

---

# HPO-B: A Large-Scale Reproducible Benchmark for Black-Box HPO based on OpenML

---

**Sebastian Pineda Arango\***  
University of Freiburg  
pineda@cs.uni-freiburg.de

**Hadi S. Jomaa\***  
University of Hildesheim  
hsjomaa@ismll.uni-hildesheim.de

**Martin Wistuba**  
Amazon Research  
marwistu@amazon.com

**Josif Grabocka**  
University of Freiburg  
grabocka@cs.uni-freiburg.de

## Abstract

1 Hyperparameter optimization (HPO) is a core problem for the machine learning  
2 community and remains largely unsolved due to the significant computational re-  
3 sources required to evaluate hyperparameter configurations. As a result, a series of  
4 recent related works have focused on the direction of transfer learning for quickly  
5 fine-tuning hyperparameters on a dataset. Unfortunately, the community does  
6 not have a common large-scale benchmark for comparing HPO algorithms. In-  
7 stead, the *de facto* practice consists of empirical protocols on arbitrary small-scale  
8 meta-datasets that vary inconsistently across publications, making reproducibility  
9 a challenge. To resolve this major bottleneck and enable a fair and fast comparison  
10 of black-box HPO methods on a level playing field, we propose **HPO-B**, a new  
11 large-scale benchmark in the form of a collection of meta-datasets. Our benchmark  
12 is assembled and preprocessed from the OpenML repository and consists of 176  
13 search spaces (algorithms) evaluated sparsely on 196 datasets with a total of 6.4  
14 million hyperparameter evaluations. For ensuring reproducibility on our bench-  
15 mark, we detail explicit experimental protocols, splits, and evaluation measures for  
16 comparing methods for both non-transfer, as well as, transfer learning HPO.

## 17 1 Introduction

18 Hyperparameter Optimization (HPO) is arguably the major open challenge for the machine learning  
19 community due to the expensive computational resources demanded to evaluate configurations.  
20 As a result, HPO and its broader umbrella research area, AutoML, have drawn particular interest  
21 over the past decade [2, 14, 26, 27]. Black-box HPO is a specific sub-problem that focuses on  
22 the case where the function to be optimized (e.g. the generalization performance of an algorithm)  
23 is unknown, non-differentiable with respect to the hyperparameters, and intermediate evaluation  
24 proxies are not computable (opposed to gray-box HPO [19] which accesses intermediate performance  
25 measurements).

26 Although black-box HPO is a core problem, existing solutions based on parametric surrogate models  
27 for estimating the performance of a configuration overfit the limited number of evaluated configu-  
28 rations. As a result, the AutoML community has recently invested efforts in resolving the sample-  
29 inefficiency of parametric surrogates via meta- and transfer-learning [10, 15, 23, 24, 29, 31, 34].

---

\*Equal contribution

30 Unfortunately, despite the promising potential of transfer-learning in black-box HPO, the impact of  
 31 such algorithms is hindered by their poor experimental reproducibility. Our personal prior research  
 32 experience, as well as the feedback from the community, highlight that reproducing and generalizing  
 33 the results of transfer-learning HPO methods is challenging. In essence, the problem arises when the  
 34 results of a well-performing method in the experimental protocol of a publication either can not be  
 35 replicated; or when the method underperforms in a slightly different empirical protocol. We believe  
 36 that a way of resolving this negative *impasse* is to propose a new public large-scale benchmark for  
 37 comparing HPO methods, where the exact training/validation/test splits of the meta-datasets, the  
 38 exact evaluation protocol, and the performance measures are well-specified. The strategy of adopting  
 39 benchmarks is a trend in related areas, such as in computer vision [7], or NAS [8, 36].

40 In this perspective, we present **HPO-B<sup>2</sup>**, the largest public benchmark of meta-datasets for black-box  
 41 HPO containing 6.4M hyperparameter evaluations across 176 search spaces (algorithms) and on 196  
 42 datasets in total. The collection is derived from the raw data of OpenML [28], but underwent an  
 43 extensive process of cleaning, preprocessing and organization (Section 5). Additionally, we offer  
 44 off-the-shelf ready variants of the benchmark that are adapted for both non-transfer, as well as transfer  
 45 HPO experiments, together with the respective evaluation protocols (Section 6). This large, diverse,  
 46 yet plug-and-play benchmark can significantly boost future research in black-box HPO.

## 47 2 Terminology

48 To help the reader navigate through our paper, we present the compact thesaurus of Table 1 for  
 49 defining the vernacular of the HPO community.

Term	Definition
Configuration	Specific settings/values of hyperparameters
Search space	The domain of a configuration: scale and range of each hyperparameter’s values
Response	The performance of an algorithm given a configuration and dataset
Surrogate	A (typically parametric) function that approximates the response
Seed	Set of initial configurations used to fit the initial surrogate model
Black-box	The response is an unknown and non-differentiable function of a configuration
Task	An HPO problem given a search space and a dataset
Evaluation	The measured response of a configuration on a dataset
Trial	An evaluation on a task during the HPO procedure
Meta-dataset	Collection of <i>recorded</i> evaluations from different tasks on a search space
Meta-instance	An evaluation in the meta-dataset for one of the tasks
Meta-feature	Descriptive attributes of a dataset
Source tasks	In a meta- or transfer-learning setup refers to the <i>known</i> tasks we <i>train from</i>
Target tasks	In a meta- or transfer-learning setup refers to the <i>new</i> tasks we <i>test on</i>
Benchmark	<b>New definition:</b> Collection of meta-datasets from different search spaces

Table 1: A thesaurus of the common HPO terminology used throughout this paper

## 50 3 Related Work

51 **Non-transfer black-box HPO:** The mainstream paradigm in HPO relies on surrogates to estimate the  
 52 performance of hyperparameter configurations. For example, [2] were the first to propose Gaussian  
 53 Processes (GP) as surrogates. The same authors also propose a Tree Parzen Estimator (TPE) for  
 54 computing the non-parametric densities of the hyperparameters given the observed performances.  
 55 Both approaches achieve a considerable lift over random [3] and manual search. To address the cubic  
 56 run-time complexity of GPs concerning the number of evaluated configurations, DNGO [26] trains  
 57 neural networks for generating adaptive basis functions of hyperparameters, in combination with a  
 58 Bayesian linear regressor that models uncertainty. Alternatively, SMAC [14] represents the surrogate  
 59 as a random forest, and BOHAMIANN [27] employs Bayesian neural networks instead of plain  
 60 neural networks to estimate the uncertainty of a configuration’s performance. For an extensive study

<sup>2</sup>The benchmark is publicly available at <https://github.com/releaunifreiburg/HPO-B>

61 on non-transfer Bayesian optimization techniques for HPO, we refer the readers to [5, 25] that study  
 62 the impact of the underlying assumptions associated with black-box HPO algorithms.

63 **Transfer black-box HPO:** To expedite HPO, it is important to leverage information from existing  
 64 evaluations of configurations from prior tasks. A common approach is to capture the similarity  
 65 between datasets using meta-features (i.e. descriptive dataset characteristics). Meta-features have  
 66 been used as a warm-start initialization technique [11, 16], or as part of the surrogate directly [1].  
 67 Transfer learning is also explored through the weighted combination of surrogates, such as in TST-  
 68 R [34], RGPE [10], and TAF-R [35]. Another direction is learning a shared surrogate across tasks.  
 69 ABLR optimizes a shared hyperparameter embedding with separate Bayesian linear regressors per  
 70 task [20], while GCP [23] maps the hyperparameter response to a shared distribution with a Gaussian  
 71 Copula process. Furthermore, FSBO [31] meta-learns a deep-kernel Gaussian Process surrogate,  
 72 whereas DMFBS incorporates the dataset context through end-to-end meta-feature networks [16].

73 **Meta-datasets:** The work by Wistuba et al. [33] popularised the usage of meta-dataset benchmarks  
 74 with pre-computed evaluations for the hyperparameters of SVM (288 configurations) and Adaboost  
 75 (108 configurations) on 50 datasets; a benchmark that inspired multiple follow-up works [10, 30].  
 76 Existing attempts to provide HPO benchmarks deal only with the non-transfer black-box HPO  
 77 setup [9]. As they contain results for one or very few datasets per search space, they cannot be  
 78 used for the evaluation of transfer black-box HPO methods. Nevertheless, there is a trend in using  
 79 evaluations of search spaces from the OpenML repository [12], which contains evaluations reported  
 80 by an open community, as well as large-scale experiments contributed by specific research labs [4, 18].  
 81 However, the choice of OpenML search spaces in publications is ad-hoc: one related work uses  
 82 SVM and XGBoost [20], a second uses GLMNet and SVM [31], while a third paper uses XGBoost,  
 83 Random Forest and SVM [21]. We assess that the community *i*) inconsistently cherry-picks (assuming  
 84 *bona fides*) search spaces, with *ii*) arbitrary train/validation/test splits of the tasks within the meta-  
 85 dataset, and *iii*) inconsistent preprocessing of hyperparameters and responses. In our experiments, we  
 86 observed that existing methods do not generalize well on new meta-datasets (Section 7).

87 **Our Novelty:** As a remedy, we propose a novel benchmark derived from OpenML [12], that resolves  
 88 the existing reproducibility issues of existing non-transfer and transfer black-box HPO methods, by  
 89 ensuring a fairly-reproducible empirical protocol. The contributions of our benchmark are multi-fold.  
 90 First of all, we remove the confounding factors induced by different meta-dataset preprocessing  
 91 pipelines (e.g. hyperparameter scaling and transformations, missing value imputations, one-hot  
 92 encodings, etc.). Secondly, we provide a specified collection of search spaces, with specified datasets  
 93 and evaluations. Furthermore, for transfer learning HPO methods, we also provide pre-defined  
 94 training/validation/testing splits of tasks. For experiments on the test tasks, we additionally provide  
 95 5 seeds (i.e. 5 sets of initial hyperparameters to fit the initial surrogate) with 5 hyperparameter  
 96 configurations, each. We also highlight recommended empirical measures for comparing HPO  
 97 methods and assessing their statistical significance in Section 6. In that manner, the results of different  
 98 papers that use our benchmark can be compared directly without fearing the confounding factors.  
 99 Table 2 presents a summary of the descriptive statistics of meta-datasets from prior literature. To the  
 100 best of our awareness, the proposed benchmark is also richer (in the number of search spaces and  
 101 their dimensionality) and larger (in the number of evaluations) than all the prior protocols.

Paper	Venue/Year	# Search Spaces	# Datasets	# HPs	# Evals.
[1]	ICML '13	1	29	2	3K
[33]	DSAA '15	2	50	2, 4	20K
[11]	AAAI '15	3	57	4, 5	93K
[32]	ECML-PKDD '15	17	59	1-7	1.3M
[20]	NeurIPS '18	2	30	4, 10	655K
[23]	ICML '20	4	26	6, 9	343K
[16]	DMKD '21	1	120	7	414K
[31]	ICLR '21	3	80	2, 4	864K
Our HPO-B-v1	-	176	196	1-53	6.39M
Our HPO-B-v2/-v3	-	16	101	2-18	6.34M

Table 2: Summary statistics for various meta-datasets considered in prior works.

## 102 4 A Brief Explanation of Bayesian Optimization Concepts

103 As we often refer to HPO methods, in this section we present a brief coverage of Bayesian optimization  
104 as the most popular HPO method for black-box optimization. HPO aims at minimizing the function  
105  $f : \mathcal{X} \rightarrow \mathbb{R}$  which maps each hyperparameter configuration  $\mathbf{x} \in \mathcal{X}$  to the validation loss obtained  
106 when training the machine learning model using  $\mathbf{x}$ . Bayesian optimization keeps track of all evaluated  
107 hyperparameter configurations in a history  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_i$ , where  $y_i \sim \mathcal{N}(f(\mathbf{x}_i), \sigma_n^2)$  is the (noisy)  
108 response. A probabilistic model, the so-called surrogate model, is used to approximate the behavior  
109 of the response function. Gaussian Processes are a common choice for the surrogate model [22].

110 Bayesian optimization is an iterative process that alternates between updating the surrogate model as  
111 described above and selecting the next hyperparameter configuration. The latter is done by finding the  
112 configuration which maximizes an acquisition function, which scores each feasible hyperparameter  
113 configuration using the surrogate model by finding a trade-off between exploration and exploitation.  
114 Arguably, the most popular acquisition function is the Expected Improvement [17].

115 The efficiency of Bayesian optimization depends on the surrogate model’s ability to approximate the  
116 response function. However, this is a challenging task since every optimization starts with no or little  
117 knowledge about the response function. To overcome this cold-start problem, transfer methods have  
118 been proposed, where the transfer scenario slightly changes the problem definition. The main objective  
119 remains finding the configuration  $\mathbf{x}_* \in \mathcal{X}^{(\bar{s})}$  which optimizes the target response function  $f^{(\bar{s}, \bar{t})}$   
120 for a given search space  $\mathcal{X}^{(\bar{s})}$ . The knowledge about the target response is now denoted by  $\mathcal{D}^{(\bar{s}, \bar{t})}$ .  
121 Additionally, side information, the so-called *meta-dataset*, for  $S$ -many search spaces and  $T$ -many  
122 datasets,  $\mathcal{D} = \bigcup_{s=1..S, t=1..T} \mathcal{D}^{(s,t)}$ , is available. Here,  $\mathcal{D}^{(s,t)} = \{(\mathbf{x}_i^{(s,t)}, y_i^{(s,t)})\}_i$ ,  $\mathbf{x}_i^{(s,t)} \in \mathcal{X}^{(s)}$ ,  
123 and  $y_i^{(s,t)}$  is the noisy observation of the source response function  $f_i^{(s,t)} = f^{(t)}(\mathbf{x}_i^{(s,t)})$ .

## 124 5 Benchmark Description

125 The benchmark HPO-B is a collection of meta-datasets collected from OpenML [12] with a diverse  
126 set of search spaces. We present three different versions of the meta-data set, as follows:

- 127 • **HPO-B-v1:** The raw benchmark of all 176 meta-datasets;
- 128 • **HPO-B-v2:** Subset of 16 meta-datasets with the most frequent search spaces;
- 129 • **HPO-B-v3:** Split of HPO-B-v2 into training, validation and testing.

130 When assembling the benchmark HPO-B-v1 we noticed that most of the evaluations are reported  
131 for a handful of popular search spaces, in particular, we noticed that 9% of the top meta-datasets  
132 include 99.3% of the evaluations. As a result, we created a second version HPO-B-v2 that includes  
133 only the frequent meta-datasets that have at least 10 datasets with at least 100 evaluations per dataset  
134 (Section 5.1). Furthermore, as we clarified in Section 3 a major reproducibility issue of the related  
135 work on transfer HPO is the lack of clear training, validation, and test splits for the meta-datasets. To  
136 resolve this issue, we additionally created HPO-B-v3 as a derivation of HPO-B-v2 with pre-defined  
137 splits of the training, validation, and testing tasks for every meta-dataset, in addition to providing  
138 initial configurations (seeds) for the test tasks. The three versions were designed to fulfill concrete  
139 purposes with regards to different types of HPO methods. For non-transfer black-box HPO methods,  
140 we recommend using HPO-B-v2 which offers a large pool of HPO tasks. Naturally, for transfer HPO  
141 tasks we recommend using HPO-B-v3 where meta-datasets are split into training, validation, and  
142 testing. We still are releasing the large HPO-B-v1 benchmark to anticipate next-generation methods  
143 for heterogeneous transfer learning techniques that meta-learn surrogates across different search  
144 spaces, where all 176 meta-datasets might be useful despite most of them having few evaluations.

145 Concretely, HPO-B-v3 contains the set of filtered search spaces of HPO-B-v2, which are specially  
146 split into *four* sets: meta-train, meta-validation, meta-test and an augmented version of the meta-train  
147 dataset. Every split contains different datasets from the same search space. We distributed the datasets  
148 per search space as 80% of the datasets to meta-train, 10% to meta-validation, and 10% to meta-test,  
149 respectively. A special, augmented version of the meta-train is created by adding all other search  
150 space evaluations from HPO-B-v1 that are not part of HPO-B-v3. On the other hand, in HPO-B-v3  
151 we also provide seeds for initializing the HPO. They are presented as five different sets of five initial

152 configurations to be used by a particular HPO method. By providing five different seeds we decrease  
 153 the random effect of the specific initial configurations. To ease the comparison among HPO methods,  
 154 we suggest using the recommended initial configurations for testing. Although, we admit that some  
 155 algorithms proposing novel warm-starting strategies might need to bypass the recommended initial  
 156 configurations.

## 157 5.1 Benchmark summary

158 The created benchmark contains 6,394,555 total evaluations across 176 search spaces that are sparsely  
 159 evaluated on 196 datasets. By accounting for the search spaces that comply with our filtering criteria  
 160 (at least 10 datasets with 100 evaluations), we obtain HPO-B-v2 with 16 different search spaces and  
 161 6,347,916 evaluations on 101 datasets. Notice that the benchmark does not include evaluations for  
 162 all datasets in every search space. The number of dimensions, datasets, and evaluations per search  
 163 space is listed in Table 3. An additional description of the rest of all the search spaces in HPO-B-v1  
 164 is presented in the Appendix. In addition, Table 3 shows the description of the meta-dataset splits  
 165 according to the HPO-B-v3.

Search Space	ID	#HPs	Meta-Train		Meta-Validation		Meta-Test	
			#Evals.	#DS	#Evals.	#DS	#Evals.	#DS
rpart.preproc(16)	4796	3	10694	36	1198	4	1200	4
svm (6)	5527	8	385115	51	196213	6	354316	6
rpart (29)	5636	6	503439	54	184204	7	339301	6
rpart (31)	5859	6	58809	56	17248	7	21060	6
glmnet (4)	5860	2	3100	27	598	3	857	3
svm (7)	5891	8	44091	51	13008	6	17293	6
xgboost (4)	5906	16	2289	24	584	3	513	2
ranger (9)	5965	10	414678	60	73006	7	83597	7
ranger (5)	5970	2	68300	55	18511	7	19023	6
xgboost (6)	5971	16	44401	52	11492	6	19637	6
glmnet (11)	6766	2	599056	51	210298	6	310114	6
xgboost (9)	6767	18	491497	52	211498	7	299709	6
ranger (13)	6794	10	591831	52	230100	6	406145	6
ranger (15)	7607	9	18686	58	4203	7	5028	7
ranger (16)	7609	9	41631	59	8215	7	9689	7
ranger (7)	5889	6	1433	20	410	2	598	2

Table 3: Description of the search spaces in HPO-B-v3; "#HPs" stands for the number of hyperparameters, "#Evals." for the number of evaluations in a search space, while "#DS" for the number of datasets across which the evaluations are collected. The search spaces are named with the respective OpenML version number (in parenthesis), and their original names are preceded by *mlr.classif*.

## 166 5.2 Preprocessing

167 The OpenML-Python API [13] was used to download the experiment data from  
 168 OpenML [12]. We have collected all evaluations (referred to as runs in OpenML) tagged  
 169 with `Verified_Supervised_Classification` available until April 15, 2021.

170 We processed the raw data as follows. While the hyperparameter configuration was di-  
 171 rectly available for many evaluations, some of them had to be parsed from WEKA argu-  
 172 ments (e.g. `weka.filters.unsupervised.attribute.RandomProjection -P 16.0 -R 42`  
 173 `-D Sparse1`). A small percentage (<0.001%) of these were too complex in structure to be automati-  
 174 cally parsed, so they were discarded. Duplicate responses for the same hyperparameter configuration  
 175 have been resolved by keeping only one random response. Finally, all tasks with fewer than five  
 176 observations were also discarded.

177 All categorical hyperparameters were one-hot encoded, taking into account all categories that occur  
 178 in the different datasets for a search space. Missing values have been replaced with zeros and the  
 179 corresponding missing indicator (a new feature) has been set to one. Hyperparameters that had  
 180 the same value for all configurations in a search space were dropped. We manually decided which

181 hyperparameters required log-scaling by inspecting the distributions of each hyperparameter in each  
182 space (considerable manual effort). Finally, the hyperparameter ranges were scaled to  $[0, 1]$ .

### 183 5.3 Benchmark JSON schema

184 The benchmark is offered as easily accessible JSON files. The first-level key of each JSON schema  
185 corresponds to the search space ID, whereas the second-level key specifies the dataset ID. By  
186 accessing the JSON schema with the search space  $s$  and the dataset  $t$ , we obtain the meta-dataset  
187  $\mathcal{D}^{(s,t)} = \{(\mathbf{x}_i^{(s,t)}, y_i^{(s,t)})\}_i, \mathbf{x}_i^{(s,t)} \in \mathcal{X}^{(s)}$ . The meta-dataset exhibits the following structure, where  
188  $N$  denotes the number of evaluations available for the specific task:

189 `{search_space_ID: {dataset_ID: {X: [[x1], ..., [xN]], y: [[y1], ..., [yN]]}}}`

190 The initialization seeds are similarly provided as a JSON schema, where the third-level subschema  
191 has 5 keys whose values are the indices of the samples to use as initial configurations.

## 192 6 Recommended Experimental Protocol

193 One of the primary purposes of HPO-B is to standardize and facilitate the comparison between HPO  
194 techniques on *a level playing field*. In this section, we provide two specific recommendations: which  
195 benchmark to use for a type of algorithm and what metrics to use for comparing results.

196 **Evaluation Metrics** We define the average normalized regret at trial  $e$  (a.k.a. average distance  
197 to the minimum) as  $\min_{x \in \mathcal{X}_e^{(s,t)}} (f^{(s,t)}(x) - y_{\min}^*) / (y_{\max}^* - y_{\min}^*)$  with  $\mathcal{X}_e^{(s,t)}$  as the set of hyperpa-  
198 rameters that have been selected by a HPO method up to trial  $e$ , with  $y_{\min}^*$  and  $y_{\max}^*$  as the best and  
199 worst responses, respectively. The average rank represents the mean across tasks of the ranks of  
200 competing methods computed using the test accuracies of the best configuration until the  $e$ -th trial.  
201 Results across different search spaces are computed by a simple mean over the search-space-specific  
202 results.

203 **Non-Transfer Black-Box HPO** Methods should be compared on all the tasks in HPO-B-v2 and  
204 for each of the five initial configurations. The authors of future papers should report the normalized  
205 regret and the mean ranks for all trials from 1 to 100 (excluding the seeds). We recommend that the  
206 authors show both aggregated and per search-space (possibly moved to the appendix) mean regret  
207 and mean rank curves for trials ranging from 1 to 100. In other words, as many runs as the number of  
208 tasks for a given space times the number of initialization seeds. To assess the statistical significance  
209 of methods, we recommend that critical difference diagrams [6] be computed for the ranks of all runs  
210 @25, @50, and @100 trials.

211 **Transfer Black-Box HPO** Methods should be compared on the meta-data splits contained in HPO-  
212 B-v3. All competing methods should use exactly the evaluations of the provided meta-train datasets  
213 for meta- and transfer-learning their method, and tune the hyper-hyperparameters on the evaluations  
214 of the provided meta-validation datasets. In the end, the competing methods should be tested on the  
215 provided evaluations of the meta-test tasks. As our benchmark does not have pre-computed responses  
216 for all possible configurations in a space, the authors need to adapt their HPO acquisitions and  
217 suggest the next configuration only from the set of the pre-computed configurations for each specific  
218 meta-test task. Additionally, we recommend that the authors present (see details in the paragraph  
219 above) regret and rank plots, besides the critical difference diagrams @25, @50, and @100 trials. If  
220 a future transfer HPO method proposes a novel strategy for initializing configurations, for the sake of  
221 reproducibility we still recommend showing additional results with our initial configurations.

## 222 7 Experimental Results

223 The benchmark is intended to serve as a new standard for evaluating non-transfer and transfer  
224 black-box HPO methods. In the following, we will compare different methods according to our  
225 recommended protocol described in Section 6. This is intended to demonstrate the usefulness of our  
226 meta-dataset, while at the same time serving as an example for the aforementioned recommendations  
227 on comparing baselines and presenting results.

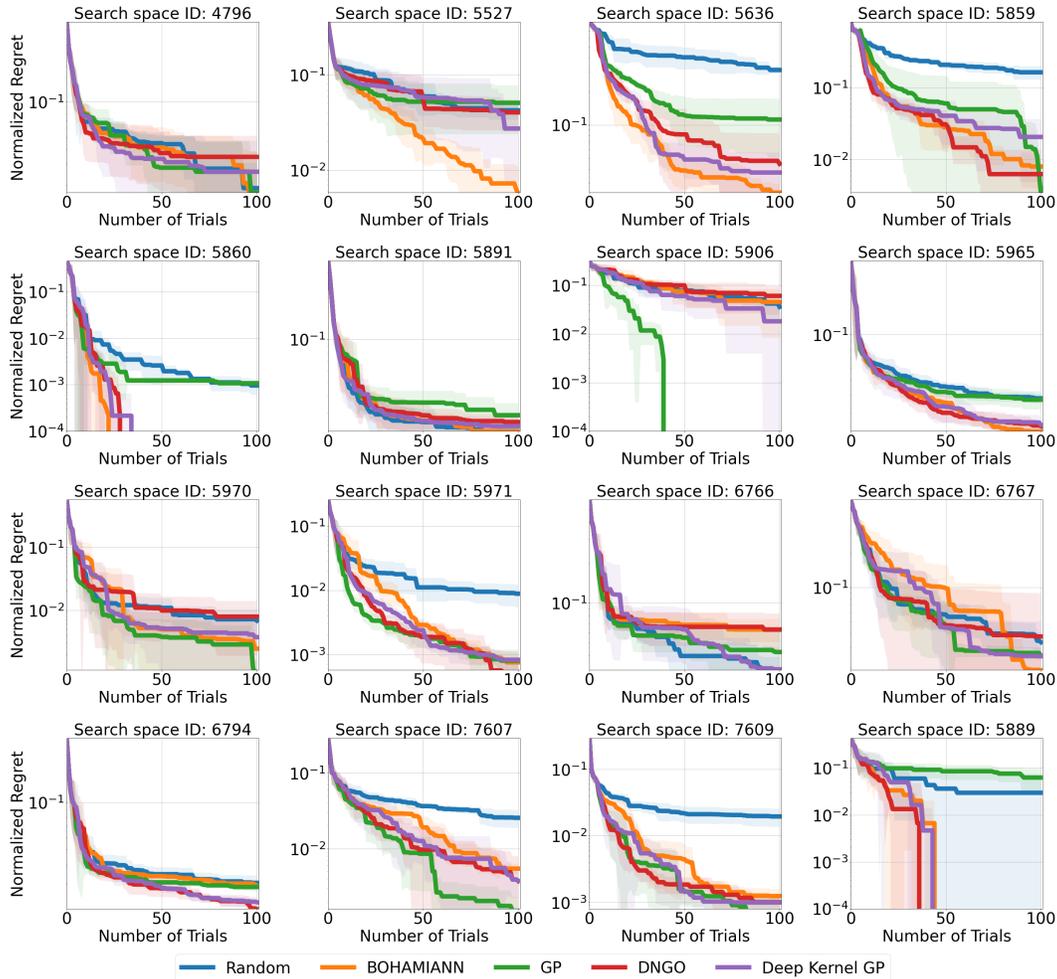


Figure 1: **Normalized regret** comparison of **non-transfer** black-box HPO methods on HPO-B-v2

## 228 7.1 Non-transfer Black-Box HPO

229 First, we compare Random Search, DNGO, BOHAMIANN, Gaussian Process (GP) with Matérn 3/2  
 230 kernel, and Deep Gaussian Process (FSBO [31] without pre-training) on HPO-B-v2 in the non-transfer  
 231 scenario. As recommended by us earlier, in Figure 1 we report aggregated results for normalized  
 232 regret, average rank, and critical difference plots. In addition, we report in Figure 2 the aggregated  
 233 normalized regret per search space. The values in the figures for the number of trials equal to 0  
 234 correspond to the result after the five initialization steps. According to Figure 2, BOHAMIANN  
 235 and Deep GP achieve comparable aggregated normalized regret across all search spaces, which  
 236 suggests that both methods are equally well-suited for the tasks. The average rank and the critical  
 237 difference plot paint a different picture, in which Deep GP and DNGO achieve better results. This  
 238 discrepancy arises because each metric measures different performance aspects on different tasks,  
 239 so it’s important to report both. As can be seen in Figure 10, Deep GP achieves better results than  
 240 the GP in most of the tasks, which leads to a better average ranking. However, as we can see in  
 241 Figure 1, the regrets are observed at heterogeneous scales that can skew the overall averages. In some  
 242 cases where BOHAMIANN outperforms Deep GP (e.g. search spaces 5527, 5859, and 5636), the  
 243 difference in normalized regret is evident, due to the nature of the search space, whereas in cases  
 244 where it is the other way around, however, the difference is only slightly less evident (e.g. search  
 245 spaces 4796, 5906, and 7609). An important aspect of HPO is the choice of the surrogate function  
 246 and acquisition. Figure 3 presents an ablation of typical combinations and shows the accuracy of the  
 247 Boosted Tree as a surrogate.

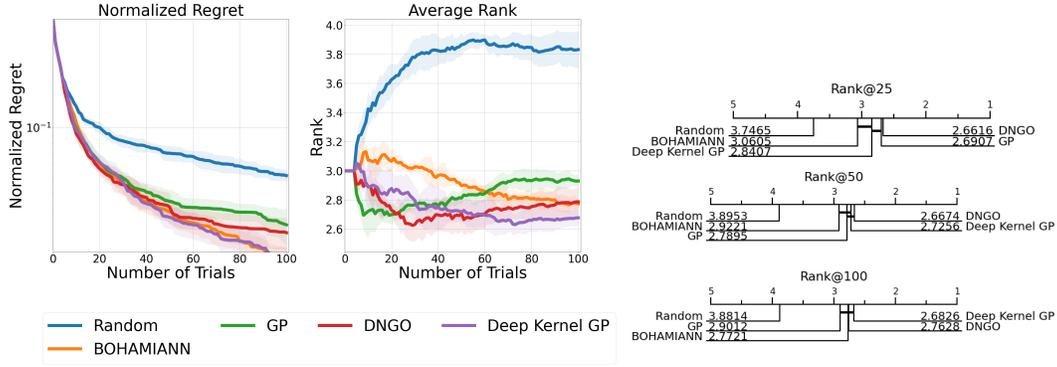


Figure 2: **Aggregated** comparisons of normalized regret and mean ranks across all search spaces for the **non-transfer** HPO methods on HPO-B-v2

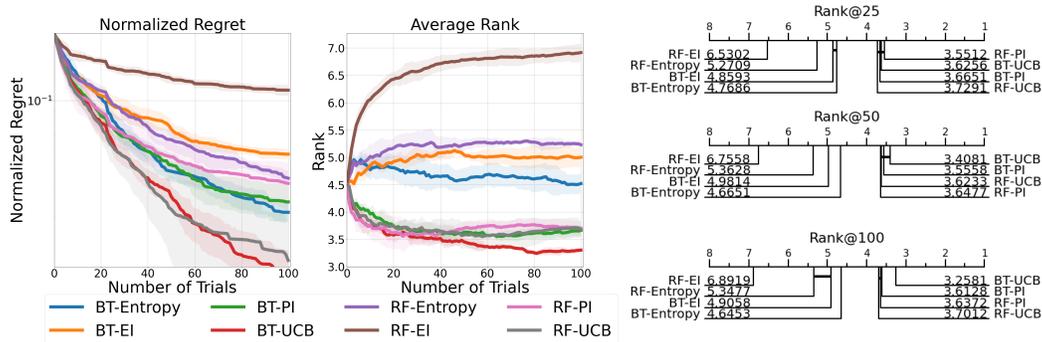


Figure 3: **Aggregated** comparisons of different surrogates and acquisition functions for **transfer** HPO methods on HPO-B-v2; BT stands for Boosted Trees, RF for Random Forests, EI for Expected Improvement, and UCB for Upper Confidence Bound.

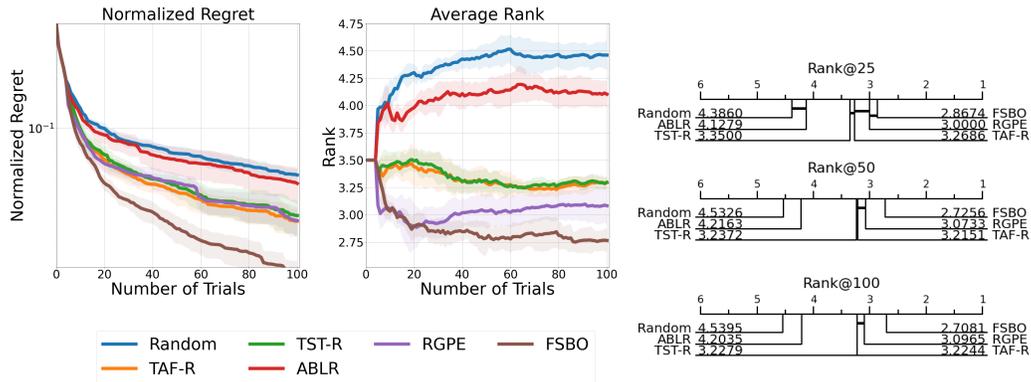


Figure 4: **Aggregated** comparisons of normalized regret and mean ranks across all search spaces for the **transfer learning** HPO methods on HPO-B-v3

## 248 7.2 Transfer Black-Box HPO

249 Finally, we compare RGPE [10], ABLR [20], TST-R [34], TAF-R [35], and FSBO [31] on HPO-B-v3  
 250 in the transfer scenario. All hyper-hyperparameters were optimized on the meta-validation datasets  
 251 and we report results aggregated across all test search spaces in terms of normalized regret and  
 252 average rank in Figure 4. The results per search space for normalized regret and average rank are  
 253 given in Figure 5 and Figure 11, respectively. FSBO shows improvements over all the compared  
 254 methods for the normalized regret metric and average rank metric. On the other hand, RGPE is  
 255 seemingly performing similar to TST-R and TAF-R for the average regret, but performs significantly

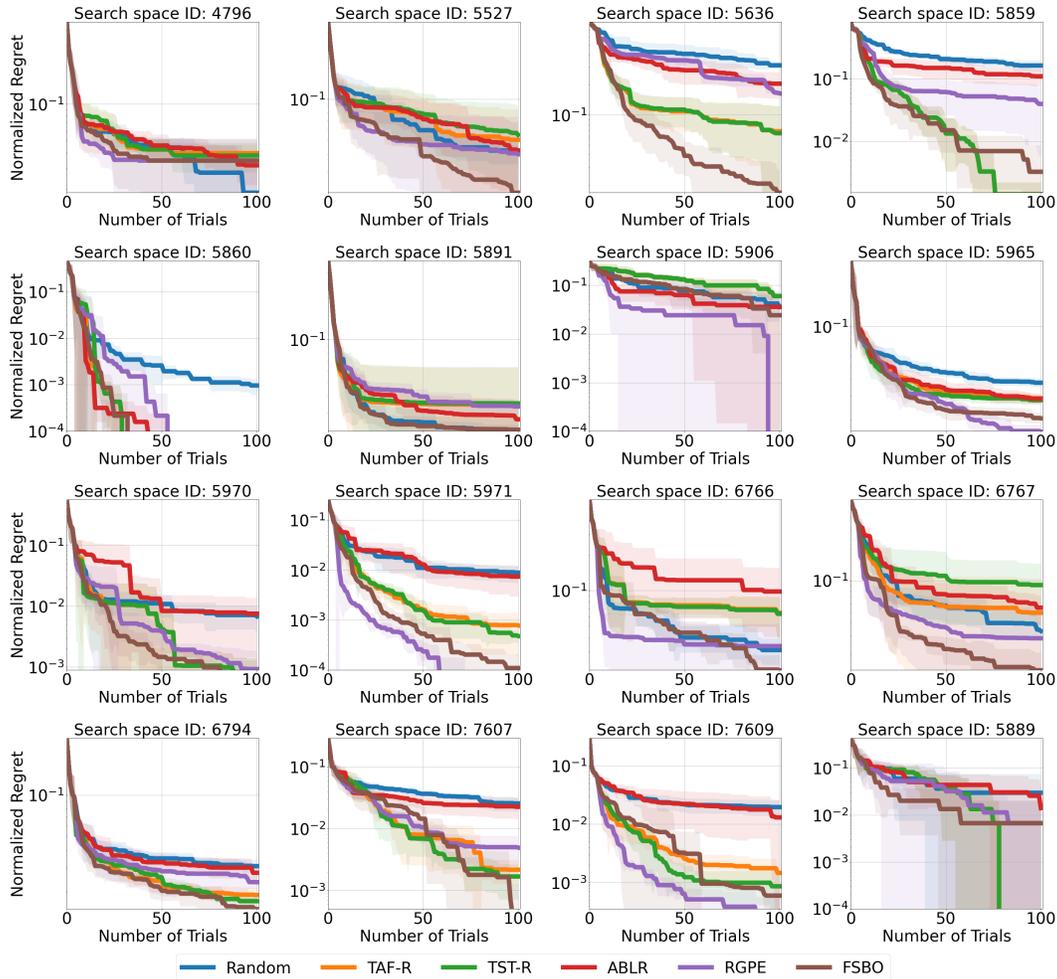


Figure 5: **Normalized regret** comparison of **transfer learning HPO** methods on HPO-B-v3

256 better for the average rank metric. The explanation is the same as for our last experiment and can  
 257 mainly be traced back to the strong performance of RGPE in search spaces 5971 and 5906. Such  
 258 behaviors strengthen our recommendations of Section 6 for showing results in terms of both the ranks  
 259 and the normalized regrets, as well as the ranks’ statistical significance.

### 260 7.3 Comparing Non-Transfer vs. Transfer Black-Box HPO

261 We provide a cumulative comparison of both non-transfer and transfer black-box methods in Figure 6,  
 262 for demonstrating the benefit of transfer learning in HPO-B-v3. We see that the transfer methods  
 263 (FSBO, RGPE, TST-R, TAF-R) achieve significantly better performances than the non-transfer tech-  
 264 niques (GP, DNGO, BOHAMIANN, Deep Kernel GP). On the average rank plot and the associated  
 265 Critical Difference diagrams, we notice that FSBO [31] achieves significantly better results than all  
 266 baselines, followed by RGPE [10]. A detailed comparison of the ranks per search-space is presented  
 267 in the supplementary material. In particular, the direct gain of transfer learning can be observed by the  
 268 dominance that FSBO has over *Deep Kernel GP*, considering that both use exactly the same surrogate  
 269 model and the same acquisition function. In comparison, the deep kernel parameters in FSBO are  
 270 initialized from the solution of a meta-learning optimization conducted on the meta-train tasks of  
 271 HPO-B-v3 (transfer), while the parameters of *Deep Kernel GP* are initialized randomly (no transfer).

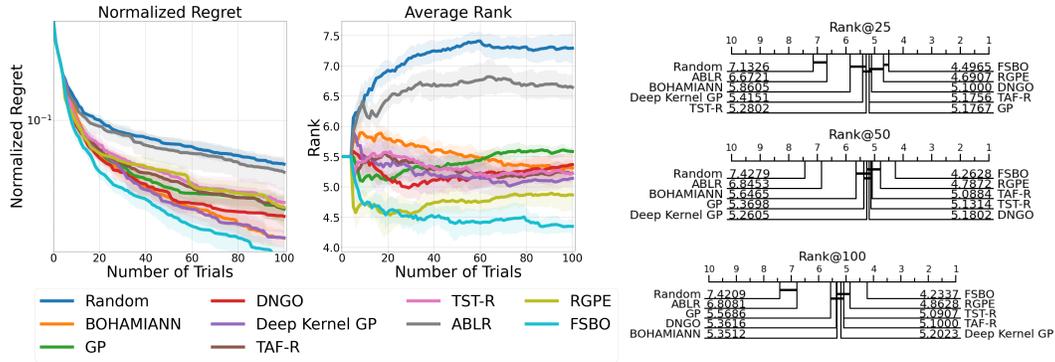


Figure 6: Comparisons of normalized regret and mean ranks across all search spaces for the **transfer learning** and **non-transfer** HPO methods on HPO-B-v3

## 272 8 Discussing the Limitations of HPO-B

273 HPO-B relies on OpenML evaluations that are not exhaustively computed for all the possible  
 274 hyperparameter configurations of a search space. In that context, an acquisition function can suggest  
 275 a next configuration only from the set of those configurations that were evaluated on a particular task.  
 276 A possible way to tackle this limitation is fitting surrogate models on the evaluated configurations,  
 277 and using the estimation of the surrogate as the response of a new configuration for which no actual  
 278 evaluation exists. However, the choice of the surrogate model might add noisy/confounding effects  
 279 to the evaluations, as there are open questions on the capacity of the surrogate model (i.e. do we  
 280 need different surrogate complexities for different tasks/datasets?), or the choice of the loss function  
 281 for training the surrogate. Moreover, as the majority of HPO-B tasks have an abundant number of  
 282 evaluations (see Appendix E for details), it is highly likely that a well-performing configuration is  
 283 already present in the set of existing evaluations.

284 Another limitation of HPO-B is that it only covers black-box HPO tasks, instead of other HPO  
 285 problems, such as grey-box HPO, or pipeline optimization for AutoML libraries. In addition, HPO-B  
 286 is restricted by the nature of search spaces found in OpenML, which contains evaluations for well-  
 287 established machine learning algorithms for tabular data, but lacks state-of-the-art deep learning  
 288 methods or tasks on image or text data.

## 289 9 Conclusions

290 Recent HPO and transfer-learning HPO papers inconsistently use different meta-datasets, arbitrary  
 291 train/validation/test splits, as well as ad-hoc preprocessing, which makes it hard to reproduce the  
 292 published results. To resolve this bottleneck, we propose HPO-B, a novel benchmark based on the  
 293 OpenML repository, that contains meta-datasets from 176 search spaces, 196 datasets, and a total of  
 294 6.4 million evaluations. For promoting reproducibility at a *level playing field* we also provide initial  
 295 configuration seeds, as well as predefined training, validation and testing splits. Our benchmark  
 296 contains pre-processed meta-datasets and a clear set of HPO tasks and exact splits, therefore, it  
 297 enables future benchmark results to be directly comparable. We believe our benchmark has the  
 298 potential to become the *de facto* standard for experimentation in the realm of black-box HPO.

## 299 Acknowledgements

300 The research of Hadi S. Jomaa is co-funded by the industry project "IIP-Ecosphere: Next Level  
 301 Ecosphere for Intelligent Industrial Production". Prof. Grabocka is also thankful to the Eva Mayr-  
 302 Stihl Foundation for their generous research grant, as well as to the Ministry of Science, Research  
 303 and the Arts of the German state of Baden-Württemberg for funding his professorship. In addition,  
 304 we thank Arlind Kadra for his assistance in interfacing with the OpenML Python package.

## References

- 305
- 306 [1] Rémi Bardenet, Mátyás Brendel, Balázs Kégl, and Michele Sebag. Collaborative hyperparameter  
307 tuning. In *International conference on machine learning*, pages 199–207, 2013.
- 308 [2] James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-  
309 parameter optimization. In *25th annual conference on neural information processing systems*  
310 (*NIPS 2011*), volume 24. Neural Information Processing Systems Foundation, 2011.
- 311 [3] James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *Journal*  
312 *of machine learning research*, 13(2), 2012.
- 313 [4] Martin Binder, Florian Pfisterer, and Bernd Bischl. Collecting empirical data about hyperparam-  
314 eters for datadriven automl. In *ICML Workshop on Automated Machine Learning*, 2020.
- 315 [5] Alexander I Cowen-Rivers, Wenlong Lyu, Zhi Wang, Rasul Tutunov, Hao Jianye, Jun Wang, and  
316 Haitham Bou Ammar. Hebo: Heteroscedastic evolutionary bayesian optimisation. *arXiv preprint*  
317 *arXiv:2012.03826*, 2020. winning submission to the NeurIPS 2020 Black Box Optimisation  
318 Challenge.
- 319 [6] Janez Demsar. Statistical comparisons of classifiers over multiple data sets. *J. Mach. Learn.*  
320 *Res.*, 7:1–30, 2006.
- 321 [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-  
322 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern*  
323 *recognition*, pages 248–255. Ieee, 2009.
- 324 [8] Xuanyi Dong and Yi Yang. Nas-bench-201: Extending the scope of reproducible neural  
325 architecture search. In *International Conference on Learning Representations*, 2020.
- 326 [9] Katharina Eggenberger, Matthias Feurer, Frank Hutter, James Bergstra, Jasper Snoek, Holger  
327 Hoos, and Kevin Leyton-Brown. Towards an empirical foundation for assessing bayesian  
328 optimization of hyperparameters. In *NIPS workshop on Bayesian Optimization in Theory and*  
329 *Practice*, volume 10, page 3, 2013.
- 330 [10] Matthias Feurer, Benjamin Letham, and Eytan Bakshy. Scalable meta-learning for bayesian  
331 optimization. *arXiv preprint arXiv:1802.02219*, 2018.
- 332 [11] Matthias Feurer, Jost Tobias Springenberg, and Frank Hutter. Initializing Bayesian hyperparam-  
333 eter optimization via meta-learning. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*,  
334 2015.
- 335 [12] Matthias Feurer, Jan N. van Rijn, Arlind Kadra, Pieter Gijsbers, Neeratyoy Mallik, Sahithya  
336 Ravi, Andreas Mueller, Joaquin Vanschoren, and Frank Hutter. Openml-python: an extensible  
337 python api for openml. *arXiv*, 1911.02490, 2020.
- 338 [13] Matthias Feurer, Jan N. van Rijn, Arlind Kadra, Pieter Gijsbers, Neeratyoy Mallik, Sahithya  
339 Ravi, Andreas Müller, Joaquin Vanschoren, and Frank Hutter. Openml-python: an extensible  
340 python api for openml. *arXiv:1911.02490*, 2019.
- 341 [14] Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. Sequential model-based optimization  
342 for general algorithm configuration. In *International conference on learning and intelligent*  
343 *optimization*, pages 507–523. Springer, 2011.
- 344 [15] Hadi S Jomaa, Josif Grabocka, and Lars Schmidt-Thieme. Hyp-rl: Hyperparameter optimization  
345 by reinforcement learning. *arXiv preprint arXiv:1906.11527*, 2019.
- 346 [16] Hadi S Jomaa, Lars Schmidt-Thieme, and Josif Grabocka. Dataset2vec: Learning dataset  
347 meta-features. *Data Mining and Knowledge Discovery*, 35(3):964–985, 2021.
- 348 [17] Donald R. Jones, Matthias Schonlau, and William J. Welch. Efficient global optimization of  
349 expensive black-box functions. *J. Global Optimization*, 13(4):455–492, 1998.
- 350 [18] Daniel Kühn, Philipp Probst, Janek Thomas, and Bernd Bischl. Automatic exploration of  
351 machine learning experiments on openml. *CoRR*, abs/1806.10961, 2018.

- 352 [19] Lisha Li, Kevin Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. Hyper-  
353 band: A novel bandit-based approach to hyperparameter optimization. *The Journal of Machine*  
354 *Learning Research*, 18(1):6765–6816, 2017.
- 355 [20] Valerio Perrone, Rodolphe Jenatton, Matthias W Seeger, and Cédric Archambeau. Scalable  
356 hyperparameter transfer learning. In *Advances in Neural Information Processing Systems*, pages  
357 6845–6855, 2018.
- 358 [21] Valerio Perrone and Huibin Shen. Learning search spaces for bayesian optimization: Another  
359 view of hyperparameter transfer learning. In *Advances in Neural Information Processing*  
360 *Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS*  
361 *2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 12751–12761, 2019.
- 362 [22] Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian processes for machine*  
363 *learning*. Adaptive computation and machine learning. MIT Press, 2006.
- 364 [23] David Salinas, Huibin Shen, and Valerio Perrone. A quantile-based approach for hyperparameter  
365 transfer learning. In *International Conference on Machine Learning*, pages 8438–8448. PMLR,  
366 2020.
- 367 [24] Nicolas Schilling, Martin Wistuba, Lucas Drumond, and Lars Schmidt-Thieme. Hyperparameter  
368 optimization with factorized multilayer perceptrons. In *Joint European Conference on Machine*  
369 *Learning and Knowledge Discovery in Databases*, pages 87–103. Springer, 2015.
- 370 [25] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical bayesian optimization of machine  
371 learning algorithms. *arXiv preprint arXiv:1206.2944*, 2012.
- 372 [26] Jasper Snoek, Oren Rippel, Kevin Swersky, Ryan Kiros, Nadathur Satish, Narayanan Sundaram,  
373 Mostofa Patwary, Mr Prabhat, and Ryan Adams. Scalable bayesian optimization using deep  
374 neural networks. In *International conference on machine learning*, pages 2171–2180. PMLR,  
375 2015.
- 376 [27] Jost Tobias Springenberg, Aaron Klein, Stefan Falkner, and Frank Hutter. Bayesian optimization  
377 with robust bayesian neural networks. In *Proceedings of the 30th International Conference on*  
378 *Neural Information Processing Systems*, pages 4141–4149, 2016.
- 379 [28] Joaquin Vanschoren, Jan N. van Rijn, Bernd Bischl, and Luis Torgo. Openml: Networked  
380 science in machine learning. *SIGKDD Explorations*, 15(2):49–60, 2013.
- 381 [29] Michael Volpp, Lukas P Fröhlich, Kirsten Fischer, Andreas Doerr, Stefan Falkner, Frank Hutter,  
382 and Christian Daniel. Meta-learning acquisition functions for transfer learning in bayesian  
383 optimization. *arXiv preprint arXiv:1904.02642*, 2019.
- 384 [30] Michael Volpp, Lukas P. Fröhlich, Kirsten Fischer, Andreas Doerr, Stefan Falkner, Frank Hutter,  
385 and Christian Daniel. Meta-learning acquisition functions for transfer learning in bayesian  
386 optimization, 2020.
- 387 [31] Martin Wistuba and Josif Grabocka. Few-shot bayesian optimization with deep kernel surrogates.  
388 In *International Conference on Learning Representations*, 2021.
- 389 [32] Martin Wistuba, Nicolas Schilling, and Lars Schmidt-Thieme. Hyperparameter search space  
390 pruning—a new component for sequential model-based hyperparameter optimization. In *Joint*  
391 *European Conference on Machine Learning and Knowledge Discovery in Databases*, pages  
392 104–119. Springer, 2015.
- 393 [33] Martin Wistuba, Nicolas Schilling, and Lars Schmidt-Thieme. Learning hyperparameter  
394 optimization initializations. In *International Conference on Data Science and Advanced*  
395 *Analytics*, pages 1–10, 2015.
- 396 [34] Martin Wistuba, Nicolas Schilling, and Lars Schmidt-Thieme. Two-stage transfer surrogate  
397 model for automatic hyperparameter optimization. In *Joint European conference on machine*  
398 *learning and knowledge discovery in databases*, pages 199–214. Springer, 2016.

- 399 [35] Martin Wistuba, Nicolas Schilling, and Lars Schmidt-Thieme. Scalable gaussian process-based  
400 transfer surrogates for hyperparameter optimization. *Machine Learning*, 107(1):43–78, 2018.
- 401 [36] Chris Ying, Aaron Klein, Eric Christiansen, Esteban Real, Kevin Murphy, and Frank Hutter.  
402 NAS-bench-101: Towards reproducible neural architecture search. In Kamalika Chaudhuri and  
403 Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine*  
404 *Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 7105–7114. PMLR,  
405 09–15 Jun 2019.

## 406 Checklist

407 The checklist follows the references. Please read the checklist guidelines carefully for information on  
408 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or  
409 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing  
410 the appropriate section of your paper or providing a brief inline description. For example:

- 411 • Did you include the license to the code and datasets? **[Yes]** See Section
- 412 • Did you include the license to the code and datasets? **[No]** The code and the data are  
413 proprietary.
- 414 • Did you include the license to the code and datasets? **[N/A]**

415 Please do not modify the questions and only use the provided macros for your answers. Note that the  
416 Checklist section does not count towards the page limit. In your paper, please delete this instructions  
417 block and only keep the Checklist section heading above along with the questions/answers below.

- 418 1. For all authors...
  - 419 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
420 contributions and scope? **[Yes]**
  - 421 (b) Did you describe the limitations of your work? **[Yes]**
  - 422 (c) Did you discuss any potential negative social impacts of your work? **[N/A]**
  - 423 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
424 them? **[Yes]**
- 425 2. If you are including theoretical results...
  - 426 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
  - 427 (b) Did you include complete proofs of all theoretical results? **[N/A]**
- 428 3. If you ran experiments (e.g. for benchmarks)...
  - 429 (a) Did you include the code, data, and instructions needed to reproduce the main ex-  
430 perimental results (either in the supplemental material or as a URL)? **[Yes]** See our  
431 repository link in the introduction.
  - 432 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
433 were chosen)? **[No]** We used pre-defined configurations from previous work.
  - 434 (c) Did you report error bars (e.g., concerning the random seed after running experiments  
435 multiple times)? **[Yes]**
  - 436 (d) Did you include the total amount of computing and the type of resources used (e.g.,  
437 type of GPUs, internal cluster, or cloud provider)? **[No]**
- 438 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - 439 (a) If your work uses existing assets, did you cite the creators? **[Yes]**
  - 440 (b) Did you mention the license of the assets? **[No]** They are included in the cited  
441 publications.
  - 442 (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]**
  - 443 (d) Did you discuss whether and how consent was obtained from people whose data you're  
444 using/curating? **[No]** We are using open-sourced assets.
  - 445 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
446 information or offensive content? **[N/A]**
- 447 5. If you used crowdsourcing or researched with human subjects...
  - 448 (a) Did you include the full text of instructions given to participants and screenshots, if  
449 applicable? **[N/A]**
  - 450 (b) Did you describe any potential participant risks, with links to Institutional Review  
451 Board (IRB) approvals, if applicable? **[N/A]**
  - 452 (c) Did you include the estimated hourly wage paid to participants and the total amount  
453 spent on participant compensation? **[N/A]**