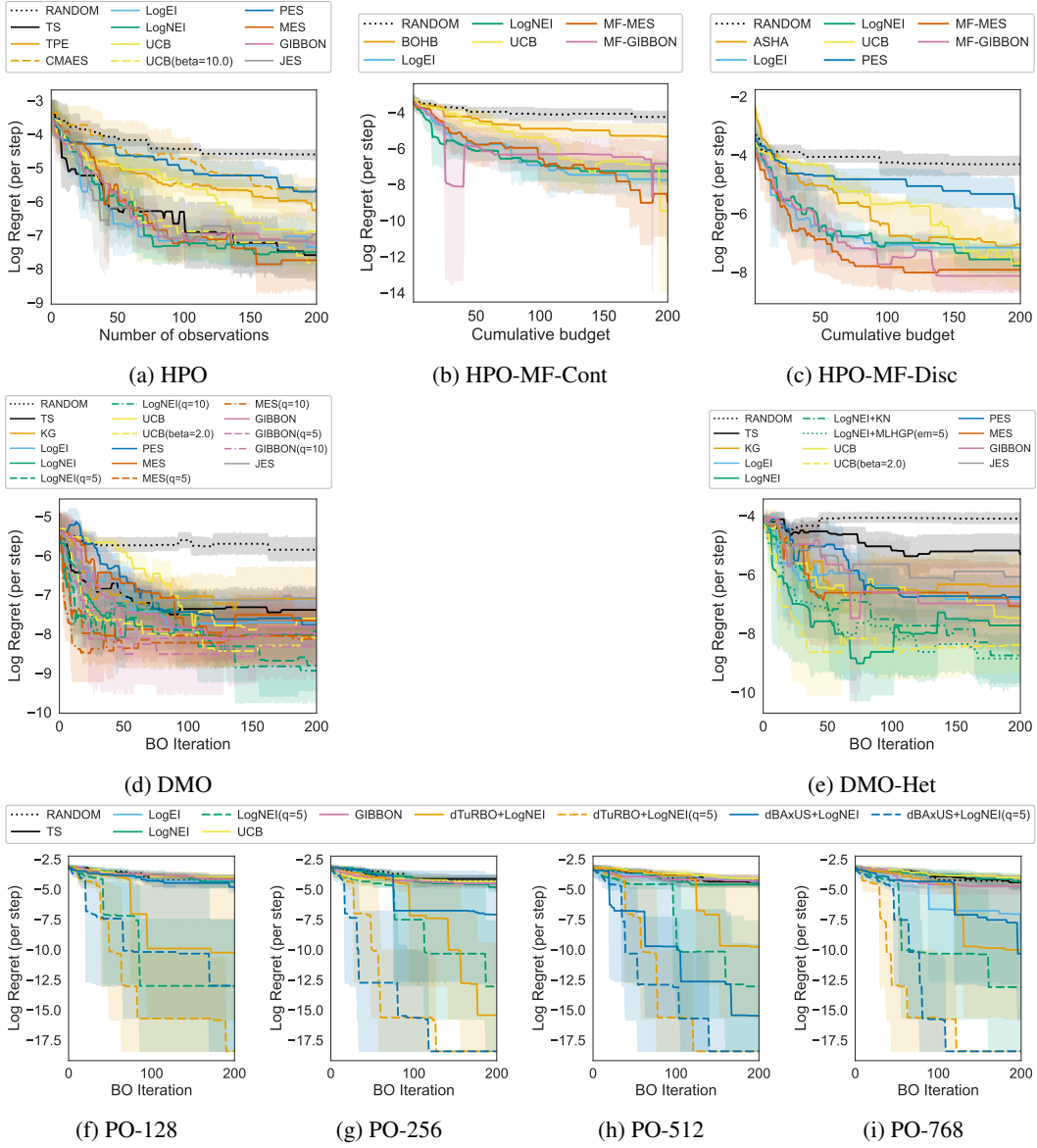


Figure 14: Results on BOLT's HPO problems using per-step simple regret.

Table 9: Mean wall-clock time per BO step — HPO and DMO.

(a) HPO		(b) DMO	
Method	Time/step	Method	Time/step
<i>HPO</i>		<i>DMO</i>	
RANDOM	0.00186 s	RANDOM	0.00509 s
TS	10.8 s	TS	2.02 s
TPE	0.0113 s	KG	21 s
LogEI	13.5 s	LogEI	3.7 s
LogNEI	12.9 s	LogNEI	5.53 s
UCB	12.5 s	LogNEI(q=5)	17.1 s
UCB(beta=10.0)	12.8 s	LogNEI(q=10)	1.8 min
PES	16.3 s	UCB	3.1 s
MES	17 s	UCB(beta=2.0)	3.15 s
GIBBON	14.6 s	PES	7.03 s
JES	17.1 s	MES	3.31 s
<i>HPO-MF-Cont</i>		MES(q=5)	7.83 s
RANDOM	0.00654 s	MES(q=10)	14.5 s
BOHB	0.665 s	GIBBON	3.17 s
LogEI	17.6 s	GIBBON(q=5)	21.4 s
LogNEI	15.8 s	GIBBON(q=10)	52.8 s
UCB	22 s	JES	4.1 s
MF-MES	15.8 s	<i>DMO-MO</i>	
MF-GIBBON	13.6 s	RANDOM	0.0246 s
<i>HPO-MF-Disc</i>		TSEMO	2.35 s
RANDOM	0.00525 s	NSGA2(pop_size=50)	8.78 s
ASHA	0.0175 s	ParEGO	4.65 s
LogEI	16.4 s	ParEGO(q=5)	19.3 s*
LogNEI	12.6 s	NEHVI	1.35 min
UCB	24.5 s	NEHVI(q=5)	6.33 min
PES	41.7 s*	PES_MO	16.6 s
MF-MES	29.4 s	MES_MO	9.14 s
MF-GIBBON	18.5 s	MES_MO(q=5)	53 s*
		JES_MO	12.8 s
		<i>DMO-Het</i>	
		RANDOM	0.00541 s
		TS	2.03 s
		KG	25.2 s
		LogEI	3.44 s
		LogNEI	6.59 s
		LogNEI+KN	5.98 s
		LogNEI+MLHGP(em=5)	7.69 s
		UCB	2.96 s
		UCB(beta=2.0)	3.38 s
		PES	4.38 s
		MES	3.11 s
		GIBBON	2.92 s
		JES	4.42 s

*Ran on H100. †Ran on H200. Otherwise ran on L40

788 contribute to making capable LLMs more accessible to smaller organizations and researchers with
789 limited resources. On the negative side, accelerating LLM development pipelines could indirectly
790 contribute to faster proliferation of capable language models, amplifying existing societal concerns
791 around misuse.

Table 10: Mean wall-clock time per BO step — PO.

Method	<i>PO-128</i>	<i>PO-256</i>	<i>PO-512</i>	<i>PO-768</i>
RANDOM	0.00746 s	0.00966 s	0.0146 s	0.0112 s
TS	8.19 s	9.04 s	8.76 s*	4.98 s [†]
LogEI	10.3 s	12 s	6.82 s*	1.76 s [†]
LogNEI	8.55 s	10.6 s	7.15 s*	6.54 s [†]
LogNEI(q=5)	30.3 s	59.4 s*	1.01 min*	40.7 s [†]
UCB	7.65 s	9.34 s	7.22 s*	4.65 s [†]
GIBBON	9.57 s	9.94 s	7.34 s*	4.89 s [†]
TuRBO+LogNEI	9.8 s	10.3 s	8.85 s*	3.84 s [†]
TuRBO+LogNEI(q=5)	14.9 s	16.1 s	13.5 s*	7.4 s [†]
BAXUS+LogNEI	7.77 s	7.04 s*	7.13 s*	4.55 s [†]
BAXUS+LogNEI(q=5)	8.15 s	7.21 s*	9.79 s*	7.68 s [†]

*Ran on H100. [†]Ran on H200. Otherwise ran on L40

Table 11: Compute Usage Summary

Stage	Model	Training	Eval
HPO	8B	29 H200 GPU-days	11 L40 GPU-days
	4B	21 H200 GPU-days	7 L40 GPU-days
DMO	4B	66 H200 GPU-days	29 L40 GPU-days
PO	14B	–	21 H200 GPU-days