
HPOBench: A Collection of Reproducible Multi-Fidelity Benchmark Problems for HPO

Katharina Eggensperger^{1*}, Philipp Müller¹, Neeratyoy Mallik¹, Matthias Feurer¹, René Sass², Aaron Klein^{3†}, Noor Awad¹, Marius Lindauer², Frank Hutter^{1,4}

¹ Albert-Ludwigs-Universität Freiburg ² Leibniz Universität Hannover

³ Amazon ⁴ Bosch Center for Artificial Intelligence

Abstract

1 To achieve peak predictive performance, hyperparameter optimization (HPO) is a
2 crucial component of machine learning and its applications. Over the last years,
3 the number of efficient algorithms and tools for HPO grew substantially. At
4 the same time, the community is still lacking realistic, diverse, computationally
5 cheap, and standardized benchmarks. This is especially the case for multi-fidelity
6 HPO methods. To close this gap, we propose *HPOBench*, which includes 7
7 existing and 5 new benchmark families, with a total of more than 100 multi-
8 fidelity benchmark problems. *HPOBench* allows to run this extendable set of
9 multi-fidelity HPO benchmarks in a reproducible way by isolating and packaging
10 the individual benchmarks in containers. It also provides surrogate and tabular
11 benchmarks for computationally affordable yet statistically sound evaluations. To
12 demonstrate *HPOBench*'s broad compatibility with various optimization tools,
13 as well as its usefulness, we conduct an exemplary large-scale study evaluating
14 13 optimizers from 6 optimization tools. We provide *HPOBench* here: <https://github.com/automl/HPOBench>.
15

16 1 Introduction

17 The plethora of design choices in modern machine learning (ML) makes research on practical and
18 effective methods for hyperparameter optimization (HPO) ever more important. In particular, ever-
19 growing models and datasets create a demand for new HPO methods that are more efficient and
20 powerful than existing black-box optimization (BBO) methods. Especially if it is only feasible
21 to evaluate very few models fully, multi-fidelity optimization methods have been shown to yield
22 impressive results by trading off cheap-to-evaluate proxies and expensive evaluations on the real
23 target [1–5]. They showed tremendous speedups, such as accelerating the search process in low-
24 dimensional ML hyperparameter spaces by a factor of 10 to 1000 [2, 5]. However, the development
25 of such methods often happens in isolation, which potentially prevents HPO research from reaching
26 its full potential. Prior publications on new HPO methods (i) often relied on artificial test functions
27 and low-dimensional toy problems, (ii) sometimes introduced a new set of problems, (iii) set up
28 on different computing environments, having different requirements and interfaces, and (iv) often
29 did not open-source their code base. All of these make it difficult to compare and develop methods,
30 necessitating an evolving set of relevant and up-to-date benchmark problems which drives continued
31 and quantifiable progress in the community.

32 While there are efforts to simplify benchmarking HPO and global optimization algorithms [6–12],
33 we are not aware of efforts to collect a diverse set of benchmarks in a single library, with a unified

*{eggensp, mallik, fh}@cs.uni-freiburg.de

†work done prior to joining Amazon

34 interface and countering potentially conflicting dependencies that may arise over time. The latter
35 is particularly important because the rapid evolution of the Python-ML ecosystem can render a
36 benchmark no longer usable for the community after a major release was published. This creates
37 a significant hurdle for contribution from the community to grow a benchmark library. To solve
38 this issue, we propose *HPOBench*, a benchmark suite for HPO problems, with a special focus
39 on multi-fidelity problems, licensed under a permissive OSS license (*Apache 2.0*) and available
40 at <https://github.com/automl/HPOBench>. *HPOBench* provides a common interface and an
41 infrastructure to isolate benchmarks in their own containers and implements 12 popular optimization
42 families, each with multiple problems and preserved with its dependencies in a container for long-term
43 use. To enable efficient comparisons, most of these benchmarks are table- or surrogate-based, enabling
44 resource efficient large-scale experiments, which we demonstrate in this work. Our contributions are:

- 45 1. The first available collection of multi-fidelity HPO problems. It contains 12 benchmark
46 families with 100+ multi-fidelity HPO problems under a unified interface, comprising
47 traditional HPO and neural architecture search (NAS). These benchmarks also define the
48 largest collection of black-box HPO problems to date.
- 49 2. The first collection of *containerized* benchmarks to ensure the longevity, maintainability
50 and extensibility of benchmarks.
- 51 3. The first set of HPO benchmarks that are available as both, the *raw* benchmark and the
52 *tabular* version.
- 53 4. The first HPO benchmark that also supports multi-objective optimization and transfer-HPO
54 across datasets (and arbitrary combinations of these with multiple fidelities).
- 55 5. We demonstrate how *HPOBench* can be used in an exemplary large-scale study with 13
56 optimizers from 6 optimization tools, assessing whether advanced methods outperform
57 random search and how effective multi-fidelity HPO is.

58 This paper is structured as follows. We first discuss background on HPO and multi-fidelity optimiza-
59 tion (Section 2). Then, we discuss related work on benchmarking (Section 3). Next, we describe
60 the challenges for an HPO benchmark and how *HPOBench* alleviates them (Section 4). Then, we
61 conduct a large-scale comparison of existing, popular HPO methods to demonstrate the usefulness of
62 *HPOBench* (in Section 5). We conclude the paper by highlighting further advantages and potential
63 future work (Section 6).

64 2 Background on Hyperparameter Optimization

65 With *HPOBench* we aim to provide benchmarks to evaluate HPO methods. In the following, we
66 briefly formalize BBO for HPO and survey multi-fidelity optimization (see Feurer and Hutter [13] for
67 a detailed overview), both with a focus on the methods used in our experiments.

68 2.1 Black-box Hyperparameter Optimization

69 Black-box optimization (BBO) aims to find a solution $\arg \min_{\lambda \in \Lambda} f(\lambda)$ where f is a black-box
70 function, for which typically no gradients are available, we cannot make any statements about its
71 smoothness, convexity and noise level. In summary, the only mode of interaction with black-box
72 functions is querying them at given inputs λ and measuring the quantity of interest $f(\lambda)$. In the
73 context of HPO, $\lambda \in \Lambda$ is a hyperparameter configuration where the domain Λ_i of a hyperparameter
74 is often bounded and continuous, but can also be integer, ordinal or categorical. There are also
75 so-called conditional hyperparameters [14, 15] defining hierarchical search spaces; however, the first
76 version of *HPOBench* focuses on flat configuration spaces first as all optimizers support this.

77 There are 3 broad families of BBO methods: (i) purely explorative approaches such as Random
78 Search (*RS*) and grid search are simple but sample-inefficient (ii) model-free Evolutionary Algorithms
79 (*EAs*) based on mutation, crossover and selection operators applied to a population of configurations
80 require comparably large resources to evaluate the entire population but can perform very well
81 given enough resources; (iii) iterative model-based methods, such as Bayesian Optimization [16],
82 which are guided by a predictive model trained on prior function evaluations are known as the most
83 sample-efficient methods. We include representative algorithms from each of these 3 families in our
84 exemplary experiments in Section 5.

85 2.2 Multi-fidelity Hyperparameter Optimization

86 To efficiently optimize today’s ever-growing ML models, multi-fidelity approaches relax the black-box
87 assumption by allowing cheaper queries at lower fidelities b as well ($\arg \min_{\lambda \in \Lambda} f(\lambda, b)$). Examples
88 for these approximations include dataset subsets [2, 17, 18], feature subsets [19] or lower number of
89 epochs [19–21]. Multi-fidelity methods have been shown to lead to speedups of up to $1000\times$ over
90 black-box methods [2, 5]. *HPOBench* will allow the community to compare different multi-fidelity
91 methods and in the following we give an overview of representative methods.

92 A popular multi-fidelity HPO approach that discretizes the fidelity space is Hyperband (*HB* [19]),
93 a very simple method with strong empirical performance. It randomly samples new configurations
94 and allocates more resources to promising configurations by repeatedly calling successive halving
95 (*SH* [4]) as a sub-algorithm. The simplicity and effectiveness of *HB* have been leveraged with
96 other popular black-box optimizers for improved performance: *BOHB* [22] combines *HB* with
97 Bayesian Optimization (*BO*) and *DEHB* [5] combines it with the evolutionary approach of Dif-
98 ferential Evolution (*DE* [23, 24]). The non-*HB*-based multi-fidelity case has also been researched
99 extensively [2, 3, 18, 20, 21, 25–28]. Not being limited to predefined fidelity values makes these
100 methods very powerful, but they rely on strong models to avoid poor choices of fidelities, often
101 making *HB*-based fidelity selection more robust. To study the efficacy of multi-fidelity optimization,
102 in our exemplary experiments in Section 5, we primarily compared black-box optimizers against their
103 multi-fidelity versions (i.e., *RS* vs. *HB*, *BO* vs. *BOHB*, and *DE* vs. *DEHB*). These experiments show
104 large speedups of multi-fidelity optimizers in the regime of small compute budgets, whereas for large
105 compute budgets multi-fidelity optimization is less useful.

106 Besides multi-fidelity optimization, a very active field of study to speed up HPO is to use transfer-
107 learning across datasets [29–33]; we note that transfer HPO methods can also be evaluated with
108 *HPOBench* by learning across the datasets within each of its families.

109 3 Related Work

110 Proper benchmarking is hard. It is important to be aware of technical and methodological pitfalls,
111 e.g. comparing implementations instead of algorithms [34, 35], comparing tuned algorithms versus
112 untuned baselines [36, 37], to not fall for an illusion of progress [38, 39] and to know which sources
113 of variance exist and control for them [40]. Also, there is a rich literature on how to empirically
114 evaluate and compare methods in various domains, e.g. evolutionary optimization [41], planning [42],
115 satisfiability and constraint satisfaction [43], algorithm configuration [44], NAS [45], and also for
116 benchmarking optimization algorithms [46]. Our goal is not to provide further recommendations on
117 how and why to benchmark, but to provide concrete benchmarks to simplify development and to
118 improve the reproducibility and comparability of HPO and in particular multi-fidelity methods.

119 Furthermore, there have been a lot of efforts to provide optimization benchmarks for the community.
120 Having a common set of benchmark problems in a unified format fosters and guides research.
121 Prominent examples in the area of HPO are ACLib [47] for algorithm configuration, COCO [9] for
122 continuous optimization, Bayesmark [8] for Bayesian optimization, Olympus [12] for optimization of
123 experiment planning tasks, and HPO-B [48] for transfer-HPO methods (for more, see Appendix B).
124 However, no benchmark so far has multi-fidelity optimization problems, supports preserving a
125 diverse set of benchmarks for the longer term (containers), supports multiple objectives, and provides
126 cheap-to-evaluate surrogate/tabular benchmarks; we hope to close this gap with *HPOBench*.

127 Besides benchmarks, competitions are another form of focusing research effort by providing a
128 common goal and incentive. Famous examples are the *AutoML* challenges [49], the *AutoDL* chal-
129 lenge [50], the GECCO BBOB workshop series based on COCO [9] and the NeurIPS 2020 BBO
130 challenge [51] (for more, see Appendix C). In contrast to these, we do not focus on defining concrete
131 experimentation protocols, but rather on providing a flexible benchmarking environment to study,
132 develop and compare optimization methods.

133 4 HPOBench: A Benchmark Suite for Multi-Fidelity Hyperparameter 134 Optimization benchmarks

135 In this section, we present *HPOBench*, a collection of HPO benchmarks defined as follows:

136 **Definition 1 (HPO Benchmark)** *An HPO benchmark consists of a function $f : \lambda \rightarrow \mathcal{R}$ to be*
137 *minimized and a (bounded) hyperparameter space Λ with hyperparameters $[\Lambda_1, \dots, \Lambda_d]$ of type*
138 *continuous, integer, categorical or ordinal. In the case of multi-fidelity benchmarks, f can be*
139 *queried at lower fidelities, $f : \lambda \times \mathbf{b} \rightarrow \mathcal{R}$, and the fidelity space \mathcal{B} describes which low-fidelities*
140 *$[B_1, \dots, B_e]$ of type continuous, integer or ordinal are available.*

141 Specifically, each benchmark consists of the implementation of that function, which returns at least
142 one loss. Since this function typically evaluates an ML algorithm, the benchmark defines all relevant
143 settings, dependencies and inputs, such as datasets, splits and how to compute the loss.

144 In the remainder of this section, we first discuss the desiderata of a benchmark that aids HPO research
145 and then highlight the features of *HPOBench* by detailing how its design fulfills these desiderata.

146 4.1 Desiderata of an HPO Benchmark

147 One of the challenges posed to standardized HPO research lies in the varied choices of the underlying
148 ML components – datasets and their splits, preprocessing, hyperparameter ranges, underlying software
149 versions, and hardware used. Moreover, the practices applied in HPO research itself can vary along
150 the lines of optimization budget, number of repetitions, metrics measured and reported. This leads to
151 inconsistencies and difficulties in comparison of different HPO methods across publications and over
152 time, affecting the reproducibility of experiments that hinders continued progress in HPO research.

153 In order to alleviate such issues and encourage participation by the research community, benchmarks
154 need to standardize these practices to allow the community to be an active stakeholder in developing
155 and re-using benchmarks. *HPOBench* is designed to both allow easy, flexible use with a minimal
156 API that is identical for all benchmarks (see Figure 2); and have a low barrier for contributing new
157 benchmark problems. We, therefore, identify 3 features of a benchmark that allow its wide-scale
158 use and long-term applicability: (i) *efficiency* by providing tabular and surrogate benchmarks for
159 quick, efficient experiments, along with the original benchmarks; (ii) *reproducibility* of results by
160 containerizing benchmarks; and (iii) *flexibility* by covering different optimization landscapes and
161 possible use cases, e.g. multi-objective, transfer-HPO, and even multi-fidelity optimization with
162 multiple fidelity variables. To our knowledge, no other existing benchmarks offer these possibilities.
163 *HPOBench* provides a framework to enable standardized, principled research and experimentation.
164 We list all benchmarks that are included in *HPOBench* in Table 1 and provide a detailed description of
165 the respective configuration spaces in Appendix D.

166 4.2 Efficiency

167 HPO benchmarks that follow Definition 1 exhibit the drawback that they evaluate a costly func-
168 tion, rendering the empirical comparison of optimization algorithms expensive and ruling out such
169 benchmarks for interactive development of new methods. To overcome this issue, beside such raw
170 benchmarks, we also provide two well-established benchmark classes which alleviate this issue:

171 **Definition 2 (Tabular Benchmark)** *A tabular benchmark returns values from a lookup table with*
172 *recorded function values of a raw HPO benchmark instead of evaluating $f(\lambda)$. The (bounded)*
173 *hyperparameter space is restricted to only contain these values and therefore bears a form of*
174 *discretization. In the case of multi-fidelity benchmarks, each tabular benchmark has a fidelity space*
175 *and the underlying table also contains the recorded function values on the low-fidelities.*

176 Tabular benchmarks are popular in the HPO community as they are easy to distribute and induce little
177 overhead [52, 30, 53–55], however, they require to discretize the hyperparameter space. Surrogate
178 benchmarks [56, 57] are an alternative since they provide the original hyperparameter space.

179 **Definition 3 (Surrogate Benchmark)** *A surrogate benchmark returns function values predicted by*
180 *an ML model trained on a tabular benchmark or recorded function values of a raw HPO benchmark.*
181 *It reuses the original hyperparameter space and can be extended to the multi-fidelity case as well.*

182 While surrogate benchmarks are similarly cheap to query, the surrogate’s internal ML model adds
183 extra complexity and the benchmark’s quality crucially depends on the quality of this model and its
184 training data. Because surrogate benchmarks yield a drop-in replacement for raw benchmarks, they
185 enjoy widespread adoption in the HPO community [22, 58, 57, 59–62].

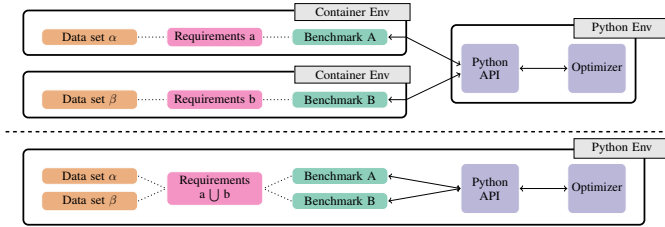


Figure 1: Overview of benchmark environments with (upper) and without (lower) using containers.

```

from hpoBench.container_benchmarks.\
nas.nasbench_101 import NASCIfor10ABenchmark

# if necessary downloads container/data
b = NASCIfor10ABenchmark(rng=1)
# get hyperparameter space
space = b.get_configuration_space(seed=1)
# sample config at random
config = space.sample_configuration()

# eval at multiple low-fidelities
res = b.objective_function(
    configuration=config,
    fidelity=("budget": 12), rng=1)
res = b.objective_function(
    configuration=config,
    fidelity=("budget": 108), rng=1)

```

Figure 2: Code example initializing and evaluating a benchmark.

186 Furthermore, while *HPOBench* puts a strong focus on multi-fidelity benchmarks, it also facilitates
 187 evaluating black-box optimization algorithms. In fact, a multi-fidelity benchmark with k different
 188 fidelity levels can be used to define k separate (yet related) benchmarks for black-box optimization.
 189 As such, *HPOBench* defines more than 400 black-box HPO benchmarks.

190 4.3 Reproducibility

191 One of the challenges that come with many new benchmarks is their one-off development and their
 192 lack of maintenance. This means that any new update to the benchmark or its dependencies can easily
 193 lead to conflicts and inconsistencies with respect to software dependencies and possibly old published
 194 results (see Appendix D.1 for examples). While in practice the very same problem, also known as
 195 *dependency hell*, can also occur on the optimizer side, in this paper we focus on the benchmark side.

196 *HPOBench* circumvents such issues through the containerization of benchmarks using Singularity [63]
 197 containers.³ Each benchmark and its dependencies are packaged as a separate container, which isolates
 198 benchmarks from each other and also from the host system. Figure 1 illustrates the advantages that
 199 containerization provides, especially when running multiple benchmarks in the same environment.
 200 Note that without containers, the environment needs to satisfy the union of all of its benchmarks’
 201 requirements (which may actually be mutually exclusive!), while with containers the dependencies
 202 for any given benchmark only need to be satisfied once: for the creation of the container. Importantly,
 203 the dependencies do not need to be satisfied again for using the benchmarks. Each benchmark is
 204 uploaded as a container to a GitLab container registry to provide the history of different versions
 205 of the benchmark. Hence, any benchmark created under the *HPOBench* paradigm remains usable
 206 without additional bookkeeping or installation overheads for long-term usage. Additionally, no effort
 207 is required for maintaining already containerized benchmarks, as long as the API does not change.
 208 Although not recommended, each benchmark can also be installed locally along with its specific
 209 dependencies without using the containers, since the interfaces of the local and the containerized
 210 version are identical and thus interchangeable. We provide a short code sample in Figure 2.

211 4.4 Flexibility

212 *HPOBench* is a flexible framework that can be used to validate existing HPO research, and develop
 213 and improve HPO algorithms, with a focus on multi-fidelity methods. It consists of two sets of
 214 benchmarks, which we describe in turn: 22 existing multi-fidelity benchmarks from 7 families that
 215 we collected from the multi-fidelity literature (Section 4.4.1); and 88 new benchmarks from 5 families
 216 we created to allow a much more flexible use of *HPOBench* (Section 4.4.2).

217 4.4.1 Existing Community Benchmarks

218 Firstly, to allow comparability with previous experiments, we collected 22 existing multi-fidelity
 219 benchmarks from 7 families from the multi-fidelity literature; *HPOBench* preserves these benchmarks
 220 by containerizing them and encapsulating them all under a common API (which was not the case
 221 before). This not only ensures important previous work to remain accessible publicly under the same
 222 umbrella, but it also bypasses dependency issues that have appeared in these benchmarks enabling
 223 long term usage (see Appendix D.1).

³We chose Singularity over the popular Docker (<https://www.docker.com/>) alternative as it (1) does not require super user access and (2) is available on the computer clusters we have access to.

Table 1: Overview of raw (✓), surrogate (✗) and tabular (✓) benchmarks. We report the number of benchmarks per family (*#benchs*), the number of continuous (*#cont*), integer (*#int*), categorical (*#cat*), ordinal (*#ord*) hyperparameters and how many are log-scaled. Furthermore, we report the fidelity, the optimization budget and the number of configurations for tabular and surrogate benchmarks.

Family	#benchs	#cont(log)	#int(log)	#cat	#ord	fidelity	type	opt. budget	#confs	Ref.
<i>Cartpole</i>	1	4(1)	3(3)	-	-	repetitions	✓	1d	-	[22]
<i>BNN</i>	2	3(1)	2(2)	-	-	samples	✓	1d	-	[22]
<i>Net</i>	6	5	1	-	-	time	✗	7d	-	[22]
<i>NBHPO</i>	4	-	-	3	6	epochs	(✓)	10 ⁷ sec	62 208	[64]
<i>NB101</i>	3	-	-	26	-	epochs	(✓)	10 ⁷ sec	423k	[54]
		21	1	5	-					
<i>NB201</i>	3	-	-	6	-	epochs	(✓)	10 ⁷ sec	15 625	[65]
<i>NB1Shot1</i>	3	-	-	9	-	epochs	(✓)	10 ⁷ sec	6 240	[66]
		-	-	9	-				29 160	
		-	-	11	-				363 648	
<i>LogReg</i>	20	2(2)	-	-	-	iter	✓, (✓)	100×	625	<i>new</i>
<i>SVM</i>	20	2(2)	-	-	-	data	✓, (✓)	average	441	<i>new</i>
<i>RandomForest</i>	20	1	3(2)	-	-	#trees	✓, (✓)	runtime on	10k	<i>new</i>
<i>XGBoost</i>	20	3(2)	1(1)	-	-	#trees	✓, (✓)	the highest	10k	<i>new</i>
<i>MLP</i>	8	2(2)	3(2)	-	-	epochs	✓, (✓)	fidelity	30k	<i>new</i>

224 Specifically, these benchmarks comprise raw benchmarks tuning a reinforcement learning agent
 225 (PPO on *Cartpole* [22]) and a Bayesian neural network (*BNN* [22]), a random forest-based surrogate
 226 benchmark tuning an MLP (*Net* [22]) and four popular NAS benchmark families (*NBHPO* [64],
 227 *NB101* [54], *NB201* [65], and *NB1Shot1* [66]). However, these existing community benchmarks also
 228 have certain limitations: they are only of limited use for transfer HPO (since there are only between 1
 229 and 6 benchmarks per family), they only offer a single fidelity dimension, and they only evaluate a
 230 single metric. We therefore augmented them with 5 new families of benchmarks we describe next.

231 4.4.2 New Benchmarks

232 To substantially increase the range of possible applications of *HPOBench*, we defined 5 new bench-
 233 mark families with up to 20 different datasets per family, comprising a total of 88 new multi-fidelity
 234 benchmarks. These new benchmarks also provide multiple metrics and multiple fidelity dimensions
 235 to go beyond the aforementioned limitations of the community benchmarks.

236 Our new benchmarks are based on the following popular ML algorithms: *SVM*, *LogReg*, *XGBoost*,
 237 *RandomForest*, and *MLP*. All of them evaluate the respective ML algorithm as implemented in
 238 scikit-learn [67] and XGBoost [68] on 20 publicly available datasets (8 for the *MLP* due to its
 239 high computational cost) from the OpenML AutoML benchmark [69]. We give the OpenML [70]
 240 task IDs in Table 7 in Appendix D, which provide fixed train-test splits; for each such task, we
 241 used 33% of the training set as the validation split, determined through stratified sampling under a
 242 fixed seed. The entire objective function then consists of preprocessing, training the model on the
 243 remaining 66% of the fixed OpenML training split, prediction on the fixed validation split, evaluating
 244 4 different metrics, and recording model fit and inference times. The fidelities are algorithm-specific
 245 if possible (number of trees, iterations, epochs) or dataset subsets otherwise (which is used for *SVM*).
 246 These benchmarks are available both as raw and tabular versions, have the same API and exist in
 247 independent, non-conflicting containers; for the tabular versions, we discretized each hyperparameter
 248 (and fidelity) and evaluated 5 different seeds for each configuration of the resulting grid.

249 Also, four of our new benchmark families (*LogReg*, *RandomForest*, *XGBoost*, *MLP*) allow up to
 250 two fidelity dimensions. This enables the development and benchmarking of methods for multi-
 251 fidelity optimization with multiple fidelity dimensions, a direction that we deem very promising yet
 252 understudied. Similarly, our tabular data collected over multiple datasets (up to 20) allows the effective
 253 use of these benchmarks for transfer-HPO, and the recording of multiple evaluation metrics also

254 allows these benchmarks to be used for multi-objective optimization. Moreover, each configuration
 255 is recorded on different fidelities with their associated costs, which further lends *HPOBench* great
 256 potential in future research in cost-based meta-learning or multi-fidelity multi-objective optimization.

257 To demonstrate the diversity of our new benchmarks, we show the *empirical cumulative distribution*
 258 *function* (ECDF) for each family in Figure 3. Each line corresponds to one dataset and shows how the
 259 objective values are distributed. From the varying amounts of well and badly performing normalized
 260 regrets we can conclude that the benchmarks yield different landscapes and thus are diverse in
 261 smoothness, which results in varying algorithm performance. Moreover, the 5 new spaces vary
 262 in their dimensionality (up to 5 for *MLP*), in the hyperparameter data types and their range (see
 263 Appendix D).

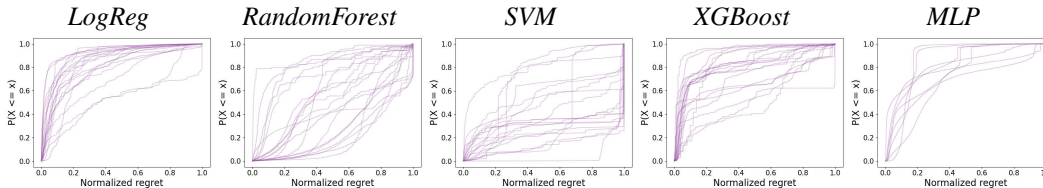


Figure 3: Empirical cumulative distribution. Each plot corresponds to one ML algorithm, and each line within a plot corresponds to one dataset. The lines show the ECDF of the normalized regret of all evaluated configurations of the respective ML algorithm on the respective dataset.

264 5 Experiments

265 Now, we turn to an exemplary use of our benchmarks in order to demonstrate some features of
 266 *HPOBench* and its utility for HPO research. We used our benchmark suite to run a large-scale
 267 empirical study comparing 13 optimization methods on our 12 benchmark families (we report
 268 detailed results in Appendix H). We first give details on the experimental setup and then study the
 269 following two exemplary research questions: **(RQ1)** *Do advanced methods improve over random*
 270 *baselines?* and **(RQ2)** *Do multi-fidelity methods improve over single-fidelity methods?*

271 5.1 Experimental Setup

272 For each benchmark and optimizer, we conducted 32 repetitions with a different seed each. For our
 273 new benchmarks, which have multiple metrics, we minimized 1–accuracy. For each run, we allowed
 274 an optimization budget as described in Table 1 and accumulate time taken by the benchmark (recorded
 275 time for tabular benchmarks, predicted time for surrogate benchmarks and wallclock time for raw
 276 benchmarks; for our new benchmarks, we used the tabular versions to avoid unnecessary compute
 277 costs and CO2 exhaustion) and the optimizer (wallclock time). We kept track of all evaluations and
 278 computed trajectories, i.e., the best-seen value at each time step, as follows: If for an evaluation we
 279 cannot find another evaluation conducted on the same or a higher fidelity, we treat it as the best-seen
 280 value; if it is on the highest fidelity evaluated so far, we treat it as the best seen value if it has a lower
 281 loss than the best-seen so far on that fidelity; otherwise, we do not consider this evaluation for the
 282 trajectory. This decision reflects the multi-fidelity setting, where a higher budget results in a better
 283 estimate of the actual value of interest but can cause jumps in the optimization trajectory, (e.g., when a
 284 configuration is the first to be evaluated on a higher budget but is worse than the best configuration on
 285 a lower budget). To aggregate and report results, we use either the *final performance* (per benchmark,
 286 see Appendix H), *performance-over-time* (per benchmark, see Appendix H) or *rank-over-time* (across
 287 multiple benchmarks). For tabular and surrogate benchmarks we report optimization regret (the
 288 difference between the best-found value and the best-known value) and for the other benchmarks, we
 289 report the actual optimized objective value.⁴

290 We give details on the hardware and required compute resources in Appendix E and F and release
 291 code for the experiments here: <https://github.com/automl/HPOBenchExperimentUtils>.

⁴Since we study optimizers, we report optimization performance (in the case of ML the validation performance, which is the objective value seen by the optimizer. We note that *HPOBench* in principle allows to compute test performance (the loss computed on a separate test set on the highest fidelity).

292 5.2 Considered Optimizers

293 We evaluated a wide set of optimizers including baselines for black-box and multi-fidelity optimization
294 which we briefly describe (for more details see Appendix G). Our selection of optimization algorithm
295 does *not* aim at finding the best optimization algorithm, but to study a broad range of different
296 implementations and tools. While there are more recent and efficient BO and black-box methods
297 available [71, 72], we focus here on multi-fidelity methods and only include basic black-box methods
298 for which there are easy-to-use tools available.

299 As black-box optimizers, which only access the highest fidelity, we considered random search (*RS*),
300 differential evolution (*DE* [23, 24]) and *BO* with different models: a Gaussian Process model (*BO_{GP}*),
301 a random forest (*BO_{RF}* [73]), a kernel density estimator (KDE) (*BO_{KDE}* [22]) and a tree Parzen
302 estimator (TPE) (*Ray_{hyp}* [14]). For multi-fidelity optimization, we used the multi-fidelity extensions
303 of the methods mentioned above: Hyperband (*HB* [19]) and its combination with KDE-based *BO*
304 (*BOHB* [22]), with RF-based *BO* (*SMAC-HB*) and with *DE* (*DEHB* [5]). Additionally, we use
305 Dragonfly [74] using a GP with multi-fidelity optimization. Furthermore, we ran combinations of
306 optimization and multi-fidelity algorithms implemented in Ray [75] and Optuna [76].⁵

307 5.3 RQ1: Do advanced methods improve over random search?

308 To demonstrate the validity of our benchmarks, we independently replicate the findings of the 1st
309 NeurIPS Blackbox Optimization challenge [51]: “*decisively showing that BO and similar methods*
310 *are superior choices over RS and grid search for tuning hyperparameters of ML models*”. While
311 this question has already been studied before [14, 77, 36, 15, 33, 72, 78], we will also study it
312 w.r.t. multi-fidelity optimization and using the popular *HB* baseline. We leave out grid search as *RS*
313 has been shown to be superior [77] and as there is no multi-fidelity version of grid search.

314 We report ranks-over-time in Figure 4, comparing black-box (*DE*, *BO_{GP}*, *BO_{RF}*, *Ray_{hyp}*, *BO_{KDE}*;
315 1st column) and multi-fidelity (*BOHB*, *DEHB*, *SMAC-HB*, *DF*, *Ray_{hyp}^{asha}*; 2nd column) optimizers
316 on *existing community* (top row) and *new* (bottom row) benchmarks. On both benchmark sets most
317 black-box and multi-fidelity optimizers clearly outperform the respective baseline (*RS* (blue) and *HB*
318 (light green)) on average. We also observe that *BO* improves over the evolutionary algorithm *DE* in
319 the beginning but loses to it in the very long run on the *existing community* benchmarks [79, 60, 5].
320 This does not happen on the *new* benchmarks, as their time limits are set more aggressively and
321 the methods developed for this setting (*BO_{GP}*, *BO_{RF}* and *SMAC-HB* [80]) achieve lower ranks.
322 Considering per-benchmark results (Appendix H), we also observe that methods which appear clearly
323 inferior in the ranking plots perform very well on individual benchmarks (e.g. *DF*⁶ on *NB201* and
324 *Ray_{hyp}^{asha}* on *Net_{Adult}*).

325 Besides qualitative measures, we also quantitatively measure whether the advanced methods outper-
326 form the respective baselines by counting the number of wins, ties and losses and using the sign test
327 to verify significance [81] on the existing community benchmarks in Table 2 (the new benchmarks
328 yield similar results; see Appendix H). We can observe that four out of five black-box methods are
329 significantly better than *RS*. In the multi-fidelity case, only two out of five methods were significantly
330 better than *HB*, one method was insignificantly better than *HB* and two methods were consistently
331 worse than *HB*. Overall, we conclude that advanced methods consistently outperform random search.

332 5.4 RQ2: Do multi-fidelity methods improve over black-box methods?

333 Next, we study whether multi-fidelity optimization methods are able to consistently improve over
334 black-box optimization methods given a fixed time budget. For this, we look again at ranking-over-
335 time in Figure 4. We first compare black-box methods with their respective multi-fidelity extension,
336 i.e., *DE* vs. *DEHB* and *BO_{KDE}* vs. *BOHB* in the two plots in the rightmost column. We can see
337 that in the beginning *HB* and the multi-fidelity optimizers perform very similarly and consistently

⁵We include these frameworks to show compatibility of *HPOBench* with popular frameworks, but note that these expect to freeze and thaw evaluations. *HPOBench* implements a stateless objective function and, thus, runs that could be thawed and continued instead get accounted the full costs of rerunning them, which slows down optimization. We defer stateful benchmarks to future work.

⁶We would like to note that the bad rank of *DF* for some benchmarks is due to its overhead which prevented it from spending sufficient budget on function evaluations; see Section F for details.

Table 2: P-value of a sign test for the hypothesis that advanced methods outperform the baseline *RS* for black-box optimization and *HB* for multi-fidelity optimization. We underline p-values that are below $\alpha = 0.05$ and also boldface p-values that are below $\alpha = 0.05$ after multiple comparison correction (dividing α by the number of comparisons, i.e. 5; boldface/underlined implies that the advanced method is better). We also give the wins/ties/losses of *RS* and *HB* against the challengers.

	<i>DE</i>	<i>BO_{GP}</i>	<i>BO_{RF}</i>	<i>Ray_{hyp}</i>	<i>BO_{KDE}</i>
p-value against <i>RS</i>	<u>0.00006</u>	<u>0.00360</u>	<u>0.00001</u>	<u>0.00043</u>	0.41591
wins/ties/losses of <i>RS</i>	1/2/19	3/3/16	0/3/19	2/2/18	9/2/11
	<i>BOHB</i>	<i>DEHB</i>	<i>SMAC-HB</i>	<i>DF</i>	<i>Ray_{hyp}^{asha}</i>
p-value against <i>HB</i>	0.26173	<u>0.00001</u>	<u>0.00011</u>	0.99640	0.97376
wins/ties/losses of <i>HB</i>	8/2/12	0/3/19	0/5/17	16/1/5	15/0/7

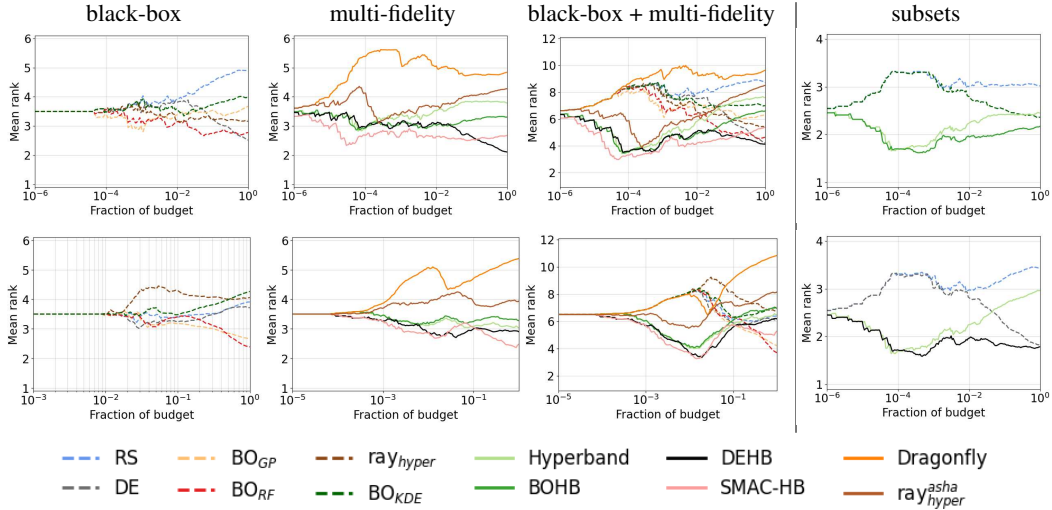


Figure 4: Mean rank-over-time across 32 repetitions of different sets of optimizers (lower is better). The left part shows rank across all *existing community* benchmarks (upper row) and *new* benchmarks (lower row). The right part report results on the *existing community* benchmarks only for subsets of optimizers.

338 outperform *RS* and the respective black-box version. After a while, the multi-fidelity versions improve
 339 over the *HB* baseline, and given enough time, the black-box versions catch up. Second, we compare
 340 all optimizers on the *existing community* (3rd column, top) and *new* (3rd column, bottom) benchmarks.
 341 Here, we can observe a similar pattern in that *HB* is a very competitive baseline in the beginning but
 342 is outperformed first by the advanced multi-fidelity methods and then also by the black-box methods.
 343 This is less pronounced on the *new* benchmarks, which we attribute to the tighter time limits.

344 Similarly to RQ1, we again counted the wins, ties and losses and used the sign test to verify
 345 significance [81] on the existing community benchmarks for 100%, 10% and 1% of the total budget
 346 in Table 3 (the new benchmarks yield similar results; see Appendix H). We can observe that only *HB*
 347 is able to outperform its black-box counterpart for all three budgets we check. For the other three
 348 multi-fidelity methods there is a significant improvement over the black-box methods for 1% of the
 349 budget. For 10% only *BOHB* is still significantly better, while *DEHB* is only insignificantly better
 350 and *SMAC-HB* ties. For the full budget we can no longer state that any of the multi-fidelity methods
 351 is statistically better than their counterpart, but judging by the wins and losses the multi-fidelity
 352 methods are still competitive.

353 Overall, multi-fidelity optimizers outperform black-box optimizers for relatively small compute
 354 budgets. Given enough budget, black-box optimizers become competitive with their multi-fidelity
 355 versions; in particular, *DE* and *BO_{RF}* performed very well in the end. However, we need to take into
 356 account that for the *existing community* benchmarks the potential catch-up (if at all) only happens after

357 a very substantial amount of (simulated) wallclock time (e.g., 10 Mio. seconds). Hence, multi-fidelity
 358 methods are crucial to efficiently tackle real, expensive optimization problems.

Table 3: P-values of a sign test for the hypothesis that multi-fidelity outperform their black-box counterparts. We boldface p-values that are below $\alpha = 0.05$ (boldface implies that the multi-fidelity method is better).

Budget		<i>RS</i> vs <i>HB</i>	<i>DE</i> vs <i>DEHB</i>	<i>BO_{KDE}</i> vs <i>BOHB</i>	<i>BO_{RF}</i> vs <i>SMAC-HB</i>
100%	p-values	0.00217	0.58409	0.14314	0.85686
	w/t/l	2/4/16	7/8/7	7/2/13	10/6/6
10%	p-values	0.00360	0.06690	0.01330	0.58409
	w/t/l	3/3/16	3/8/11	4/3/15	8/6/8
1%	p-values	0.00074	0.00845	0.02624	0.03918
	w/t/l	2/3/17	4/2/16	4/4/14	5/3/14

358

359 To conclude, in general when low-fidelities are available and they are representative of the true
 360 objective function, multi-fidelity methods are clearly beneficial. In practice, we found that *DEHB*
 361 and *SMAC-HB* are reliable multi-fidelity optimizers that work well across the whole collection of
 362 benchmarks, while other multi-fidelity optimizers are not able to improve over *HB* consistently. By
 363 exploring a very broad range of benchmarks, we also found an existence proof that black-box methods
 364 can outperform multi-fidelity methods for very high budgets and that even advanced methods can
 365 be outperformed by *RS* in individual benchmarks. We pose it as a challenge to the field to develop
 366 methods that do not exhibit poor performance in *any* of the many benchmarks in *HPOBench*.

367 6 Discussion and Future Work

368 We proposed *HPOBench*, a library for multi-fidelity HPO benchmarks. It serves two purposes: (a) to
 369 provide benchmarks with a unified API, and (b) to make them easy to install and use by containerizing
 370 them and thus enable rapid prototyping and the development of new multi-fidelity methods that
 371 are crucial for ML research and applications. Finally, our library is open-source and we welcome
 372 contributions of new benchmarks to keep the library up-to-date and evolve it.

373 On the technical side, we so far focused on developing a benchmark library, but we see a large
 374 potential in connecting our library with other benchmarking frameworks (e.g. COCO [9] and
 375 Bayesmark [8]), optimization frameworks (e.g. Nevergrad [10] and Sherpa [82]) and extending it
 376 with further benchmarks [11, 60, 83, 62, 12] to increase diversity and to simplify evaluation and
 377 comparison of optimizers. For this it would be interesting to also containerize the optimizers since
 378 they can suffer from the same issues as benchmarks. Furthermore, so far *HPOBench* only contains
 379 stateless benchmarks starting a single container. We would like to extend the library to also support
 380 optimizers requiring stateful benchmarks (to freeze and thaw evaluations) or running in parallel.

381 Our set of benchmarks already covers raw, tabular, and surrogate benchmarks, but it would be
 382 useful to have all three versions available for all benchmarks, and to automatically generate tabular
 383 and surrogate-based benchmarks from raw benchmarks. Also, our new benchmarks can be used
 384 to evaluate multi-objective (multiple metrics) and meta-learning (across datasets) methods or even
 385 meta-learned multi-fidelity multi-objective methods. We hope for the community to play a large
 386 role in defining the protocols for the different special cases; e.g., budgets need to be set differently
 387 for black-box multi-objective optimization and single-objective hyperparameter transfer learning.
 388 Additionally, it would be interesting to study hierarchical search spaces to cover work in *AutoML*.
 389 Furthermore, there is a large potential in automatically creating multi-fidelity benchmarks from any
 390 ML algorithm by using data subsets as a low-fidelity.

391 We also conducted a large-scale study evaluating 13 algorithm implementations to demonstrate
 392 compatibility with a wide range of optimization tools, and we thus believe that our library is well
 393 suited for future research on multi-fidelity optimization. We showed that advanced HPO methods are
 394 preferable over *RS* and *HB* baselines, and that multi-fidelity extensions of popular optimizers improve
 395 over their black-box version. Lastly, to reduce computational effort, we would like to study whether
 396 we can learn which benchmarks are hard and whether it suffices by executing only a representative
 397 subset of the benchmarks [84].

398 **References**

- 399 [1] A. Forrester, A. Sobester, and A. Keane. Multi-fidelity optimization via surrogate modelling.
400 *Proceedings of The Royal Society A Mathematical Physical and Engineering Sciences*, 463:
401 3251–3269, 2007.
- 402 [2] A. Klein, S. Falkner, S. Bartels, P. Hennig, and F. Hutter. Fast Bayesian optimization of machine
403 learning hyperparameters on large datasets. In A. Singh and J. Zhu, editors, *Proceedings of*
404 *the Seventeenth International Conference on Artificial Intelligence and Statistics (AISTATS)*,
405 volume 54. Proceedings of Machine Learning Research, 2017.
- 406 [3] K. Kandasamy, G. Dasarathy, J. Schneider, and B. Póczos. Multi-fidelity Bayesian Optimisation
407 with Continuous Approximations. In D. Precup and Y. Teh, editors, *Proceedings of the 34th*
408 *International Conference on Machine Learning (ICML'17)*, volume 70, pages 1799–1808.
409 Proceedings of Machine Learning Research, 2017.
- 410 [4] K. Jamieson and A. Talwalkar. Non-stochastic best arm identification and hyperparameter
411 optimization. In A. Gretton and C. Robert, editors, *Proceedings of the Seventeenth International*
412 *Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 51. Proceedings of
413 Machine Learning Research, 2016.
- 414 [5] N. Awad, N. Mallik, and F. Hutter. DEHB: Evolutionary hyperband for scalable, robust and
415 efficient hyperparameter optimization. In Z.-H. Zhou, editor, *Proceedings of the Thirtieth Inter-*
416 *national Joint Conference on Artificial Intelligence, IJCAI-21*, pages 2147–2153. International
417 Joint Conferences on Artificial Intelligence Organization, 2021.
- 418 [6] K. Eggensperger, M. Feurer, F. Hutter, J. Bergstra, J. Snoek, H. Hoos, and K. Leyton-Brown.
419 Towards an empirical foundation for assessing Bayesian optimization of hyperparameters. In
420 *NeurIPS Workshop on Bayesian Optimization in Theory and Practice (BayesOpt'13)*, 2013.
- 421 [7] C. Doerr, H. Wang, F. Ye, S. van Rijn, and T. Bäck. Iohprofiler: A benchmarking and profiling
422 tool for iterative optimization heuristics. *arXiv:1810.05281 [cs.NE]*, 2018.
- 423 [8] R. Turner and D. Eriksson. Bayesmark: Benchmark framework to easily compare bayesian
424 optimization methods on real machine learning tasks. github.com/uber/bayesmark, 2019.
- 425 [9] N. Hansen, A. Auger, R. Ros, O. Mersman, T. Tušar, and D. Brockhoff. COCO: A platform for
426 comparing continuous optimizers in a black-box setting. *Optimization Methods and Software*,
427 2020.
- 428 [10] J. Rapin and O. Teytaud. Nevergrad - A gradient-free optimization platform. <https://GitHub.com/FacebookResearch/Nevergrad>, 2018.
- 430 [11] L. Bliiek, A. Guijt, R. Karlsson, S. Verwer, and M. de Weerd. EXPObench: Benchmarking
431 surrogate-based optimisation algorithms on expensive black-box functions. *arXiv:2106.04618*
432 *[cs.LG]*, 2021.
- 433 [12] F. Häse, M. Aldeghi, R. Hickman, L. Roch, M. Christensen, E. Liles, J. Hein, and A. Aspuru-
434 Guzik. Olympus: a benchmarking framework for noisy optimization and experiment planning.
435 *Machine Learning: Science and Technology*, 2(3), 2021.
- 436 [13] M. Feurer and F. Hutter. Hyperparameter optimization. In F. Hutter, L. Kotthoff, and J. Van-
437 schoren, editors, *Automated Machine Learning: Methods, Systems, Challenges*, volume 5 of
438 *The Springer Series on Challenges in Machine Learning*, chapter 1, pages 3–38. Springer, 2019.
439 Available for free at <http://automl.org/book>.
- 440 [14] J. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl. Algorithms for hyper-parameter optimization.
441 In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K. Weinberger, editors, *Proceedings*
442 *of the 24th International Conference on Advances in Neural Information Processing Systems*
443 *(NeurIPS'11)*, pages 2546–2554. Curran Associates, 2011.

- 444 [15] C. Thornton, F. Hutter, H. Hoos, and K. Leyton-Brown. Auto-WEKA: combined selection and
445 hyperparameter optimization of classification algorithms. In I. Dhillon, Y. Koren, R. Ghani,
446 T. Senator, P. Bradley, R. Parekh, J. He, R. Grossman, and R. Uthurusamy, editors, *The 19th*
447 *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'13)*,
448 pages 847–855. ACM Press, 2013.
- 449 [16] B. Shahriari, K. Swersky, Z. Wang, R. Adams, and N. de Freitas. Taking the human out of the
450 loop: A review of Bayesian optimization. *Proceedings of the IEEE*, 104(1):148–175, 2016.
- 451 [17] J. Petrak. Fast subsampling performance estimates for classification algorithm selection. Tech-
452 nical Report TR-2000-07, Austrian Research Institute for Artificial Intelligence, 2000.
- 453 [18] K. Swersky, J. Snoek, and R. Adams. Multi-task Bayesian optimization. In C. Burges, L. Bottou,
454 M. Welling, Z. Ghahramani, and K. Weinberger, editors, *Proceedings of the 26th International*
455 *Conference on Advances in Neural Information Processing Systems (NeurIPS'13)*, pages 2004–
456 2012. Curran Associates, 2013.
- 457 [19] L. Li, K. Jamieson, G. DeSalvo, A. Rostamizadeh, and A. Talwalkar. Hyperband: A novel
458 bandit-based approach to hyperparameter optimization. *Journal of Machine Learning Research*,
459 18(185):1–52, 2018.
- 460 [20] K. Swersky, J. Snoek, and R. Adams. Freeze-thaw Bayesian optimization. *arXiv:1406.3896*
461 *[stats.ML]*, 2014.
- 462 [21] T. Domhan, J. Springenberg, and F. Hutter. Speeding up automatic hyperparameter optimization
463 of deep neural networks by extrapolation of learning curves. In Q. Yang and M. Wooldridge, edi-
464 tors, *Proceedings of the 25th International Joint Conference on Artificial Intelligence (IJCAI'15)*,
465 pages 3460–3468, 2015.
- 466 [22] S. Falkner, A. Klein, and F. Hutter. BOHB: Robust and efficient hyperparameter optimization at
467 scale. In J. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on*
468 *Machine Learning (ICML'18)*, volume 80, pages 1437–1446. Proceedings of Machine Learning
469 Research, 2018.
- 470 [23] R. Storn and K. Price. Differential evolution – a simple and efficient heuristic for global
471 optimization over continuous spaces. *Journal of Global Optimization*, 11:341–359, 1997.
- 472 [24] N. Awad, N. Mallik, and F. Hutter. Differential evolution for neural architecture search. In
473 *Proceedings of the 1st workshop on neural architecture search@ICLR'20*, 2020.
- 474 [25] D. Golovin, B. Solnik, S. Moitra, G. Kochanski, J. Karro, and D. Sculley. Google Vizier: A
475 service for black-box optimization. In S. Matwin, S. Yu, and F. Farooq, editors, *Proceedings of*
476 *the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*
477 *(KDD)*, pages 1487–1495. ACM Press, 2017.
- 478 [26] J. Wu, S. Toscano-Palmerin, P. Frazier, and A. Wilson. Practical multi-fidelity Bayesian
479 optimization for hyperparameter tuning. In R. Adams and V. Gogate, editors, *Proceedings*
480 *of The 35th Uncertainty in Artificial Intelligence Conference (UAI'20)*, volume 115, pages
481 788–798. PMLR, 2020.
- 482 [27] S. Takeno, H. Fukuoka, Y. Tsukada, T. Koyama, M. Shiga, I. Takeuchi, and M. Karasuyama.
483 Multi-fidelity Bayesian optimization with max-value entropy search and its parallelization.
484 In H. Daume III and A. Singh, editors, *Proceedings of the 37th International Conference on*
485 *Machine Learning (ICML'20)*, volume 98, pages 9334–9345. Proceedings of Machine Learning
486 Research, 2020.
- 487 [28] J. Song, Y. Chen, and Y. Yue. A general framework for multi-fidelity Bayesian optimization
488 with Gaussian processes. In K. Chaudhuri and M. Sugiyama, editors, *Proceedings of the 22nd*
489 *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 3158–3167.
490 Proceedings of Machine Learning Research, 2019.

- 491 [29] J. Vanschoren. Meta-learning. In F. Hutter, L. Kotthoff, and J. Vanschoren, editors, *Auto-*
492 *ated Machine Learning: Methods, Systems, Challenges*, volume 5 of *The Springer Series*
493 *on Challenges in Machine Learning*, pages 35–61. Springer, 2019. Available for free at
494 <http://automl.org/book>.
- 495 [30] R. Bardenet, M. Brendel, B. Kégl, and M. Sebag. Collaborative hyperparameter tuning. In
496 S. Dasgupta and D. McAllester, editors, *Proceedings of the 30th International Conference on*
497 *Machine Learning (ICML'13)*, pages 199–207. Omnipress, 2013.
- 498 [31] D. Yogatama and G. Mann. Efficient transfer learning method for automatic hyperparameter
499 tuning. In S. Kaski and J. Corander, editors, *Proceedings of the Seventeenth International*
500 *Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 33 of *JMLR Workshop*
501 *and Conference Proceedings*, pages 1077–1085, 2014.
- 502 [32] M. Feurer, J. Springenberg, and F. Hutter. Initializing Bayesian hyperparameter optimization
503 via meta-learning. In B. Bonet and S. Koenig, editors, *Proceedings of the Twenty-ninth National*
504 *Conference on Artificial Intelligence (AAAI'15)*, pages 1128–1135. AAAI Press, 2015.
- 505 [33] M. Wistuba, N. Schilling, and L. Schmidt-Thieme. Scalable Gaussian process-based transfer
506 surrogates for hyperparameter optimization. *Machine Learning*, 107(1):43–78, 2018.
- 507 [34] H.-P. Kriegel, E. Schubert, and A. Zimek. The (black) art of runtime evaluation: Are we
508 comparing algorithms or implementations? *Knowl. Inf. Syst.*, 52(2):341–378, 2017.
- 509 [35] S. Narang, H. Chung, Y. Tay, W. Fedus, T. Fevry, M. Matena, K. Malkan, N. Fiedel, N. Shazeer,
510 Z. Lan, Y. Zhou, W. Li, N. Ding, J. Marcus, A. Roberts, and C. Raffel. Do transformer
511 modifications transfer across implementations and applications? *arxiv:2102.11972 [cs.LG]*,
512 2021.
- 513 [36] J. Bergstra, D. Yamins, and D. Cox. Making a science of model search: Hyperparameter opti-
514 mization in hundreds of dimensions for vision architectures. In S. Dasgupta and D. McAllester,
515 editors, *Proceedings of the 30th International Conference on Machine Learning (ICML'13)*,
516 pages 115–123. Omnipress, 2013.
- 517 [37] G. Melis, C. Dyer, and P. Blunsom. On the state of the art of evaluation in neural language
518 models. In *Proceedings of the International Conference on Learning Representations (ICLR'18)*,
519 2018. Published online: iclr.cc.
- 520 [38] D. Hand. Classifier Technology and the Illusion of Progress. *Statistical Science*, 21(1):1 – 14,
521 2006.
- 522 [39] M. Dacrema, P. Cremonesi, and D. Jannach. Are we really making much progress? a worrying
523 analysis of recent neural recommendation approaches. In *RecSys '19: Proceedings of the*
524 *13th ACM Conference on Recommender Systems*, page 101–109, New York, NY, USA, 2019.
525 Association for Computing Machinery.
- 526 [40] X. Bouthillier, P. Delaunay, M. Bronzi, A. Trofimov, B. Nichyporuk, J. Szeto, N. Sepah, E. Raff,
527 K. Madan, V. Voleti, S. Kahou, V. Michalski, D. Serdyuk, T. Arbel, C. Pal, G. Varoquaux,
528 and P. Vincent. Accounting for variance in machine learning benchmarks. *arXiv:2103.03098*
529 *[cs.LG]*, 2021.
- 530 [41] T. Weise, R. Chiong, and K. Tang. Evolutionary optimization: Pitfalls and booby traps. *Journal*
531 *of Computer Science and Technology*, 27:907–936, 2012.
- 532 [42] A. Howe and E. Dahlman. A critical assessment of benchmark comparison in planning. *Journal*
533 *of Artificial Intelligence Research*, 17:1–33, 2002.
- 534 [43] I. Gent, S. Grant, E. MacIntyre, P. Prosser, P. Shaw, B. Smith, and T. Walsh. How not to do it.
535 Technical Report 97.92, University of Leeds, 1997.
- 536 [44] K. Eggensperger, M. Lindauer, and F. Hutter. Pitfalls and best practices in algorithm configura-
537 tion. *Journal of Artificial Intelligence Research*, pages 861–893, 2019.

- 538 [45] M. Lindauer and F. Hutter. Best practices for scientific research on neural architecture search.
539 *Journal of Machine Learning Research*, 21:1–18, 2020.
- 540 [46] T. Bartz-Beielstein, C. Doerr, J. Bossek, S. Chandrasekaran, T. Eftimov, A. Fischbach,
541 P. Kerschke, M. López-Ibáñez, K. Malan, J. Moore, B. Naujoks, P. Orzechowski, V. Volz,
542 M. Wagner, and T. Weise. Benchmarking in optimization: Best practice and open issues.
543 *arXiv:2007.03488v2 [cs.NE]*, 2020.
- 544 [47] F. Hutter, M. López-Ibáñez, C. Fawcett, M. Lindauer, H. Hoos, K. Leyton-Brown, and T. Stützle.
545 ACLib: a benchmark library for algorithm configuration. In P. Pardalos and M. Resende, editors,
546 *Proceedings of the Eighth International Conference on Learning and Intelligent Optimization*
547 (*LION'14*), Lecture Notes in Computer Science, pages 36–40. Springer, 2014.
- 548 [48] S. Pineda, H. Jomaa, M. Wistuba, and J. Grabocka. HPO-B: A large-scale reproducible
549 benchmark for black-box HPO based on OpenML. *arXiv:2106.06257 [cs.LG]*, 2021.
- 550 [49] I. Guyon, L. Sun-Hosoya, M. Boullé, H. Escalante, S. Escalera, Z. Liu, D. Jajetic, B. Ray,
551 M. Saeed, M. Sebag, A. Statnikov, W. Tu, and E. Viegas. Analysis of the AutoML challenge
552 series 2015-2018. In F. Hutter, L. Kotthoff, and J. Vanschoren, editors, *Automated Machine*
553 *Learning: Methods, Systems, Challenges*, volume 5 of *The Springer Series on Challenges in*
554 *Machine Learning*. Springer, 2019. Available for free at <http://automl.org/book>.
- 555 [50] Z. Liu, A. Pavao, Z. Xu, S. Escalera, F. Ferreira, I. Guyon, S. Hong, F. Hutter, R. Ji, J. Jacques
556 Junior, G. Li, M. Lindauer, Z. Luo, M. Madadi, T. Nierhoff, K. Niu, C. Pan, D. Stoll, S. Treguer,
557 J. Wang, P. Wang, C. Wu, , Y. Xiong, A. Zela, and Y. Zhang. Winning solutions and post-
558 challenge analyses of the ChaLearn AutoDL challenge 2019. *hal-02957135*, 2020.
- 559 [51] R. Turner, D. Eriksson, M. McCourt, J. Kiili, E. Laaksonen, Z. Xu, and I. Guyon. Bayesian
560 optimization is superior to random search for machine learning hyperparameter tuning: Analysis
561 of the black-box optimization challenge 2020. *arXiv:2104.10201 [cs.LG]*, 2021.
- 562 [52] J. Snoek, H. Larochelle, and R. Adams. Practical Bayesian optimization of machine learning
563 algorithms. In P. Bartlett, F. Pereira, C. Burges, L. Bottou, and K. Weinberger, editors, *Pro-*
564 *ceedings of the 25th International Conference on Advances in Neural Information Processing*
565 *Systems (NeurIPS'12)*, pages 2960–2968. Curran Associates, 2012.
- 566 [53] M. Wistuba, N. Schilling, and L. Schmidt-Thieme. Two-stage transfer surrogate model for
567 automatic hyperparameter optimization. In B. Berendt, B. Bringmann, E. Fromont, G. Garriga,
568 P. Miettinen, N. Tatti, and V. Tresp, editors, *Machine Learning and Knowledge Discovery in*
569 *Databases (ECML/PKDD'16)*, volume 9851 of *Lecture Notes in Artificial Intelligence*, pages
570 199–214. Springer, 2016.
- 571 [54] C. Ying, A. Klein, E. Real, E. Christiansen, K. Murphy, and F. Hutter. NAS-Bench-101:
572 Towards reproducible neural architecture search. In K. Chaudhuri and R. Salakhutdinov, editors,
573 *Proceedings of the 36th International Conference on Machine Learning (ICML'19)*, volume 97,
574 pages 7105–7114. Proceedings of Machine Learning Research, 2019.
- 575 [55] L. Metz, N. Maheswaranathan, R. Sun, C. Freeman, B. Poole, and J. Sohl-Dickstein. Using
576 a thousand optimization tasks to learn hyperparameter search strategies. *arXiv:2002.11887*
577 *[cs.LG]*, 2020.
- 578 [56] K. Eggensperger, F. Hutter, H.H. Hoos, and K. Leyton-Brown. Efficient benchmarking of
579 hyperparameter optimizers via surrogates. In B. Bonet and S. Koenig, editors, *Proceedings of*
580 *the Twenty-ninth National Conference on Artificial Intelligence (AAAI'15)*, pages 1114–1120.
581 AAAI Press, 2015.
- 582 [57] K. Eggensperger, M. Lindauer, H. H. Hoos, F. Hutter, and K. Leyton-Brown. Efficient bench-
583 marking of algorithm configurators via model-based surrogates. *Machine Learning*, 107(1):
584 15–41, 2018.
- 585 [58] V. Perrone, R. Jenatton, M. Seeger, and C. Archambeau. Scalable hyperparameter transfer
586 learning. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett,
587 editors, *Proceedings of the 31st International Conference on Advances in Neural Information*
588 *Processing Systems (NeurIPS'18)*, pages 12751–12761. Curran Associates, 2018.

- 589 [59] R. Martinez-Cantin. Funneled Bayesian optimization for design, tuning and control of au-
590 tonomous systems. *IEEE Transactions on Cybernetics*, 49(4):1489–1500, 2019.
- 591 [60] A. Klein, Z. Dai, F. Hutter, N. Lawrence, and J. Gonzalez. Meta-surrogate benchmarking for
592 hyperparameter optimization. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc,
593 E. Fox, and R. Garnett, editors, *Proceedings of the 32nd International Conference on Advances
594 in Neural Information Processing Systems (NeurIPS’19)*, pages 6267–6277. Curran Associates,
595 2019.
- 596 [61] E. Daxberger, A. Makarova, M. Turchetta, and A. Krause. Mixed-variable bayesian optimization.
597 In C. Bessiere, editor, *Proceedings of the Twenty-Ninth International Joint Conference on
598 Artificial Intelligence (IJCAI-20)*, pages 2633–2639. ijcai.org, 2020.
- 599 [62] J. Siems, L. Zimmer, A. Zela, J. Lukasik, M. Keuper, and F. Hutter. NAS-Bench-301 and the
600 case for surrogate benchmarks for neural architecture search. *arXiv:2008.09777 [cs.LG]*, 2020.
- 601 [63] G. Kurtzer, V. Sochat, and M. Bauer. Singularity: Scientific containers for mobility of compute.
602 *PLOS ONE*, 12(5), 2017.
- 603 [64] A. Klein and F. Hutter. Tabular benchmarks for joint architecture and hyperparameter optimiza-
604 tion. *arXiv:1905.04970 [cs.LG]*, 2019.
- 605 [65] X. Dong and Y. Yang. NAS-Bench-201: Extending the scope of reproducible neural architecture
606 search. In *Proceedings of the International Conference on Learning Representations (ICLR’20)*,
607 2020. Published online: iclr.cc.
- 608 [66] A. Zela, J. Siems, and F. Hutter. NAS-Bench-1Shot1: Benchmarking and dissecting one-
609 shot neural architecture search. In *Proceedings of the International Conference on Learning
610 Representations (ICLR’20)*, 2020. Published online: iclr.cc.
- 611 [67] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,
612 P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher,
613 M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine
614 Learning Research*, 12:2825–2830, 2011.
- 615 [68] T. Chen and C. Guestrin. Xgboost: A scalable tree boosting system. In B. Krishnapuram,
616 M. Shah, A. Smola, C. Aggarwal, D. Shen, and R. Rastogi, editors, *Proceedings of the 22nd
617 ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*,
618 pages 785–794. ACM Press, 2016.
- 619 [69] P. Gijbbers, E. LeDell, S. Poirier, J. Thomas, B. Bischl, and J. Vanschoren. An open source
620 automl benchmark. In K. Eggenberger, M. Feurer, F. Hutter, and J. Vanschoren, editors, *ICML
621 workshop on Automated Machine Learning (AutoML workshop 2019)*, 2019.
- 622 [70] J. Vanschoren, J. N. Van Rijn, B. Bischl, and L. Torgo. Openml: networked science in machine
623 learning. *ACM SIGKDD Explorations Newsletter*, 15(2):49–60, 2014.
- 624 [71] A. Cowen-Rivers, W. Lyu, R. Tutunov, Z. Wang, A. Grosnit, R. Griffiths, A. Maraval, H. Jianye,
625 J. Wang, J. Peters, and H. Ammar. An empirical study of assumptions in bayesian optimisation.
626 *arXiv:2012.03826 [cs.LG]*, 2021.
- 627 [72] D. Eriksson, M. Pearce, J. Gardner, R. Turner, and M. Poloczek. Scalable global optimization
628 via local bayesian optimization. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc,
629 E. Fox, and R. Garnett, editors, *Proceedings of the 32nd International Conference on Advances
630 in Neural Information Processing Systems (NeurIPS’19)*. Curran Associates, 2019.
- 631 [73] F. Hutter, H. Hoos, and K. Leyton-Brown. Sequential model-based optimization for general
632 algorithm configuration. In C. Coello, editor, *Proceedings of the Fifth International Conference
633 on Learning and Intelligent Optimization (LION’11)*, volume 6683 of *Lecture Notes in Computer
634 Science*, pages 507–523. Springer, 2011.
- 635 [74] K. Kandasamy, K. Vysyaraju, W. Neiswanger, B. Paria, C. Collins, J. Schneider, B. Póczos,
636 and E. Xing. Tuning hyperparameters without grad students: Scalable and robust bayesian
637 optimisation with Dragonfly. *Journal of Machine Learning Research*, 21(81):1–27, 2020.

- 638 [75] P. Moritz, R. Nishihara, S. Wang, A. Tumanov, R. Liaw, E. Liang, M. Elibol, Z. Yang,
639 W. Paul, M. Jordan, and I. Stoica. Ray: A distributed framework for emerging AI applications.
640 *arXiv:1712.05889v2 [cs.DC]*, 2018.
- 641 [76] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama. Optuna: A next-generation hyper-
642 parameter optimization framework. In A. Teredesai, V. Kumar, Y. Li, R. Rosales, E. Terzi,
643 and G. Karypis, editors, *Proceedings of the 25th ACM SIGKDD International Conference on*
644 *Knowledge Discovery & Data Mining, KDD'19*, pages 2623–2631. ACM Press, 2019.
- 645 [77] J. Bergstra and Y. Bengio. Random search for hyper-parameter optimization. *Journal of*
646 *Machine Learning Research*, 13:281–305, 2012.
- 647 [78] M. Balandat, B. Karrer, D. Jiang, S. Daulton, B. Letham, A. Wilson, and E. Bakshy. BoTorch:
648 A Framework for Efficient Monte-Carlo Bayesian Optimization. In H. Larochelle, M. Ranzato,
649 R. Hadsell, M.F. Balcan, and H. Lin, editors, *Proceedings of the 33rd International Conference*
650 *on Advances in Neural Information Processing Systems (NeurIPS'20)*. Curran Associates, 2020.
- 651 [79] F. Hutter and M. Osborne. A kernel for hierarchical parameter spaces. *arXiv:1310.5738v1*
652 *[stats.ML]*, 2013.
- 653 [80] M. Lindauer, M. Feurer, K. Eggensperger, A. Biedenkapp, and F. Hutter. Towards assessing the
654 impact of bayesian optimization's own hyperparameters. In P. De Causmaecker, M. Lombardi,
655 and Y. Zhang, editors, *IJCAI 2019 DSO Workshop*, 2019.
- 656 [81] J. Demšar. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine*
657 *Learning Research*, 7:1–30, 2006.
- 658 [82] L. Hertel, J. Collado, P. Sadowski, J. Ott, and P. Baldi. Sherpa: Robust hyperparameter
659 optimization for machine learning. *SoftwareX*, 12:100591, 2020.
- 660 [83] Y. Xiao, E. Xing, and W. Neiswanger. Amortized auto-tuning: Cost-efficient transfer optimiza-
661 tion for hyperparameter recommendation. *arXiv:2106.09179 [cs.LG]*, 2021.
- 662 [84] L. Cardoso, V. Santos, R. Francês, R. Prudêncio, and R. Alves. Data vs classifiers, who wins?
663 *arXiv:2107.07451 [cs.LG]*, 2021.
- 664 [85] T. Gebru, J. Morgenstern, B. Vecchione, J. Vaughan, H. Wallach, H. Daumé III, and K. Crawford.
665 Datasheets for datasets. *arXiv:1803.09010 [cs.DB]*, 2020.
- 666 [86] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba.
667 OpenAI Gym. *arXiv:1606.01540 [cs.LG]*, 2016.
- 668 [87] B. Bischl, G. Casalicchio, M. Feurer, F. Hutter, M. Lang, R. Mantovani, J. N. van Rijn, and
669 J. Vanschoren. Openml benchmarking suites. *arXiv:1708.03731v2*, 2019.
- 670 [88] D. Molina, A. Latorre, and F. Herrera. An insight into bio-inspired and evolutionary algorithms
671 for global optimization: Review, analysis, and lessons learnt over a decade of competitions.
672 *Cognitive Computation*, 10:517–544, 2018.
- 673 [89] A. Kuhnle, M. Schaarschmidt, and K. Fricke. Tensorforce: a TensorFlow library for ap-
674 plied reinforcement learning. Web page, 2017. URL [https://github.com/tensorforce/](https://github.com/tensorforce/tensorforce)
675 [tensorforce](https://github.com/tensorforce/tensorforce).
- 676 [90] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
677 algorithms. *arXiv:1707.06347 [cs.LG]*, 2017.
- 678 [91] T. Chen, E.B. Fox, and C. Guestrin. Stochastic gradient Hamiltonian Monte Carlo. In E. Xing
679 and T. Jebara, editors, *Proceedings of the 31th International Conference on Machine Learning,*
680 *(ICML'14)*. Omnipress, 2014.
- 681 [92] J. Springenberg, A. Klein, S. Falkner, and F. Hutter. Bayesian optimization with robust Bayesian
682 neural networks. In D. Lee, M. Sugiyama, U. von Luxburg, I. Guyon, and R. Garnett, editors,
683 *Proceedings of the 29th International Conference on Advances in Neural Information Processing*
684 *Systems (NeurIPS'16)*. Curran Associates, 2016.

- 685 [93] D. Dua and C. Graff. UCI machine learning repository, 2019. URL <http://archive.ics.uci.edu/ml>.
686
- 687 [94] S. Dieleman, J. Schlüter, C. Raffel, E. Olson, S. Sønderby, D. Nouri, D. Maturana, M. Thoma,
688 E. Battenberg, J. Kelly, J. De Fauw, M. Heilman, diogo149, B. McFee, H. Weideman, takacs84,
689 peterderivaz, Jon, instagibbs, K. Rasul, CongLiu, Britefury, and J. Degraeve. Lasagne: First
690 release., 2015.
- 691 [95] Theano Development Team. Theano: A Python framework for fast computation of mathematical
692 expressions. *arXiv:1605.02688 [cs.SC]*, 2016.
- 693 [96] D. Jones, M. Schonlau, and W. Welch. Efficient global optimization of expensive black box
694 functions. *Journal of Global Optimization*, 13:455–492, 1998.
- 695 [97] F. Hutter. *Automated Configuration of Algorithms for Solving Hard Computational Problems*.
696 PhD thesis, University of British Columbia, Department of Computer Science, Vancouver,
697 Canada, 2009.
- 698 [98] M. Lindauer, K. Eggensperger, M. Feurer, S. Falkner, A. Biedenkapp, and F. Hutter. SMACv3:
699 Algorithm configuration in Python. github.com/automl/SMAC3, 2017.
- 700 [99] L. Li, K. Jamieson, A. Rostamizadeh, E. Gonina, J. Bentzur, M. Hardt, B. Recht, and A. Tal-
701 walkar. A system for massively parallel hyperparameter tuning. In I. Dhillon, D. Papailiopoulos,
702 and V. Sze, editors, *Proceedings of Machine Learning and Systems 2*, volume 2, pages 230–246,
703 2020.

704 Checklist

- 705 1. For all authors...
- 706 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
707 contributions and scope? [Yes]
- 708 (b) Did you describe the limitations of your work? [Yes] See Section 6
- 709 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 710 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
711 them? [Yes]
- 712 2. If you are including theoretical results...
- 713 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 714 (b) Did you include complete proofs of all theoretical results? [N/A]
- 715 3. If you ran experiments (e.g. for benchmarks)...
- 716 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
717 mental results (either in the supplemental material or as a URL)? [Yes] Yes, we release
718 all scripts here: <https://github.com/automl/HPOBenchExperimentUtils>
- 719 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
720 were chosen)? [Yes] See Section 5.1 and the appendix.
- 721 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
722 ments multiple times)? [Yes] We do where applicable (see figures in the appendix),
723 but we also compute statistical tests to quantify the statistical significance of our results
724 (see Tables 2, Table 3 and Appendix H).
- 725 (d) Did you include the total amount of compute and the type of resources used (e.g., type
726 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix E and Appendix F.
- 727 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 728 (a) If your work uses existing assets, did you cite the creators? [Yes] See Table 1
- 729 (b) Did you mention the license of the assets? [Yes] See Section 1.
- 730 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
731 Our benchmarks are available here: <https://github.com/automl/HPOBench>

- 732 (d) Did you discuss whether and how consent was obtained from people whose data you're
733 using/curating? [N/A]
- 734 (e) Did you discuss whether the data you are using/curating contains personally identifiable
735 information or offensive content? [N/A]
- 736 5. If you used crowdsourcing or conducted research with human subjects...
- 737 (a) Did you include the full text of instructions given to participants and screenshots, if
738 applicable? [N/A]
- 739 (b) Did you describe any potential participant risks, with links to Institutional Review
740 Board (IRB) approvals, if applicable? [N/A]
- 741 (c) Did you include the estimated hourly wage paid to participants and the total amount
742 spent on participant compensation? [N/A]