

Automated Benchmark-Driven Design and Explanation of Hyperparameter Optimizers

Julia Moosbauer^{1b}, Martin Binder, Lennart Schneider^{1b}, Florian Pfisterer^{1b}, Marc Becker, Michel Lang, Lars Kotthoff^{1b}, and Bernd Bischl^{1b}

Abstract—Automated hyperparameter optimization (HPO) has gained great popularity and is an important component of most automated machine learning frameworks. However, the process of designing HPO algorithms is still an unsystematic and manual process: new algorithms are often built on top of prior work, where limitations are identified and improvements are proposed. Even though this approach is guided by expert knowledge, it is still somewhat arbitrary. The process rarely allows for gaining a holistic understanding of which algorithmic components drive performance and carries the risk of overlooking good algorithmic design choices. We present a principled approach to automated benchmark-driven algorithm design applied to multifidelity HPO (MF-HPO). First, we formalize a rich space of MF-HPO candidates that includes, but is not limited to, common existing HPO algorithms and then present a configurable framework covering this space. To find the best candidate automatically and systematically, we follow a programming-by-optimization approach and search over the space of algorithm candidates via Bayesian optimization. We challenge whether the found design choices are necessary or could be replaced by more naive and simpler ones by performing an ablation analysis. We observe that using a relatively simple configuration (in some ways, simpler than established methods) performs very well as long as some critical configuration parameters are set to the right value.

Index Terms—Algorithm analysis, algorithm design, automated machine learning (AutoML), hyperparameter optimization (HPO), multifidelity.

Manuscript received 21 September 2021; revised 15 March 2022 and 28 July 2022; accepted 15 September 2022. Date of publication 6 October 2022; date of current version 1 December 2022. The work of Lennart Schneider was supported by the Bavarian Ministry of Economic Affairs, Regional Development and Energy through the Center for Analytics-Data-Applications (ADACenter) within the Framework of BAYERN DIGITAL II under Grant 20-3410-2-9-8. The work of Michel Lang was supported in part by the Research Center Trustworthy Data Science and Security (<https://rc-trust.ai>). The work of Lars Kotthoff was supported by NSF under Grant 1813537. (Julia Moosbauer and Martin Binder contributed equally to this work.) (Corresponding author: Julia Moosbauer.)

Julia Moosbauer, Martin Binder, Lennart Schneider, Florian Pfisterer, Marc Becker, and Bernd Bischl are with the Department of Statistics, Ludwig-Maximilians-Universität München, 80539 Munich, Germany (e-mail: julia.moosbauer@stat.uni-muenchen.de; martin.binder@stat.uni-muenchen.de; lennart.schneider@stat.uni-muenchen.de; florian.pfisterer@stat.uni-muenchen.de; marc.becker@stat.uni-muenchen.de; bernd.bischl@stat.uni-muenchen.de).

Michel Lang is with the Research Center Trustworthy Data Science and Security, 44227 Dortmund, Germany (e-mail: michel.lang@tu-dortmund.de).

Lars Kotthoff is with the Department of Computer Science, University of Wyoming, Laramie, WY 82071 USA (e-mail: larsko@uwyo.edu).

This article has supplementary downloadable material available at <https://doi.org/10.1109/TEVC.2022.3211336>, provided by the authors.

Digital Object Identifier 10.1109/TEVC.2022.3211336

I. INTRODUCTION

MACHINE learning (ML) is, in many regards, an optimization problem, and many ML methods can be expressed as algorithms that perform loss minimization with respect to a given objective function. The higher-level task of selecting the ML method and its configuration is often framed as an optimization problem as well, sometimes referred to as a hyperparameter optimization (HPO) [1] or combined algorithm selection and HPO (CASH) problem [2]. Successfully addressing this problem can lead to large performance gains compared to simply using defaults, and in the context of automated ML (AutoML), the use of HPO can make ML more accessible to nonexperts. Because of their potential benefits to ML performance and usability, it is of particular interest to design optimization algorithms that perform particularly well on the HPO problem.

Optimization problems arise in many fields of science and engineering, but as the no-free-lunch theorem states, there is no one optimization algorithm that solves all problems equally well [3]. To design suitable optimizers, it is therefore important to understand the characteristics of HPO.

- 1) *Black-Box*: The objective usually provides no analytical information [4], such as a gradient. Thus, the application of many traditional optimization methods, such as BFGS, is rendered inappropriate or at least questionable.
- 2) *Complex Search Space*: The search space of the optimization problem is often high-dimensional and may contain continuous, integer-valued, and categorical dimensions. Often, there are dependencies between dimensions or even specific hyperparameter values [5].
- 3) *Expensive*: A single evaluation of the objective function may take hours or days. Thus, the total number of possible function evaluations is often severely limited [4].
- 4) *Low-Fidelity Approximations Possible*: An approximation of the true objective value at lower expense can often be obtained, for example, through a partial evaluation [6].
- 5) *Low Effective Dimensionality*: The landscape of the objective function can usually be approximated well by a function of a small subset of all dimensions [7].

Recent HPO and AutoML research has focused on finding and improving optimization algorithms that work particularly well under these conditions. A common approach is to tackle HPO by estimating a local or global structure of the objective landscape by some form of the predictive model.

This introduces additional overhead and complexity with the aim of reducing the overall number of expensive objective evaluations necessary to find an approximate optimum. Typical representatives of this approach are Bayesian optimization (BO) [8] algorithms and frameworks based on BO, which are global optimization schemes based on a nonlinear regression model, e.g., a Gaussian process or random forest. They have shown significant improvements in performance compared to other methods [9] but carry a significant overhead. Furthermore, BO is somewhat difficult to parallelize due to its sequential nature, although many variants exist (e.g., [10], [11], [12], and [13]).

Multifidelity HPO (MF-HPO) algorithms aim to accelerate the optimization process by exploiting cheaper proxy functions of the objective function itself (e.g., by training ML models on a smaller subsample of the available training data, or by running fewer training iterations). Bandit-based algorithms like Hyperband (HB) [14] have become particularly popular because of their good tradeoff between optimization performance and simplicity.

Progress in the field of HPO often consists of iterative improvements of established algorithms. Considerable work exists, for example, to improve the limitations of HB: asynchronous successive halving (ASHA) [15] proposes a sophisticated way to make efficient use of parallel resources, BO HB (BOHB) [16] improves performance during later parts of a run by incorporating surrogate assistance into HB, and asynchronous BOHB (A-BOHB) [17] unites a bandit-based optimization scheme using model-based guidance with asynchronous parallelization.

While these conceptual extensions of HPO all have their respective merit, it is often somewhat overlooked that the simplicity of an optimization algorithm (i.e., how difficult modifications and extensions are, and on how many dependencies a system relies [18]) heavily influences its adoption in practice. Random search (RS), for example, still enjoys great popularity, as it is extremely simple to implement and parallelize, has almost no overhead, and is able to take advantage of the aforementioned low effective dimensionality [7]. Furthermore, algorithmic developments identify and address limitations of prior research, but rarely question core algorithmic choices that have been made in the original implementation. Many multifidelity algorithms, for example, are extensions and further developments of HB that take the fixed successive halving (SH) schedule [19] for granted. The process of designing a good MF-HPO optimizer in practice—and many other algorithmic solutions in science in general—can therefore often feel somewhat like a “manual stochastic local search on the meta level.” The drawback of this manual procedure is that the design space of all HPO algorithms is not systematically searched, and parts of the design space are excluded by prior algorithmic decisions. If “established” algorithms are not challenged, there is a risk that algorithms that work well will be overlooked, and it is often hard to identify what algorithmic components make a difference. In particular, it is possible that overly complicated algorithms are developed by extending “established” designs, only some of which contribute meaningfully to performance

gains. Sometimes certain technical components of an algorithm, which are neither exposed nor discussed in detail, may also influence performance significantly.

A. Contributions

We make a principled demonstration of how HPO algorithm design can be performed systematically and automatically with a benchmark-driven approach following the programming-by-optimization paradigm [20]. In particular, the contributions of this work are as follows.

- 1) *Formalization*: We formalize the design space of MF-HPO algorithms and demonstrate that established MF-HPO algorithms represent instances within this space.
- 2) *Framework*: Based on this formalization, we present a rich, configurable framework for MF-HPO algorithms, whose software implementation we call surrogate model-assisted HB (SMASHY).
- 3) *Configuration*: Based on the formalization and framework, we follow an empirical approach to design an MF-HPO algorithm by optimization, given a large benchmark suite. This configuration procedure does not only consider performance but also, e.g., the simplicity of the design.
- 4) *Benchmark*: As in general any HPO algorithm will be applied in a diverse set of application scenarios, we evaluate the performance of our newly designed algorithm on a representative set of problems that were not previously used for its configuration (i.e., a clean test-set approach on the meta-level) and compare them with established implementations of HPO methods.
- 5) *Explanation*: For the resulting MF-HPO system, we systematically assess and explain the effect of different design choices on overall algorithmic performance. Furthermore, we investigate the behavior of algorithmic design components in the context of specific problem scenarios; i.e., we investigate which algorithmic components lead to performance improvements for simple HPO with numeric hyperparameters, AutoML pipeline configuration, and neural architecture search.

II. RELATED WORK

HPO is one of the most essential components of current AutoML methods [1], and MF-HPO has recently become more prominent, given that cheap, low-fidelity evaluations have proven useful to speed up optimization, especially for expensive HPO of complex ML algorithms on larger data sets [14]. While AutoML tools have historically relied on a limited set of HPO methods, we argue that the optimal HPO method depends on problem characteristics, and therefore a systematic development of HPO methods under consideration of problem characteristics is required. Approaches toward such systematic development have often relied on a *high-level* language or template that allows expressing solution algorithms for a given problem class, e.g., to solve constraint satisfaction problems [21], [22], [23], satisfiability problems [24], scheduling problems [25], or general multiobjective combinatorial problems [26], [27].

Even if a high-level language is available, manual configuration of such frameworks is laborious and requires expert knowledge. This motivates the design philosophy of “programming by optimization” [20] (PBO), which advocates for allowing algorithmic choices in a software system (instead of fixing them at the time of implementation) and automatic configuration by optimization for a given problem context.

As one approach to automatic and efficient algorithm configuration, racing-based strategies have been used to design optimization algorithms. For example, iterated F-RACE [28] has been used for the automatic design of multiobjective ant colony optimization algorithms [26]. Similarly, IRACE [29] has been used for the automatic design multiobjective evolutionary algorithms [27] or to meta-configure the parameters IRACE itself [30]. Another commonly used framework is SMAC [5], which extends the sequential model-based optimization paradigm (SMBO, see also Section IV-A2) to an algorithm configuration setting. This is achieved through the use of an intensification procedure that governs across how many problem instances each configuration is evaluated, trading off computational cost against confidence regarding the superiority of a given configuration. While such intensification mechanisms have been used in other work before [31], [32], SMAC also uses instance features describing properties of a problem instance are used to train the empirical performance model predicting the performance of a configuration on a new problem instance. Besides racing and sequential model-based approaches, genetic algorithms have also been used to evolve optimal solvers [33].

We argue that the design of HPO algorithms can be seen as an instance of PBO. However, while there are many approaches that focus on individual algorithmic choices (e.g., the choice of a surrogate model for BO [34]), we are not aware of many cases where PBO is applied to designing HPO systems themselves. One exception is [35], who use SMACv3 [36] to automatically configure BO for HPO from a flexible search space of components. We take a similar approach here in that the algorithmic choices are exposed as hyperparameters that can be tuned. However, unlike [35], we do not configure an established HPO method (such as BO) with a predefined structure and associated control parameters (e.g., varying the surrogate model of BO). Instead, we introduce a *new* configurable algorithmic framework, which covers many different MF-HPO structures, including well-established principles for multifidelity handling (e.g., SH) as well as new approaches (e.g., equal batch size in all proposals).

In addition to designing well-performing algorithms, it is equally important to facilitate an understanding of the effects of all considered design choices. The field of *sensitivity analysis* (SA) comprises a multitude of methods to assess the importance of input factors on the output of a mathematical model [37]. Functional ANOVA (fANOVA) methods, which decompose the response of a (mathematical) model or function into lower-order components, are a widely studied method in the field of SA, dating back to [38]. This class of methods has also become popular in the field of ML to analyze the importance of hyperparameters [39].

Popular ways of analyzing effects of algorithmic effects in ML and algorithm configuration are *ablation studies* [40]. This involves measuring the performance when removing one or more of algorithmic subcomponents to understand the relative contribution of the ablated components to overall performance. There are different ways of performing an ablation analysis; probably the most common approach is *leave-one-component-out* (LOCO) ablation [41]. In the context of algorithm configuration, Fawcett and Hoos [40] proposed an ablation approach that links a source configuration (e.g., the default) to a target (e.g., the optimized configuration) through an ablation path.

Nevertheless, many existing works that propose or improve HPO or algorithm configuration systems do not analyze the algorithmic choices of an optimized system, and the ones that do perform relatively straightforward analyses. For example, Minton [21] compared the designs and their approach finds automatically to the designs expert humans generated. López-Ibáñez and Stützle [42] performed ANOVA and non-parametric Friedman tests to investigate in detail the effects that algorithmic choices, found through automatic configuration [26], have on the performance of multiobjective ant colony optimization algorithms. de Nobel et al. [43] presented a modular framework for CMA-ES variants on which they perform optimization; in particular, they investigate how the optimized configuration changes when the search space is enlarged by introducing new components.

III. METHODOLOGY

A. Supervised Machine Learning

Supervised ML typically deals with a dataset (which is, mathematically speaking, a tuple) $\mathcal{D} = ((\mathbf{x}^{(i)}, y^{(i)})) \in (\mathcal{X} \times \mathcal{Y})^n$ of n observations, assumed to be drawn i.i.d. from a data-generating distribution \mathbb{P}_{xy} . An ML model is a function $\hat{f} : \mathcal{X} \rightarrow \mathbb{R}^g$ that assigns a prediction to a feature vector from \mathcal{X} .¹ \hat{f} is itself constructed by an *inducer* function \mathcal{I} , i.e., the model-fitting algorithm. The inducer $\mathcal{I} : (\mathcal{D}, \lambda) \mapsto \hat{f}$ uses training data \mathcal{D} and a vector of *hyperparameters* $\lambda \in \Lambda$ that govern its behavior. The overall goal of supervised ML is to derive a model \hat{f} from a data set \mathcal{D} so that \hat{f} predicts data sampled from \mathbb{P}_{xy} best. The quality of a prediction is measured as the discrepancy between predictions and ground truth. This is operationalized by the loss function $L : \mathcal{Y} \times \mathbb{R}^g \rightarrow \mathbb{R}_0^+$, which is to be minimized during model fitting. In contrast to the optimization problems that we will define in Sections III-B and III-C, we term this the “first-level” optimization problem.

The expectation of the loss value of predictions made for data samples drawn from \mathbb{P}_{xy} is the *generalization error*

$$\text{GE} := \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_{xy}} [L(y, \hat{f}(\mathbf{x}))] \quad (1)$$

which cannot be computed directly if \mathbb{P}_{xy} is not known beyond the available data \mathcal{D} . Therefore, one often uses so-called *resampling* techniques that fit models on N_{iter} subsamples $\mathcal{D}[J_j]$ and evaluate them on complements $\mathcal{D}[-J_j]$ of these

¹where g allows handling of multioutput regression, as well as multiclass classification with g classes by returning decision scores.

subsets to obtain an estimate of the generalization error

$$\widehat{\text{GE}}(\mathcal{I}, \lambda, \mathbf{J}) = \frac{1}{N_{\text{iter}}} \sum_{j=1}^{N_{\text{iter}}} L(y[-J_j], \mathcal{I}(\mathcal{D}[J_j], \lambda)(x[-J_j])). \quad (2)$$

Depending on the resampling method, the inducer \mathcal{I} , and the quantity of data in \mathcal{D} , estimating the generalization error $\widehat{\text{GE}}(\mathcal{I}, \lambda, \mathbf{J})$ can require large amounts of computational resources.

B. Hyperparameter Optimization

The goal of HPO is to identify a hyperparameter configuration that performs well in terms of the estimated generalization error in (2). Often, optimization only concerns a subspace of available hyperparameters because some hyperparameters might be set based on prior knowledge or due to other constraints. One would therefore split up the space of hyperparameters Λ into a subspace of hyperparameters Λ_S over which optimization takes place, and the remaining hyperparameters $\Lambda_C = \Lambda / \Lambda_S$ for which values λ_C are given exogenously. We define the HPO problem as

$$\lambda_S^* \in \underset{\lambda_S \in \Lambda_S}{\operatorname{argmin}} c(\lambda_S) = \underset{\lambda_S \in \Lambda_S}{\operatorname{argmin}} \widehat{\text{GE}}(\mathcal{I}, (\lambda_S, \lambda_C), \mathbf{J}). \quad (3)$$

Here, λ_S^* denotes a theoretical optimum, and $c(\lambda_S)$ is a shorthand for the estimated generalization error in (2). We refer to Problem 3 as the “second-level” optimization problem.

Hyperparameters can be either continuous, discrete, or categorical, and search spaces are often a mix of the different types. The search space may be hierarchical, i.e., some subordinate hyperparameters can only be set in a meaningful way if another parent hyperparameter takes a certain value. In particular, many AutoML frameworks perform optimization over a hierarchical hyperparameter space that represents the components of a complex ML pipeline [1].

Many HPO algorithms can be characterized by how they handle two different tradeoffs: 1) the exploration versus exploitation tradeoff refers to how much budget an optimizer spends on either trying to directly exploit the currently available knowledge base by evaluating very close to the currently best candidates (e.g., local search) or whether it explores the search space to gather new knowledge (e.g., RS) and 2) the inference versus search tradeoff refers to how much time and overhead is spent to induce a model from the currently available archive data in order to exploit past evaluations as much as possible. Other relevant aspects that HPO algorithms differ in are: *Parallelizability*, i.e., how many configurations a tuner can (reasonably) propose at the same time; *global versus local* behavior of the optimizer, i.e., if updates are always quite close to already evaluated configurations; *noise handling*, i.e., if the optimizer takes into account that the estimated generalization error is noisy; *search space complexity*, i.e., if and how hierarchical search spaces can be handled; *multifidelity*, i.e., if the optimizer uses cheaper evaluations to infer performance on the full data.

Multifidelity methods make use of the fact that the resampling procedure in (2) can be modified in multiple ways to make evaluation cheaper: one can 1) reduce the training sizes

$|J_j|$ via subsampling, as model evaluation complexity is often at least linear in training set size or 2) change some components in λ in a way that makes model fits cheaper. Examples of 2) are reducing the overall number of training cycles performed by a neural network fitting process or reducing the number of base learner fits in a bagging or boosting method. These modifications can both increase the variance of $\widehat{\text{GE}}$ and introduce an (often pessimistic) bias, as models trained on smaller datasets or with values of λ that make fitting cheaper often have worse generalization errors.

We introduce a *fidelity* parameter $r \in (0, 1]$ that influences the resource requirements of the evaluation of $\widehat{\text{GE}}$ and define

$$c(\lambda_S; r) := \widehat{\text{GE}}(\mathcal{I}, (\lambda_S, \lambda_C(r)), \mathbf{J}(r)). \quad (4)$$

With this definition, we make the choice that r should influence the evaluation cost of $\widehat{\text{GE}}$ only by modifying the resampling, $\mathbf{J}(r)$ or by modifying a hyperparameter $\lambda_C(r)$. Typically, r only affects one of these aspects at a time, and if it affects λ_C , it only affects a single hyperparameter dimension.

Note that we normally assume that a higher fidelity r returns a better model in terms of the estimate of the generalization error, and the best estimate is returned for $r = 1$. Therefore, r enters the expression in a way where it can influence performance but is not searched over. We define $c(\lambda_S) := c(\lambda_S; 1)$ as in [44], and the optimization problem remains as in (3).

This assumption may be violated in some scenarios, and model performance could worsen for a higher value of r (e.g., a neural network, which may overfit on a small dataset if trained for too many epochs). In this case, we define the optimization problem as $(\lambda_S^*, r^*) \in \underset{\lambda_S \in \Lambda_S, r \in (0, 1]}{\operatorname{argmin}} c(\lambda_S; r)$.

The resource requirements of evaluating $c(\lambda; r)$ can have a complicated relationship with λ and r ; in practice, r is chosen in such a way that it has an overwhelming and linear influence on resource demand. The overall cost of optimization up to a given point in the optimization process is therefore assumed to be the cumulative sum of the values of r of all evaluations of $c(\lambda; r)$ up to that point. We can also interpret r as the fraction of the budget of a single full fidelity model evaluation that must be spent for evaluating $c(\lambda; r)$.

Given the definition of the HPO problem, we present an (MF-)HPO algorithm for a single, synchronous worker in its most generic form in Algorithm 1. Until a predetermined budget is exhausted, such an algorithm decides in every iteration 1) which configuration(s) λ_S to evaluate next and 2) which fidelity r to use for evaluation; nonmultifidelity algorithms set this to $r = 1$ as default. The algorithm makes use of an *archive* \mathcal{A} , a database recording previously proposed hyperparameter configurations and, if available, their evaluation results. This database can be shared among multiple worker processes that optimize concurrently.

The optimization process can be accelerated by making efficient use of parallel resources. We distinguish between *synchronous* and *asynchronous* scheduling. The former starts multiple evaluations synchronously at the same time and waits until all of these have finished. To be more precise, a number of $k > 1$ configurations are proposed in line 2 and evaluated

Algorithm 1 Generic HPO Algorithm

```

1: while budget is not exhausted do
2:   Propose  $(\lambda_S^{(i)}, r^{(i)})$ ,  $i = 1, \dots, k$ , based on archive  $\mathcal{A}$ 
3:   Write proposals into a shared archive  $\mathcal{A}$ 
4:   Estimate generalization error(s)  $c(\lambda_S^{(i)}; r^{(i)})$ 
5:   Write results into shared archive  $\mathcal{A}$ 
6: end while
7: Wait for workers to synchronize
8: Return best configuration in archive  $\mathcal{A}$ 

```

in parallel in line 4, all within the inner loop of Algorithm 1. Given K available parallel resources, it should be ensured that the number k of configurations scheduled in parallel is not significantly smaller than K and that the evaluation runtimes amongst these k configurations do not differ significantly in order to avoid unnecessarily idling single parallel resources. In contrast, for asynchronous scheduling, Algorithm 1 is run individually in K separate worker processes. Given a shared archive that is synchronized between the workers, every worker can independently schedule new configurations to evaluate.

C. Algorithm Design and Configuration

Our goal will be to design and configure a new HPO algorithm based on a superset of design choices included in previously published HPO methods. We are interested in finding a configuration (or making design choices) based on a set of training instances that works across a broad set of future problem instances. This problem is called *algorithm configuration* [5], [45]. It is quite similar to HPO; a major difference is that algorithm configuration optimizes the configuration of an arbitrary algorithm over a diverse set of often heterogeneous instances for optimal average performance, while HPO performs a per-instance configuration of an ML inducer for a single data set. We introduce the following notation for consistency with the relevant literature: γ denotes configuration parameters controlling our optimizer A , while λ denotes hyperparameters optimized by our optimizer, controlling our inducer \mathcal{I} . The algorithm configuration problem can be formally stated as follows. Given an algorithm $A : \Omega \times \Gamma \rightarrow \Lambda$ parameterized by $\gamma \in \Gamma$ and a distribution \mathbb{P}_Ω over problem instances Ω together with a cost metric ζ , we must find a parameter setting γ^* that minimizes the expected $\zeta(A)$ over \mathbb{P}_Ω

$$\gamma^* \in \underset{\gamma \in \Gamma}{\operatorname{argmin}} \mathbb{E}_{\omega \sim \mathbb{P}_\Omega} [\zeta(A(\omega, \gamma))]. \quad (5)$$

In our example, Γ corresponds to the space of possible components of our HPO method and Ω corresponds to a class of HPO problems (i.e., ML methods and datasets on which they are evaluated) for which their configuration should be optimal. Based on a training set of representative instances $\{\omega_i\}$ drawn from \mathbb{P}_Ω , a configuration γ^* that minimizes c across these instances should be chosen through optimization. When necessary, we refer to this process as the “third-level” optimization problem to distinguish it from the optimization performed by the HPO algorithm A , i.e., the second-level optimization.

Algorithm 2 SMASHY Algorithm

Configuration Parameters: batch size schedule $\mu(b)$, number of fidelity stages s , survival rate η_{surv} , fidelity rate η_{fid} , SAMPLE method (either SAMPLETOURNAMENT or SAMPLEPROGRESSIVE), *batch_method* (one of equal, SH, or HB), total budget B ; further configuration parameters of SAMPLE: \mathcal{I}_{sur} , $\mathbb{P}_\lambda(\mathcal{A})$, $\rho(t)$, $(N_s^0(t), N_s^1(t))$, n_{trn} .

State Variables: Expended budget fraction $t \leftarrow 0$, bracket counter $b \leftarrow 1$ (remains 1 for *batch_method* \in {equal, SH}), current fidelity $r \leftarrow 1$, batch of proposed configurations $C \leftarrow \emptyset$

```

1: while  $t < 1$  do
2:   if  $r = 1$  then ▷ Generate new batch of configurations
3:      $r \leftarrow (\eta_{\text{fid}})^{b-s}$ 
4:      $C \leftarrow \text{SAMPLE}(\mathcal{A}, \mu(b), r; \mathcal{I}_{\text{sur}}, \mathbb{P}_\lambda(\mathcal{A}),$ 
5:        $\rho(t), (N_s^0(t), N_s^1(t)), n_{\text{trn}})$ 
6:     if batch_method = HB then
7:        $b \leftarrow (b \bmod s) + 1$ 
8:     end if
9:   else ▷ Progress fidelity
10:     $r \leftarrow r \cdot \eta_{\text{fid}}$ 
11:     $C \leftarrow \text{SELECT\_TOP}(C, |C|/\eta_{\text{surv}})$ 
12:    if batch_method = equal then
13:       $\tilde{\mu} \leftarrow \mu(b) - |C|$ 
14:       $C \leftarrow C \cup \text{SAMPLE}(\mathcal{A}, \tilde{\mu}, r; \mathcal{I}_{\text{sur}}, \mathbb{P}_\lambda(\mathcal{A}),$ 
15:         $\rho(t), (N_s^0(t), N_s^1(t)), n_{\text{trn}})$ 
16:    end if
17:    end if
18:    Evaluate configuration(s)  $c(\lambda_S; r)$  for all  $\lambda_S \in C$ 
19:    Write results into shared archive  $\mathcal{A}$ 
20:     $t \leftarrow t + r \cdot |C|/B$  ▷ Update budget spent
21: end while

```

IV. FORMALIZING BROAD CLASS OF MF-HPO ALGORITHMS

We aim to find an HPO algorithm that performs particularly well in the multifidelity setting. To design an algorithm by optimization, we propose a framework and search space of HPO algorithm candidates that cover a large class of possible algorithms and focus on a subclass of algorithms similar to HB because of their favorable properties. This subclass focuses on multifidelity algorithms that use a predefined schedule of geometrically increasing fidelity evaluations containing algorithms like HB [14] and BOHB [16].

The basis of this framework is presented in Algorithm 2, which can be configured by combining algorithmic building blocks in novel ways. The main difference to Algorithm 1 is that the *Propose* part is specified more explicitly. At its core, Algorithm 2 consists of two parts: 1) sampling new configurations at low fidelities (lines 2–7) and 2) increasing the fidelity for existing configurations (lines 8–14). In contrast to Algorithm 1, Algorithm 2 makes use of state variables t , b , and r to account for optimization progress. However, these variables are only shown in Algorithm 2 for clarity and can, in principle, be inferred from the archive \mathcal{A} . As argued in Section III, every single worker instance of Algorithm 1 can, in principle, be scheduled asynchronously, but we do not consider this in this work.

TABLE I

RS, BO, SH, HB, AND BOHB AS INSTANCES OF ALGORITHM 2. η , ρ , AND N_s ARE CONFIGURATION PARAMETERS OF THE RESPECTIVE ALGORITHMS. “—” DENOTES THAT THE VALUE HAS NO INFLUENCE ON THE ALGORITHM IN THIS CONFIGURATION. *: BO AND BOHB USE INDUCERS THAT PRODUCE NONSTANDARD MODEL FUNCTIONS, WHICH DO NOT AIM TO PREDICT THE ACTUAL PERFORMANCE OF CONFIGURATIONS, AND INSTEAD CALCULATE THE VALUE OF AN ACQUISITION FUNCTION SUCH AS EI [4] (FOR BO) OR THE RATIO OF TWO KERNEL DENSITY ESTIMATOR (KDE) MODELS (FOR BOHB). †: IN A SMALL DEPARTURE FROM BOHB, ALGORITHM 2 USES THE KDE ESTIMATE OF GOOD POINTS FOR ALL SAMPLED POINTS, EVEN WHEN RANDOMLY INTERLEAVED. BOHB RANDOMLY INTERLEAVES FROM A UNIFORM DISTRIBUTION

Algorithm	$\mu(b)$	s	η_{surv}	η_{budget}	$\mathcal{I}_{f_{\text{sur}}}$	ρ	N_s	<i>batch_mode</i>	$\mathbb{P}_{\lambda}(\mathcal{A})$
RS	—	1	—	—	—	1	—	—	uniform
BO	1	1	—	—	e.g. GP+EI*	ρ	N_s	—	uniform
SH	μ	$\lfloor -\log_{\eta}(r_{\min}) \rfloor + 1$	η	η	—	1	—	SH	uniform
HB	$\lceil s \cdot \frac{\eta^{s-b}}{s-b+1} \rceil$	$\lfloor -\log_{\eta}(r_{\min}) \rfloor + 1$	η	η	—	1	—	HB	uniform
BOHB	$\lceil s \cdot \frac{\eta^{s-b}}{s-b+1} \rceil$	$\lfloor -\log_{\eta}(r_{\min}) \rfloor + 1$	η	η	TPE*	ρ	N_s	HB	KDE†

In its first iteration, Algorithm 2 uses a SAMPLE-subroutine to initialize the initial batch C of μ solution candidates. The fidelity of the evaluation of the proposed configurations is refined iteratively; when all configurations in the batch have been evaluated with given fidelity r , the top $1/\eta_{\text{surv}}$ fraction of configurations is evaluated with a fidelity that is increased by a factor of η_{fid} . When the fidelity cannot be further increased for a batch because all of its configurations were evaluated at full fidelity $r = 1$, they are set aside, and a new batch of configurations is sampled.

The SAMPLE subroutine creates new configurations to be evaluated, possibly using information from the archive to propose points that are likely to perform well. We allow that any inducer $\mathcal{I}_{f_{\text{sur}}}$ that produces a surrogate model f_{sur} can be used for model-assisted sampling. The subroutine works by at first sampling a number of points from a given generating distribution $\mathbb{P}_{\lambda}(\mathcal{A})$. The performance of these points is then predicted using the surrogate model, and points with unfavorable predictions are discarded in a process we refer to as *filtering*. This process is repeated until the requested number μ of nondiscarded points is obtained. N_s and ρ have the same function as in [16] (see Section IV-A5), with the filter factor N_s controlling the number of sampled points needed for each of the μ points returned, and ρ controlling the fraction of points that are not filtered. Thus, the configuration space of sampling methods also includes purely random sampling, as in HB, by setting $\rho = 1$. The influence of the surrogate model on sampled candidates is larger when 1) the number of sampled configurations N_s is large or 2) the fraction ρ of candidates sampled at random is small. We present two slightly different SAMPLE algorithms: SAMPLETOURNAMENT (Algorithm 3) and SAMPLEPROGRESSIVE (Algorithm 4) based on this principle (see Appendix A in the supplementary material). Both allow to use different N_s values for different points they sample, parameterized by N_s^0 and N_s^1 .

While hyperparameters λ_S are proposed by one of the two SAMPLE methods, the fidelity hyperparameter r follows a fixed schedule similar to SH [19] and HB [14], with a few extensions. For one, the survivor factor η_{surv} can be a different value from the fidelity scaling factor η_{fid} . Furthermore, the algorithm allows three scheduling modes, controlled by *batch_method*: SH does SH. The HB mode evaluates brackets, as performed by HB. While $\mu(b)$ is, in principle, a free configuration parameter

for every value of b , we choose to set $\mu(b)$ so that total budget expenditure is approximately equal between all brackets. This follows the principle used in HB, but the dependency on η_{surv} and η_{fid} is more complex and determined dynamically. Finally, *equal batch_method* uses equal batch sizes for every evaluation. Individuals that perform badly at low fidelity are removed, as in SH, but new individuals are sampled to fill up batches to the original size. Because new individuals are added to the batches at all fidelity steps, it is not necessary to use brackets with different initial fidelities, and therefore, only a single repeating bracket $b = 1$ is used. The *equal* method is an original contribution of this work and was designed to be similar to HB while using parallel resources more efficiently; the two batch scheduling methods are illustrated in Fig. 1.

If the exploration–exploitation tradeoff is not balanced properly, the optimization progress can either stagnate or function evaluations are wasted due to too much exploration of uninteresting regions of the search space. However, the relative importance of exploration and exploitation can change throughout the course of optimization, where exploration performed later during the optimization is not as useful as during the beginning. The given configuration space makes it possible to make the exploration–exploitation tradeoff dependent on optimization progress by providing the option to make $\rho(t)$ and $(N_s^0(t), N_s^1(t))$ dependent on the proportion of exhausted total budget at every configuration proposal step. It is likely that large values of $\rho(t)$ /small values of $N_s^*(t)$ perform better when t is small. Conversely, it is likely that small $\rho(t)$ /large $N_s^*(t)$ work well for large t .

A. Common MF-HPO Algorithms Covered by Algorithm 2

The following describes a few common HPO algorithms that can be instantiated within this framework; see Table I for specific configuration parameter settings within Algorithm 2 that correspond to these algorithms.

1) *Random Search*: Configurations λ_S are drawn (uniformly) at random, and every configuration is evaluated with full fidelity $r = 1$. Parallelization is straightforward, as configurations are drawn independently.

2) *Bayesian Optimization* [8]: The configuration that maximizes an acquisition function $a(\lambda)$ (e.g., expected improvement, EI [4]) is proposed and evaluated with the full fidelity $r = 1$. $a(\lambda)$ is based on a surrogate model trained on the

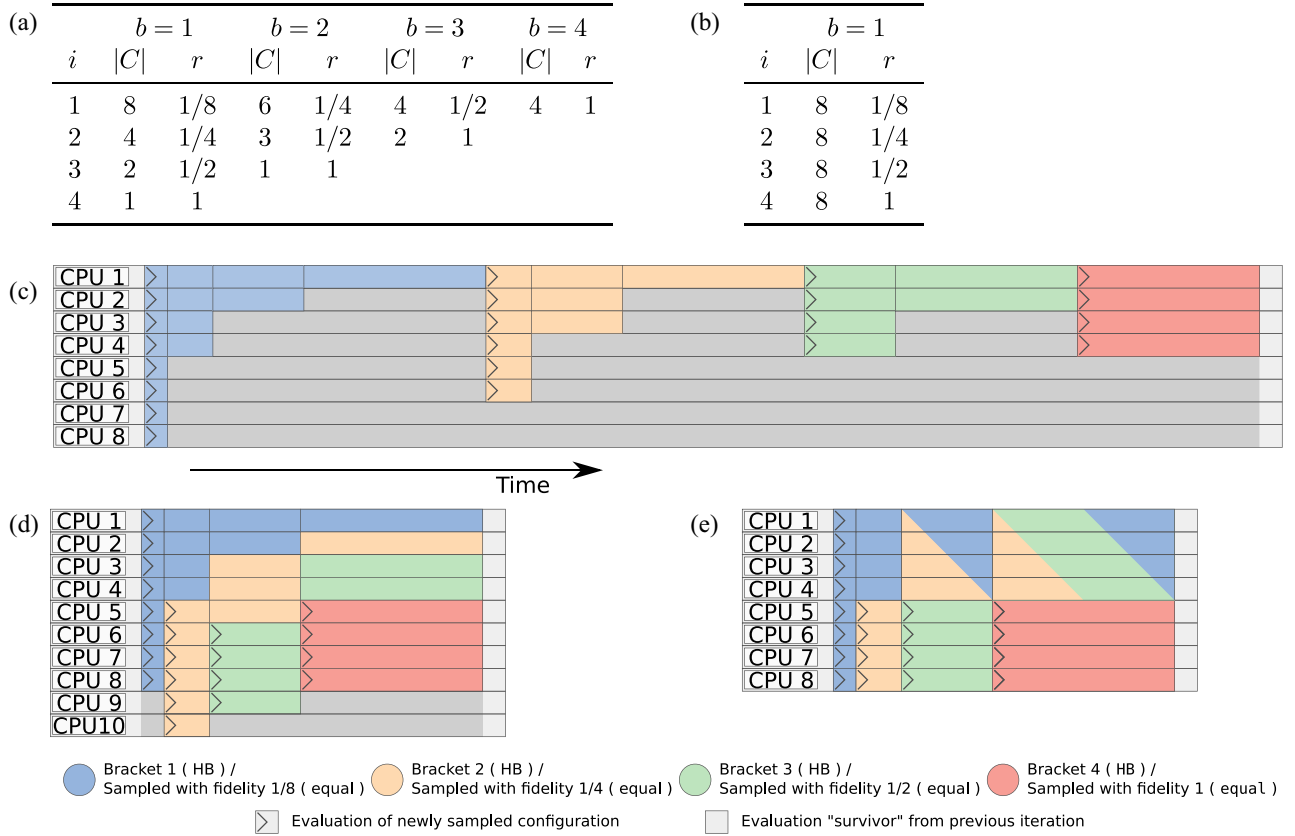


Fig. 1. Illustration of the different *batch_methods* used, corresponding to the values of $\eta_{\text{fid}} = \eta_{\text{surv}} = 2$, $s = 4$, and $\mu = 8$. The tables show the (a) HB method and (b) *equal* method. Shown are the number $|C|$ and fidelity value r of configurations being evaluated in the iterations i of the various brackets counted by b . Except for i , the variables are the same as in Algorithm 2. Subfigures (c)–(e) illustrate resource utilization by the batch methods, given availability of parallel resources. (c) Naively scheduling the configuration evaluations one batch after another can make use of available parallel resources but leaves many of them idle. (d) Hypothetical way of scheduling configuration evaluations of different brackets at the same time so that all configurations with the same r -value are scheduled together utilizes resources more efficiently, but the number of evaluations in each batch still varies. (e) Simpler *equal* batch scheduling method always evaluates the same number of configurations within each batch and, therefore, makes optimal use of available parallel resources.

archive \mathcal{A} . BO can be parallelized by either using methods that can propose multiple points at the same time using a single surrogate model or, alternatively, by fitting a surrogate model on the anticipated outcome of configurations that were proposed but not yet evaluated [11]. BO can be represented in Algorithm 2 by using an inducer $\mathcal{I}_{\text{furr}}$ that produces a function f_{furr} equal to the composition of model prediction and acquisition function. In its basic form, BO is not an MF algorithm and therefore always sets $r = 1$.

3) *Successive Halving* [19]: SH, also called sequential halving [46], is a simple multifidelity optimization algorithm that combines the random sampling of configurations with a fixed schedule for r . At the beginning, a batch of μ configurations is sampled randomly and evaluated with an initial fidelity $r_{\text{min}} < 1$. This is followed by repeated “halving” steps, where the top fraction η^{-1} of configurations is kept and evaluated after r is increased by a factor of η , until the maximum fidelity value is reached. The schedule is chosen to keep the total sum of all evaluated r constant in each batch. Both η_{surv} and η_{fid} in Algorithm 2 correspond to SH’s η -parameter.

4) *Hyperband* [14]: Similar to SH, HB uses a fixed schedule for the fidelity parameter r , but it augments SH by using multiple *brackets* b of SH runs starting at different $r_{\text{min}}(b)$ and with different $\mu(b)$. The number of brackets is

set to

$$s = \lfloor \log_{\eta}(1/r_{\text{min}}) \rfloor + 1 \quad (6)$$

which coincides with the number of fidelity steps that can be performed on a geometric scale on the interval $[r_{\text{min}}, 1]$. In bracket $b \in \{1, 2, \dots, s\}$, a number of $\mu(b)$ samples are initially sampled and evaluated with initial fidelity $r = \eta^{s-b}$. $\mu(b)$ is chosen such that each bracket needs an approximately similar amount of budget: $\mu(b) = \lceil s \cdot (\eta^{s-b}/s - b + 1) \rceil$.

5) *Bayesian Optimization Hyperband* [16]: Model-based methods outperform HB when a relatively large amount of budget is available and many objective function evaluations can be performed. BOHB was created to overcome this drawback. This method iterates through SH brackets like HB, but, instead of sampling new configurations randomly, it uses information from the archive to propose points that are likely to perform well. A total number of N_s configurations are proposed for evaluation; ρ are sampled at random, and the rest are chosen based on a surrogate model induced on the evaluated configurations in \mathcal{A} . The models used by BOHB are a pair of KDEs of the top and bottom configurations in \mathcal{A} , similar to the process in [47]. To implement BOHB in Algorithm 2, one, therefore, needs to use an inducer $\mathcal{I}_{\text{furr}}$ that

TABLE II
THREE BENCHMARK COLLECTIONS OF YAHPO GYM USED IN OUR BENCHMARK

Scenario	Target Metric	d	Hyperparameter Types				# Instances	# Training Set
			Cont.	Integer	Categ.	Hierarchical		
<i>lcbench</i> : HPO of a neural network	cross entropy loss	7	6	1	0	✗	35	8
<i>rbv2_super</i> : AutoML pipeline configuration	log loss	38	20	11	7	✓	89	30
<i>nb301</i> : Neural architecture search	validation accuracy	34	0	0	34	✓	1	—

produces a function that calculates the ratio of kernel densities, an unusual kind of regression model.

B. Limitations and Further MF-HPO Algorithms

The following lists notable HPO algorithms not currently covered by the optimization space of Algorithm 1. They were excluded because they differ in too substantial ways from the other algorithms considered here.

1) *FABOLAS* [48]: Fabolas is a continuous multifidelity BO method, where the conditional validation error is modeled as a Gaussian process using a complex kernel-capturing covariance with the training set fraction $r \in (0, 1]$ to allow for adaptive evaluation at different resource levels.

2) *Asynchronous Successive Halving* [15] and *Asynchronous Hyperband*: HB, as well as SH, have the drawback that batch sizes decrease throughout the stages of an SH run, preventing efficient utilization of parallel resources. ASHA is an effective method to parallelize SH by an asynchronous parallelization scheme. A shared archive across a number of different workers is maintained. Instead of waiting until all n configurations of a batch have been evaluated for fidelity r , every free worker queries the shared archive \mathcal{A} for “promotable” configurations (i.e., configurations that belong to the fraction of top η^{-1} configurations evaluated with the same fidelity). Asynchronous HB works similarly.

3) *Asynchronous BOHB* [17]: A-BOHB, an asynchronous extension of BOHB where configurations are sampled from a joint Gaussian Process, explicitly capturing correlations across fidelities. In contrast to ASHA and asynchronous versions of BOHB in the original BOHB publication [16], A-BOHB does not perform synchronization after each stage but instead uses a stopping rule [49] to asynchronously determine whether a configuration should continue to run or be terminated.

V. EXPERIMENTAL ANALYSIS

Given the formalization of the framework in Section IV, our goal is to find the best representative (out of this class of algorithms) by solving the third-level optimization problem in (5), and explain the role of specific algorithmic components in a benchmark-driven approach. We aim to answer the following research questions.

- RQ1: How does the optimal configuration of our MF-HPO framework differ between problem scenarios, i.e., do different problem scenarios benefit from different HPO algorithms?
- RQ2: How does our optimized MF-HPO algorithm compare to other established HPO implementations?
- RQ3: Does the successive-halving fidelity schedule have an advantage over the simpler equal-batch-size schedule?

RQ4: What is the effect of using multifidelity methods in general?

RQ5a: Does changing SAMPLE configuration parameters throughout the optimization process offer an advantage?

RQ5b: Does (more complicated) surrogate-assisted sampling in SAMPLE provide an advantage over using simple random sampling with surrogate filtering?

RQ6: What effect do different surrogate models (or using no model at all) have on performance?

RQ7: Does the equal-batch-size schedule give an advantage over established methods when parallel resources are available?

We rely on benchmark scenarios of the YAHPO Gym benchmark suite [50], each of which provides a number of related *instances* of optimization problems. The benchmark scenarios we have chosen cover three important application areas of AutoML: HPO of a neural network (*lcbench*), AutoML pipeline configuration (*rbv2_super*), and neural architecture search (*nb301*). These classes of problems do not only represent common and relevant tasks for researchers and practitioners in the field; as presented in Table II, they are also quite different with regards to: 1) the dimensionality of the search space; 2) hyperparameter types (categorical, integer, and continuous); and 3) whether there are hierarchical dependencies between hyperparameters. More details on the characteristics of the problem classes are given in Appendix B in the supplementary material. To avoid an optimistic bias in the analysis caused by over-adaption to the random peculiarities of the particular instances used during configuration, we are using meta-holdout splits on the level of HPO problem instances (see Appendix D in the supplementary material). This means that for analyzing the performance of a configured candidate of Algorithm 2, we are evaluating this candidate by running it on instances that were not seen during configuration. Algorithm 2 is always run with a budget limit corresponding to $30 \cdot d$ full fidelity evaluations (where d is the dimension of the problem instance).

A. Algorithm Design via Configuration

First, we describe the experiments we conducted to configure Algorithm 2 via optimization.

We follow the PBO principle and configure Algorithm 2 by optimizing separately for different HPO scenarios, namely, for *lcbench* and *rbv2_super*, resulting in two optimized configurations $\boldsymbol{y}^{*lcbench}$ and $\boldsymbol{y}^{*rbv2_super}$, respectively. The *nb301* scenario is not used for configuration, but exclusively for subsequent analysis.

For the algorithm configuration of our framework (third level), the performance objective $\mathbb{E}_{\omega \sim \mathbb{P}_{\Omega}}[\zeta(A(\omega, \boldsymbol{y}))]$ for a

TABLE III

SUMMARY OF EXPERIMENT. SHOWN ARE THE VARIOUS OPTIMIZER CONFIGURATIONS γ THAT WERE OBTAINED FROM OPTIMIZATIONS WITH DIFFERENT CONSTRAINTS. “NAME”: THE NAME BY WHICH WE REFER TO THE CONFIGURATION IN THE TEXT. “RQ”: THE RESEARCH QUESTION THAT MAINLY RELATES TO THE CONFIGURATION. “OPTIMIZE”: WHETHER THE GIVEN CONFIGURATION WAS OBTAINED BY CONDUCTING A (POSSIBLY CONSTRAINED) OPTIMIZATION (\checkmark), OR BY SUBSTITUTING VALUES INTO THE GLOBAL OPTIMUM γ^*

Name	RQ	Optimize	Design Modification
γ^*	1, 2, 3	\checkmark	none (global optimization)
γ_1	4	\times	$\eta_{\text{fid}} \rightarrow \infty$
γ_2	5a	\checkmark	$n_{\text{trn}}(0) = n_{\text{trn}}(1)$, $N_s^0(0) = N_s^0(1)$, $N_s^1(0) = N_s^1(1)$, $\rho(0) = \rho(1)$
γ_3	5b	\checkmark	$\text{filter_method} \rightarrow \text{tournament}$, $n_{\text{trn}} \rightarrow 1$, $N_s^0(0) = N_s^0(1) = N_s^1(0) = N_s^1(1)$, $\rho(0) = \rho(1)$
γ_4	6	\times	$\text{batch_method} \rightarrow \text{equal}$, $\mathcal{I}_{f_{\text{sur}}} \rightarrow *$
γ_5	6	\times	$\text{batch_method} \rightarrow \text{equal}$, $\rho \rightarrow 0$
γ_6	6	\times	$\text{batch_method} \rightarrow \text{equal}$, $\rho \rightarrow 0$, $\mathbb{P}_{\lambda}(\mathcal{A}) \rightarrow \text{uniform}$
γ_7	7	\times	$\text{batch_method} \rightarrow \text{equal}$, $\mu \rightarrow 32$, quadruple budget

configuration γ in (5) is estimated by running Algorithm 2 (i.e., second-level optimization) configured by γ on a set of problem instances and taking the average of observed performances. For this, all problem instances included in the respective benchmark scenario that has not been held out for subsequent analysis are used. As a configuration for our framework, we use BO with the lower confidence bound acquisition function [51] with interleaved random configurations every three evaluations.² Configuration is repeated three times for each scenario, each running for 60 h, with different random seeds. To get the overall best configuration, the set of all evaluated configurations γ (i.e., the third-level optimization archive) is combined into a single data set for each scenario. To estimate the actual best configuration, a common *identification criterion* [52] is used: a surrogate model is fitted on the combined datasets and the optimum among the in-sample predictions of this model is used ($\gamma^{*lcbench}$ and γ^{*rbv2_super} , respectively). We also store the (surrogate-smoothed) optima of all three individual optimization runs and record the range of configuration parameter values to obtain an estimate of the uncertainty of the overall optimal configurations.

The search space used for the optimization of Algorithm 2 is shown in Table V in Appendix C in the supplementary material. While the batch size μ is constant in the *equal batch_method*, it changes for every bracket when *batch_method* is HB. The batch sizes $\mu(2)$, $\mu(3)$, \dots are constructed from $\mu(1)$ dynamically as described in Section IV. The search space contains several surrogate learners: Random forests [53] (RF), K -nearest-neighbors with k set to 1 (KNN1), kernelized K -nearest-neighbors with “optimal” weighting [54] (KKN7), and the ratio of density predictions of good and bad points, similar to tree parzen estimators [47] without a hierarchical structure as in BOHB [16] (TPE). For the prefiltering sample distribution $\mathbb{P}_{\lambda}(\mathcal{A})$, we evaluate both uniform sampling (*uniform*), and sampling from the estimated density of good points as done in BOHB [16] (KDE). *filter_mb* determines whether the surrogate model makes predictions assuming the highest fidelity value r observed (TRUE), as opposed to assuming the fidelity of the points being sampled; in the framework of the SAMPLE Algorithms 3 and 4 in Appendix A in the supplementary material, this influences the behavior of $\mathcal{I}_{f_{\text{sur}}}$.

Note that the maximum number of fidelity steps per batch s is not part of the search space and instead inferred automatically from η_{fid} and the lower bound for r that is given as part of the optimization problem instance. As in HB, it is set to the largest number of stages that is possible given η_{fid} and the lower bound on r according to (6).

B. Algorithm Analysis

Our goal in this work is not only to determine configurations of Algorithm 2 that perform well on the respective benchmarking scenarios but also to determine what effect individual components have on performance. However, performing a complete SA would be prohibitively computationally expensive, as it would require evaluation of the objective (i.e., running Algorithm 2) in an experimental design of different configurations. Instead, we evaluate the performance of the candidate configurations found in Section V-A and alternative configurations—which are chosen in a way to allow for answering our research questions—on the benchmark test instances which were held out during configuration. A simple method to answer many of these questions is to take the optimized configuration of Algorithm 2 and swap components of it for simpler components (or removing them completely), thereby performing a one-factor-at-a-time analysis or an ablation study. However, the optimal values of some components may interact strongly with other components. We, therefore, auto-configure the framework several times under certain constraints dictated by our particular research question at hand. For example, to investigate the effect of varying n_{trn} and N_s over t , we run the optimization of Algorithm 2 with the constraint $n(0)$ to be equal to $n(1)$ and compare the resulting configuration to the overall optimum γ . Table III lists the different values of γ we generate under different constraints. For each value of γ , we run the, respectively, configured HPO algorithm on both the *lcbench* and the *rbv2_super* scenario, and (unless stated otherwise) once each for *batch_method* set to *equal* and HB. We refer to an optimized configuration that was obtained on the *lcbench* scenario with *batch_method* set to *equal* as $\gamma^{*lcbench}[\text{equal}]$, and to the overall optimum (i.e., the better of $\gamma^{*lcbench}[\text{equal}]$ and $\gamma^{*lcbench}[\text{HB}]$) as $\gamma^{*lcbench}$, similar for *rbv2_super*.

²Note that this optimizer used for third-level optimization is not an instance of Algorithm 2.

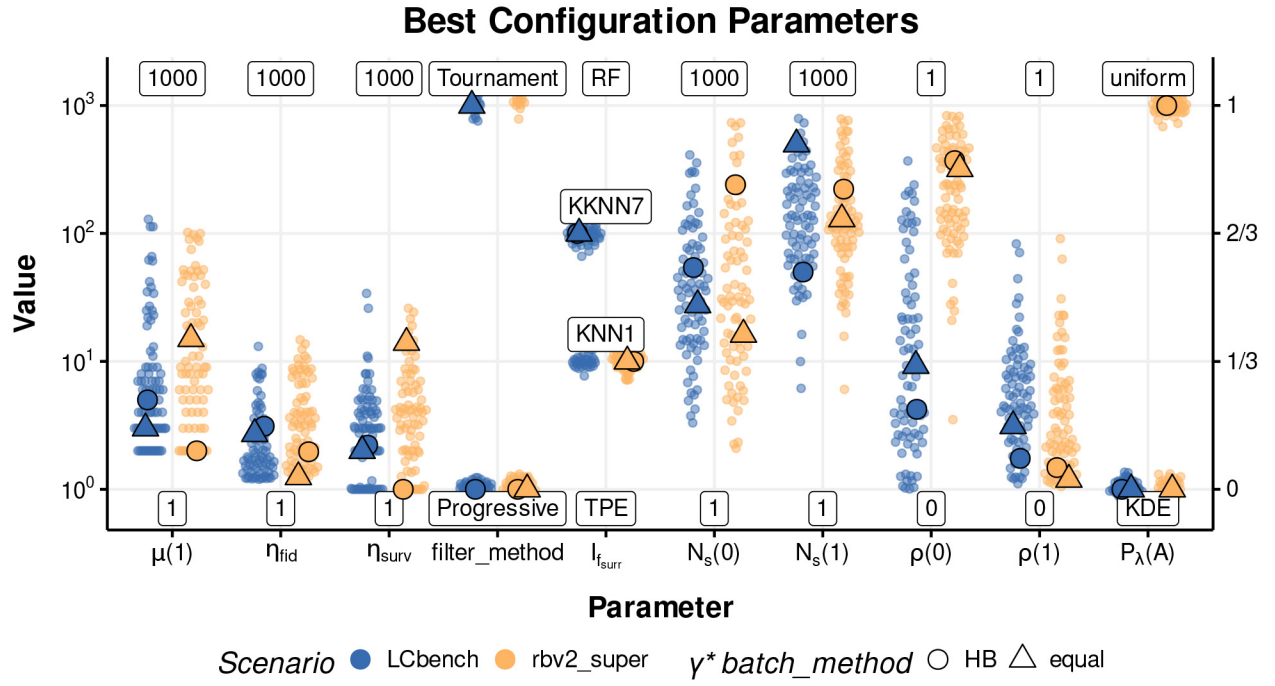


Fig. 2. Beeswarm plot of the best configurations according to the surrogate model over the meta-optimization archive of γ^* . Shown are the top 80 configuration points (according to the surrogate-model-predicted performance) that were evaluated during optimization. Levels of discrete parameters are shown. Most numeric parameters are on a log-scale (left axis), except for $\rho(0)$ and $\rho(1)$, which are on a linear scale (right axis). Instead of showing both $N_s^0(t)$ and $N_s^1(t)$, their geometric mean $N_s(t)$ is shown. The highlighted large points are $\gamma^*[\text{HB}]$ and $\gamma^*[\text{EQUAL}]$, which were found on both benchmark scenarios.

Every evaluation of a framework configuration, i.e., a complete HPO run on a problem instance, is repeated 30 times (with different random seeds) to allow for statistical analysis.

The analysis of our research questions is based on the following tables and visualizations. Table VI in Appendix D in the supplementary material shows the configuration parameters that were selected for each benchmark scenario with various search space restrictions. We perform all optimization runs constrained to the fidelity scheduling `equal` and `HB`, respectively, and denote the resulting optimal configurations $\gamma^*[\text{equal}]$ and $\gamma^*[\text{HB}]$. Fig. 2 shows the configuration values of the top 80 evaluated points according to their surrogate-predicted performance. The ranges covered by the bee swarms are again an indicator of approximate ranges of configuration values that can be expected to work well. Fig. 4 shows the final performance at $30 \cdot d$ full-budget evaluations for all optimization runs that were performed. The standard error shown is the estimated standard deviation of the mean of benchmark-instance-wise performance, representing uncertainty about the “true” performance mean if an infinite number of benchmark instances of the given class of problems were available.

We now describe in more detail how we operationalize each of the research question RQ1–RQ7 and report results.

RQ1: How does the optimal configuration differ between problem scenarios, i.e., do different problem scenarios benefit from different HPO algorithms?

Setup: We investigate the difference in the values that $\gamma^*[\text{LCbench}]$ and $\gamma^*[\text{rbv2_super}]$ take, and put this difference in perspective by comparing it to the uncertainty of these values.

To evaluate how well $\gamma^*[\text{LCbench}]$ and $\gamma^*[\text{rbv2_super}]$ generalize to other problem scenarios, we evaluate them on the respective instances of scenarios that they were not configured on.

Results: As can be seen in Table VI in the supplementary material and in Fig. 2, many of the selected components of the γ^* are relatively close to each other across the two scenarios on which they were optimized, relative to their uncertainty ranges. $\mathcal{I}_{f_{\text{surr}}}$ is chosen as `KNN1` on `rbv2_super`, but can also use `KKNN7` on `LCbench`, which in fact seems to be slightly preferred. This is interesting as KNN-based models are rarely considered in surrogate-based HPO; the typically preferred random forest model was not selected. $\mathbb{P}_\lambda(\mathcal{A})$ takes any of the two values for `rbv2_super`, but is chosen to be `KDE` in `LCbench`. Finally, $\rho(0)$ is close to 1 in the beginning on `rbv2_super`, and closer to 0 (although still greater than $\rho(1)$) for `LCbench`.

The degree to which the differences in γ^* influence the outcome can be observed in Fig. 4. The optimized results generalize well to test instances from the same scenario as they were configured on. Fig. 3 shows the optimization progress (on unseen test instances) of configurations if configured on the same scenario versus configurations that were configured on a different scenario. We see, for example, a clear advantage of the configurations that we obtained by optimizing directly on `LCbench` when we evaluate them on their respective held out test instances. We suspect that this difference in performance is mainly due to the different choices of surrogate model classes $\mathcal{I}_{f_{\text{surr}}}$ as well as the random interleave fraction ρ (cf. Fig. 2), and that specific settings for these two algorithmic components are needed for `LCbench` to reach optimal performance.

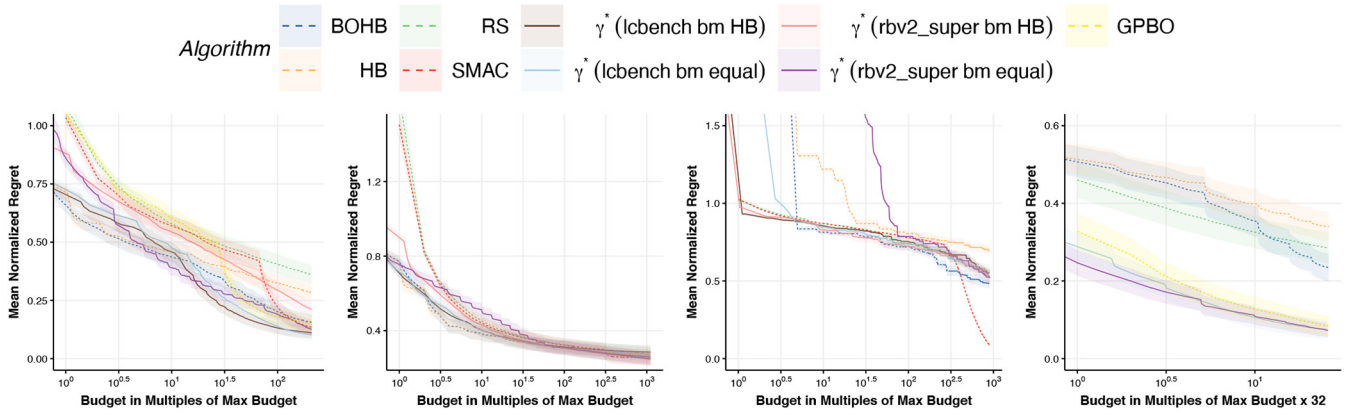


Fig. 3. Optimization progress (mean normalized regret) of serial evaluation on each benchmark scenario as well as $32\times$ parallel evaluation on *lcbench*. Different configurations of Algorithm 2 are executed on benchmark functions that have not been used for the meta-optimization itself, and the progress of these algorithm runs is shown. “ γ^* (lcbench bm equal)” is the configuration obtained from optimizing on *lcbench* with *batch_method*=equal, other labels are constructed similarly. Shown is the mean over 30 evaluations, averaged over all available test benchmark instances for each of the three scenarios. The uncertainty bands show the standard error over the test instances. Note the log-scale on the *x*-axis. Regret is calculated as the difference between the best evaluation performance so far and the overall best value found on each benchmark instance over all experiments; normalized such that 1 corresponds to the median of the performance of all randomly sampled full-fidelity evaluations. We plot performance values observed by the HPO algorithm which depend on evaluation fidelity. This is the reason for the initially “slow” convergence of algorithms that makes their first full-fidelity evaluation late. Note that μ of γ^* [equal] was set to 32 for the parallel evaluations, and HB and BOHB were only naïvely parallelized to simulate a synchronous “single optimizer, multiple workers” environment. See Fig. 6 in Appendix E in the supplementary material for a larger version.

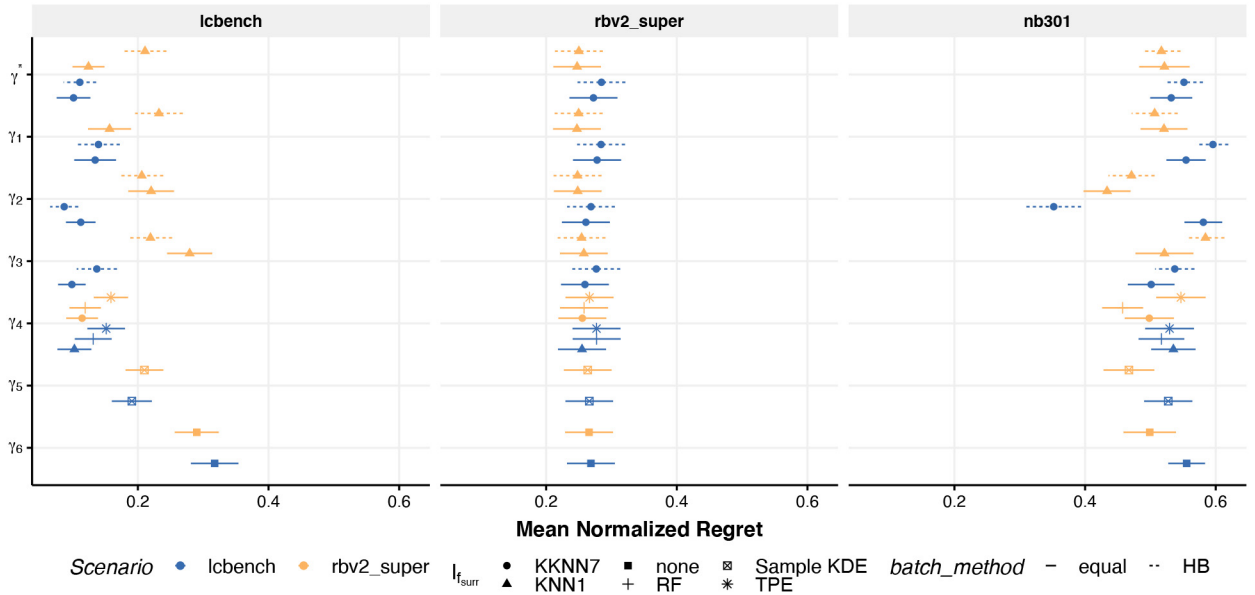


Fig. 4. Mean normalized regret of final performance on “test” benchmark instances for the configuration, shown in Table III. Shown is the mean over 30 evaluations, averaged over all available test benchmark instances for each of the three scenarios. The uncertainty bands show the standard error over instance means. Regret is calculated as the difference between the best evaluation performance so far and the overall best value on each benchmark instance over all experiments; normalized such that 1 corresponds to the median of the performance of all randomly sampled full-fidelity evaluations.

This is not the case for the *rbv2_super* scenario, where none of the different algorithms seem to clearly exploit the problem structure of *rbv2_super* better than others.

RQ2: How does the optimized algorithm compare to other established HPO implementations?

Setup: We evaluate several well-known HPO algorithms in their default configuration on the same benchmark instances: for BOHB [16], we use the implementation found in HpBandSter³ (version 0.7.4); for HB [14], we use

mlr3hyperband⁴ (version 0.1.2); and for SMAC [5], we use the SMACv3 package⁵ (version 1.0.1). We also construct a traditional Gaussian process-based BO (GPBO) [4] with mlrMBO⁶ (version 1.1.5). As GPBO works best with numerical search spaces, we only evaluate it on *lcbench*. Note that GPBO, SMAC, and RS are not multifidelity algorithms and therefore always evaluate points with maximum fidelity 1.

⁴<https://cran.r-project.org/package=mlr3hyperband>

⁵<https://github.com/automl/SMAC3>

⁶<https://cran.r-project.org/package=mlrMBO>

³<https://github.com/automl/HpBandSter>

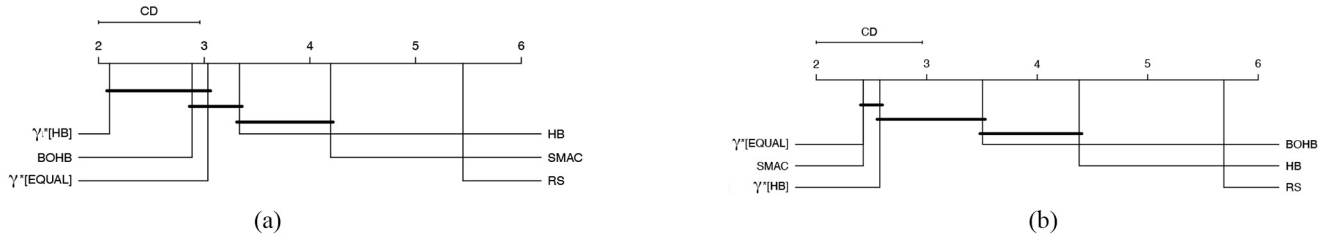


Fig. 5. Critical difference plot [55] comparing the performance of different algorithms across all instances and scenarios. For each of the three scenarios, the mean performance (across replications) for each of the six algorithms is computed ($\gamma^*[\text{HB}]$ is equal to $\gamma^*[\text{icbench}][\text{HB}]$ for instances of the *icbench* scenario, and to $\gamma^*[\text{rbv}][\text{HB}]$ for the *rbv2_super* scenario; same for $\gamma^*[\text{EQUAL}]$). The critical difference test is based on the ranks of the algorithms computed per scenario and instance. Lower ranks are better. Horizontal bold bars indicate that there is no significant difference between algorithms ($\alpha = 1\%$). GPBO, which was not evaluated on all scenarios, is not included. (a) Intermediate optimization budget of 100 full evaluations. (b) Full evaluation budget (final performance).

Results: The performance curves for the mean normalized regret are shown in Fig. 3, and the final performance values at $30 \cdot d$ full-fidelity evaluations are shown in Fig. 4. A critical difference plot and test can be seen in Fig. 5(b). The behavior of RS, HB, BOHB, and SMAC is not surprising; initially, RS and SMAC perform the same, as SMAC evaluates an initial random design. After this, the performance of SMAC improves quickly. HB and BOHB initially both perform better than RS or SMAC because of their multifidelity evaluations, but there is little difference between them. After a while, BOHB starts to outperform HB because of its surrogate-based sampling, which aligns with the observations in [16]. Therefore, BOHB performs well for most budgets, often being the best optimizer for a budget of one as well as for 100 full-fidelity evaluations. Given its multifidelity characteristics, HB is a good choice for low budgets, while SMAC is well suited for larger optimization budgets. Our framework is very competitive on both *icbench* and *rbv2_super*, but is outperformed by SMAC on *nb301*. We assume that this is because Algorithm 2 was not explicitly optimized for the *nb301* scenario.

Although our framework was only optimized for performance at $30 \cdot d$ evaluations, it is also competitive with BOHB after fewer evaluations, as seen in Fig. 5(b).

RQ3: Does the successive-halving fidelity schedule have an advantage over the (simpler) equal-batch-size schedule?

Setup: It is likely that the type of fidelity scheduling used interacts with other configuration parameters. Therefore, we investigate the difference of resulting optimal configurations $\gamma^*[\text{equal}]$ and $\gamma^*[\text{HB}]$.

Results: In both scenarios, the batch method HB is ultimately selected for the optimum γ^* , although Fig. 5(a) and (b) shows that the difference to batch size equal is not statistically significant at $\alpha = 1\%$. We observe that the equal fidelity scheduling mode has several advantages: it is much simpler than HB as it does not need to keep track of SH brackets and does not need to adapt $\mu(b)$ to make the expended budget at each bracket approximately equal. As another benefit, it allows for easy parallel scheduling of evaluations (see also Fig. 1). This is because it always schedules the same number of function evaluations at a time, which can therefore be run synchronously.

RQ4: What is the effect of using multifidelity methods in general?

Setup: We evaluate the performance of a modified γ^* where the number of fidelity stages s is set to 1, thus ensuring that configurations are only evaluated with maximum fidelity 1.⁷

Results: Our results show the superiority of MF-HPO methods compared to HPO methods that do not make use of lower-fidelity approximations. Fig. 5(a) suggests that multifidelity methods are significantly better than their nonmultifidelity counterparts if optimization is stopped at an intermediate overall budget corresponding to 100 full-fidelity evaluations. To be more precise, we see that BOHB as well as both optimized variants $\gamma^*[\text{equal}]$ and $\gamma^*[\text{HB}]$ (optimized for the respective scenario, respectively) significantly outperform SMAC under this strict budget constraint. In line with [14], HB significantly outperforms RS for this budget. On the other hand, Fig. 5(b) provides evidence that multifidelity methods can achieve performance on the same level as state-of-the-art methods that do not make use of low-fidelity approximations (e.g., SMAC) for larger budgets. We conclude that a properly designed multifidelity mechanism provides substantial improvements of anytime performance without affecting performance for larger budgets negatively. In our opinion, the gain in anytime performance justifies the additional algorithmic complexity that is introduced by multifidelity methods.

RQ5a and RQ5b: Does changing SAMPLE configuration parameters throughout the optimization process offer an advantage? Does (more complicated) surrogate-assisted sampling in SAMPLE provide an advantage over using simple random sampling with surrogate filtering?

Setup: To investigate RQ5a (i.e., the effect of the dependence of ρ , n_{tm} and the N_s configuration parameters on t), we performed an optimization where this t -dependence was removed. As these parameters are interpolated between the values at $t = 0$ and $t = 1$, this corresponds to restricting the search space to where these values are equal, as shown for γ_2 in Table III. In addition to this, we ran another optimization where we further restricted N_s^0 and N_s^1 to be equal, n_{tm} to be 1, and only the tournament filter_method be used for RQ5b. The performance of the resulting configurations gives an indication of the performance that is lost for the gain in simplicity.

⁷Because s is not part of the search space Γ and is instead given by 6, this is achieved by setting η_{fid} to ∞ .

Results: The observations made for γ_2 (forbidding change over time) and γ_3 (forbidding change over time and within each batch) are slightly contradictory. In particular, the *nb301* performance of $\gamma_2^{lcbench}$ [HB] is a visible outlier with regards to optimization performance. There is no obvious explanation from inspecting the configuration parameters of $\gamma_2^{lcbench}$ [HB], but it is possible that it is an accidental “good fit” of configuration parameters to the specific landscape of *nb301*.

On *lcbench* and *rbv2_super*, the impact of restricting the search space is smaller and within the uncertainty of the performance of a single configuration. However, we note that both changing configuration parameters over time and within each batch sample introduce significant complexity to the algorithm; thus we prefer the restricted optimization results over γ^* .

RQ6: What effect do different surrogate models (or using no model at all) have on performance?

Setup: We evaluate the overall result γ^* [equal] with \mathcal{I}_{sur} set to each of the inducers in the original search space (see Table V in the supplementary material). Furthermore, γ^* [equal] is evaluated with ρ set to 1 (i.e., all points are sampled randomly from a distribution that may be nonuniform), and finally, with $\rho = 1$ and $P_\lambda(\mathcal{A}) = \text{uniform}$ (i.e., all points are sampled completely uniformly at random).

Results: Surprisingly, the simple k -nearest-neighbors algorithm seems to be chosen consistently by the algorithm configuration for both *lcbench* and *rbv2_super* (see Fig. 2), either with a value of $k = 1$ or $k = 7$. This result is in line with what we already speculated for RQ1. Our ablation experiments suggest that the performance of the optimizer is on average best when using this surrogate learner, even though the differences do not seem to be significant. KNN1 is therefore a reasonable, and simpler, alternative to more complex surrogate learners like the TPE-based method proposed for the original BOHB algorithm.

RQ7: Does the equal-batch-size schedule give an advantage over established methods when parallel resources are available?

Setup: The optimization of ML methods that are expensive to evaluate is often done in parallel; we evaluate the performance of our method and other methods in a (simulated) parallel setting. We evaluate γ^* [equal] with μ set to 32 and with an optimization budget of $30 \cdot 4 \cdot d$, where d is the dimensionality of the optimization problem. We compare it to GPBO with qLCB [10] for 32 parallel evaluations and simulate parallel execution of RS by running $30 \cdot 4 \cdot d$ random evaluations. Both BOHB and SMAC offer parallelized versions, but the YAHPO Gym benchmark package does not yet provide support for asynchronous parallel evaluations [50]. However, since HB and BOHB propose evaluations in batches, we compared HB and BOHB by accounting for submitted batches in increments of 32, essentially simulating a single HB/BOHB optimizer sending evaluations to 32 parallel workers and waiting for their completion synchronously.

Results: Fig. 3 shows that our algorithm is competitive with GPBO—a state-of-the-art synchronously parallel optimization algorithm—when evaluated with 32 parallel resources. This result also shows the main advantage that the equal fidelity

schedule has over scheduling like HB, as synchronously parallelizing HB or BOHB puts them at a great disadvantage over even RS. For HB and BOHB, it is necessary to use asynchronously parallelized methods [15], [17] or use an archive shared between multiple workers [16] to obtain competitive results. However, synchronous objective evaluations are much easier to implement in many environments than asynchronous communication between workers, making the advantage of the simplicity of the equal schedule even more pronounced.

C. Reproducibility and Open Science

The implementation of the framework in Algorithm 2 and reproducible scripts for the algorithm configuration and analysis are available in public repositories.⁸ All data that were generated by our analyses are available as well.

VI. CONCLUSION

We presented a principled approach and framework to benchmark-driven algorithm design and applied it to generic MF-HPO. We formalized the search space of multifidelity hyperparameter optimizers and created a rich and configurable optimization framework. Given the search space, we used BO for meta-optimization of our framework on two different problem scenarios within the field of AutoML and evaluated the result on held out test problems and an entirely held out test scenario. We evaluated the configured optimizers and compared to BOHB, HB, SMAC, and a simple RS as reference. We performed an extensive analysis of the effect of different algorithmic components on performance, while also considering the additional algorithmic complexity they introduce. Our configured framework showed equal and in some cases superior performance to widely used HPO algorithms.

The additional algorithmic complexity introduced by multifidelity evaluations provides substantial benefits. However, based on our experiments, we argue that design choices made by established multifidelity optimizers like BOHB can be replaced by simpler choices: For example, the (more complex) SH schedule is not significantly better than a schedule using equal batch sizes, which allows for more efficient parallelization.

A KDE-based sampling of points to propose, whether filtered by a surrogate model or not, was consistently chosen by our framework. This detail, which is not usually presented as the main feature of BOHB, seems to have an unexpectedly large impact. On the other hand, our optimization results suggest that a surprisingly simple surrogate learner (knn, $k = 1$) can perform even better.

Some components of our search space with large algorithmic complexity have not shown much benefits. Optimization on *rbv2_super* did choose time-varying random interleaving, and overall, more aggressive filtering late during an optimization run ($N_s(1) > N_s(0)$) was slightly favored, but the results did not consistently outperform a configuration obtained from a restricted optimization that excluded time-varying configuration parameters.

⁸<https://github.com/mlr-org/smashy>,

https://github.com/compstat-lmu/paper_2021_benchmarking_special_issue

Our analysis of the set of best observed performances during optimization indicates that there is a large agreement between benchmark scenarios about what the optimal γ^* configuration should be, with parameters that control (model-based) sampling and the surrogate model being the notable exception. This suggests that there may be a set of configuration parameters that are either generally good for many ML problems, or have little impact on performance and can therefore be set to the simplest value. However, some configuration parameters should be adapted to the properties of the particular given optimization problem. The meta-optimization framework presented in this work can be used in future work to investigate the relationship between features of optimization problems and related optimal configurations.

Other fruitful directions for future work include the more in-depth evaluation of asynchronous evaluations; asynchronous methods are important nowadays where parallel resources are plentiful, but current widely used surrogate-based benchmarks do not allow for easy asynchronous evaluations. Suggested methods, such as waiting with a sleep-timer for an appropriate amount [16], are impractical for meta-optimization.

REFERENCES

- [1] B. Bischl et al., "Hyperparameter optimization: Foundations, algorithms, best practices and open challenges," 2021, *arXiv:2107.05847*.
- [2] L. Kotthoff, C. Thornton, H. H. Hoos, F. Hutter, and K. Leyton-Brown, "Auto-WEKA: Automatic model selection and hyperparameter optimization in WEKA," in *Automated Machine Learning: Methods, Systems, Challenges*, F. Hutter, L. Kotthoff, and J. Vanschoren, Eds. Cham, Switzerland: Springer Int., 2019, pp. 81–95.
- [3] D. H. Wolpert and W. G. Macready, "No free lunch theorems for optimization," *IEEE Trans. Evol. Comput.*, vol. 1, no. 1, pp. 67–82, Apr. 1997.
- [4] D. R. Jones, M. Schonlau, and W. J. Welch, "Efficient global optimization of expensive black-box functions," *J. Global Optim.*, vol. 13, no. 4, pp. 455–492, 1998.
- [5] F. Hutter, H. H. Hoos, and K. Leyton-Brown, "Sequential model-based optimization for general algorithm configuration," in *Learning and Intelligent Optimization*, C. A. C. Coello, Ed. Berlin, Germany: Springer, 2011, pp. 507–523.
- [6] K. Swersky, J. Snoek, and R. P. Adams, "Freeze-thaw Bayesian optimization," 2014, *arXiv:1406.3896*.
- [7] J. Bergstra and Y. Bengio, "Random search for hyper-parameter optimization," *J. Mach. Learn. Res.*, vol. 13, pp. 281–305, Feb. 2012.
- [8] J. Snoek, H. Larochelle, and R. P. Adams, "Practical Bayesian optimization of machine learning algorithms," in *Proc. 25th Int. Conf. Neural Inf. Process. Syst. Vol. 2*, 2012, pp. 2951–2959.
- [9] R. Turner et al., "Bayesian optimization is superior to random search for machine learning hyperparameter tuning: Analysis of the black-box optimization challenge 2020," in *Proc. NeurIPS Competition Demonstration Track*, vol. 133. Vancouver, BC, Canada, Dec. 2020, pp. 3–26.
- [10] F. Hutter, H. H. Hoos, and K. Leyton-Brown, "Parallel algorithm configuration," in *Learning and Intelligent Optimization*, Y. Hamadi and M. Schoenauer, Eds. Berlin, Germany: Springer, 2012, pp. 55–70.
- [11] B. Bischl, S. Wessing, N. Bauer, K. Friedrichs, and C. Weihs, "MOI-MBO: Multiobjective infill for parallel model-based optimization," in *Proc. 8th Int. Conf. Learn. Intell. Optim.*, Gainesville, FL, USA, Feb. 2014, pp. 173–186.
- [12] J. González, Z. Dai, P. Hennig, and N. D. Lawrence, "Batch Bayesian optimization via local penalization," in *Proc. 19th Int. Conf. Artif. Intell. Stat. (AISTATS)*, vol. 51. Cadiz, Spain, May 2016, pp. 648–657.
- [13] C. Chevalier and D. Ginsbourger, "Fast computation of the multi-points expected improvement with applications in batch selection," in *Learning and Intelligent Optimization*. Berlin, Germany: Springer, 2013, pp. 59–69.
- [14] L. Li, K. G. Jamieson, G. DeSalvo, A. Rostamizadeh, and A. Talwalkar, "Hyperband: A novel bandit-based approach to hyperparameter optimization," *J. Mach. Learn. Res.*, vol. 18, pp. 1–52, Jan. 2017.
- [15] L. Li et al., "A system for massively parallel hyperparameter tuning," in *Proc. Int. Conf. Mach. Learn. Syst. (MLSys)*, Austin, TX, USA, Mar. 2020, pp. 230–246.
- [16] S. Falkner, A. Klein, and F. Hutter, "BOHB: Robust and efficient hyperparameter optimization at scale," in *Proc. 35th Int. Conf. Mach. Learn. (ICML)*, vol. 80. Stockholm, Sweden, Jul. 2018, pp. 1436–1445.
- [17] L. C. Tiao, A. Klein, C. Archambeau, and M. W. Seeger, "Model-based asynchronous hyperparameter optimization," 2020, *arXiv:2003.10865*.
- [18] D. Sculley et al., "Hidden technical debt in machine learning systems," in *Proc. Annu. Conf. Neural Inf. Process. Syst.*, Montreal, QC, Canada, Dec. 2015, pp. 2503–2511.
- [19] K. G. Jamieson and A. Talwalkar, "Non-stochastic best arm identification and hyperparameter optimization," in *Proc. 19th Int. Conf. Artif. Intell. Stat. (AISTATS)*, vol. 51. Cadiz, Spain, May 2016, pp. 240–248.
- [20] H. H. Hoos, "Programming by optimization," *Commun. Assoc. Comput. Mach.*, vol. 55, no. 2, pp. 70–80, Feb. 2012.
- [21] S. Minton, "Automatically configuring constraint satisfaction programs: A case study," *Constraints*, vol. 1, pp. 7–43, Sep. 1996.
- [22] S. J. Westfold and D. R. Smith, "Synthesis of efficient constraint satisfaction programs," *Knowl. Eng. Rev.*, vol. 16, no. 1, pp. 69–84, 2001.
- [23] D. Balasubramaniam, L. de Silva, C. A. Jefferson, L. Kotthoff, I. Miguel, and P. Nightingale, "Dominion: An architecture-driven approach to generating efficient constraint solvers," in *Proc. 9th Workshop IEEE/IFIP Conf. Softw. Archit.*, Jun. 2011, pp. 228–231.
- [24] A. R. KhudaBukhsh, L. Xu, H. H. Hoos, and K. Leyton-Brown, "SATenstein: Automatically building local search SAT solvers from components," in *Proc. 21st Int. Joint Conf. Artif. Intell.*, San Francisco, CA, USA, 2009, pp. 517–524.
- [25] J.-N. Monette, Y. Deville, and P. van Hentenryck, "Aeon: Synthesizing scheduling algorithms from high-level models," in *Operations Research and Cyber-Infrastructure*. Boston, MA, USA: Springer, 2009, pp. 43–59.
- [26] M. López-Ibáñez and T. Stützle, "The automatic design of multiobjective ant colony optimization algorithms," *IEEE Trans. Evol. Comput.*, vol. 16, no. 6, pp. 861–875, Dec. 2012.
- [27] L. C. T. Bezerra, M. López-Ibáñez, and T. Stützle, "Automatic component-wise design of multiobjective evolutionary algorithms," *IEEE Trans. Evol. Comput.*, vol. 20, no. 3, pp. 403–417, Jun. 2016.
- [28] M. Birattari, Z. Yuan, P. Balaprakash, and T. Stützle, "F-race and iterated F-race: An overview," in *Experimental Methods for the Analysis of Optimization Algorithms*. Berlin, Germany: Springer, 2010, pp. 311–336.
- [29] M. López-Ibáñez, J. Dubois-Lacoste, L. P. Cáceres, M. Birattari, and T. Stützle, "The irace package: Iterated racing for automatic algorithm configuration," *Oper. Res. Perspect.*, vol. 3, pp. 43–58, Jan. 2016.
- [30] N. Dang, L. P. Cáceres, P. D. Causmaecker, and T. Stützle, "Configuring irace using surrogate configuration benchmarks," in *Proc. Genet. Evol. Comput. Conf. (GECCO)*, Berlin, Germany, Jul. 2017, pp. 243–250.
- [31] F. Hutter, H. H. Hoos, K. Leyton-Brown, and T. Stützle, "ParamILS: An automatic algorithm configuration framework," *J. Artif. Intell. Res.*, vol. 36, pp. 267–306, Sep. 2009.
- [32] O. Maron and A. W. Moore, "The racing algorithm: Model selection for lazy learners," *Artif. Intell. Rev.*, vol. 11, nos. 1–5, pp. 193–225, 1997. [Online]. Available: <https://doi.org/10.1005/556606079>
- [33] S. van Rijn, H. Wang, M. van Leeuwen, and T. Bäck, "Evolving the structure of evolution strategies," in *Proc. IEEE Symp. Ser. Comput. Intell. (SSCI)*, 2016, pp. 1–8.
- [34] G. Malkomes and R. Garnett, "Automating Bayesian optimization with Bayesian optimization," in *Proc. Int. Conf. Adv. Neural Inf. Process. Syst.*, vol. 31, 2018, pp. 5984–5994.
- [35] M. Lindauer, M. Feuer, K. Eggenberger, A. Biedenkapp, and F. Hutter, "Towards assessing the impact of Bayesian optimization's own hyperparameters," 2019, *arXiv:1908.06674*.
- [36] M. Lindauer et al., "SMAC3: A versatile Bayesian optimization package for hyperparameter optimization," 2021, *arXiv:2109.09831*.
- [37] A. Saltelli, "Sensitivity analysis for importance assessment," *Risk Anal.*, vol. 22, no. 3, pp. 579–590, 2002.
- [38] W. Hoeffding, "A class of statistics with asymptotically normal distribution," *Ann. Math. Stat.*, vol. 19, no. 3, pp. 293–325, 1948.
- [39] F. Hutter, H. Hoos, and K. Leyton-Brown, "An efficient approach for assessing hyperparameter importance," in *Proc. 31st Int. Conf. Mach. Learn.*, vol. 32. Beijing, China, Jun. 2014, pp. 754–762.
- [40] C. Fawcett and H. H. Hoos, "Analysing differences between algorithm configurations through ablation," *J. Heuristics*, vol. 22, no. 4, pp. 431–458, 2016.

- [41] S. Sheikholeslami, M. Meister, T. Wang, A. H. Payberah, V. Vlassov, and J. Dowling, "AutoAblation: Automated parallel ablation studies for deep learning," in *Proc. 1st Workshop Mach. Learn. (EuroMLSys@EuroSys)*, Edinburgh, U.K., Apr. 2021, pp. 55–61.
- [42] M. López-Ibáñez and T. Stützle, "An experimental analysis of design choices of multi-objective ant colony optimization algorithms," *Swarm Intell.*, vol. 6, pp. 207–232, Jul. 2012.
- [43] J. de Nobel, D. Vermetten, H. Wang, C. Doerr, and T. Bäck, "Tuning as a means of assessing the benefits of new ideas in interplay with existing algorithmic modules," in *Proc. Genet. Evol. Comput. Conf. Companion*, 2021, pp. 1375–1384.
- [44] A. Klein, L. C. Tiao, T. Lienart, C. Archambeau, and M. Seeger, "Model-based asynchronous hyperparameter and neural architecture search," 2020, *arXiv:2003.10865*.
- [45] M. Birattari, *Tuning Metaheuristics—A Machine Learning Perspective* (Studies in Computational Intelligence), vol. 197. Berlin, Germany: Springer, 2009.
- [46] Z. Karnin, T. Koren, and O. Somekh, "Almost optimal exploration in multi-armed bandits," in *Proc. Int. Conf. Mach. Learn.*, 2013, pp. 1238–1246.
- [47] J. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl, "Algorithms for hyperparameter optimization," in *Proc. Int. Conf. Adv. Neural Inf. Process. Syst.*, vol. 24, 2011, pp. 2546–2554.
- [48] A. Klein, S. Falkner, S. Bartels, P. Hennig, and F. Hutter, "Fast Bayesian optimization of machine learning hyperparameters on large datasets," in *Proc. Int. Conf. Artif. Intell. Stat.*, 2017, pp. 528–536.
- [49] D. Golovin, B. Solnik, S. Moitra, G. Kochanski, J. Karro, and D. Sculley, "Google vizier: A service for black-box optimization," in *Proc. 23rd ACM SIGKDD Int. Conf. Knowl. Discov. Data Min.*, 2017, pp. 1487–1495.
- [50] F. Pfisterer, L. Schneider, J. Moosbauer, M. Binder, and B. Bischl, "YAHPO gym—Design criteria and a new multifidelity benchmark for hyperparameter optimization," 2021, *arXiv:2109.03670*.
- [51] D. R. Jones, "A taxonomy of global optimization methods based on response surfaces," *J. Global Optim.*, vol. 21, no. 4, pp. 345–383, 2001.
- [52] H. Jalali, I. Van Nieuwenhuysse, and V. Picheny, "Comparison of kriging-based algorithms for simulation optimization with heterogeneous noise," *Eur. J. Oper. Res.*, vol. 261, no. 1, pp. 279–301, 2017.
- [53] L. Breiman, "Random forests," *Mach. Learn.*, vol. 45, no. 1, pp. 5–32, 2001.
- [54] R. J. Samworth, "Optimal weighted nearest neighbour classifiers," *Ann. Statist.*, vol. 40, no. 5, pp. 2733–2763, Oct. 2012.
- [55] J. Demšar, "Statistical comparisons of classifiers over multiple data sets," *J. Mach. Learn. Res.*, vol. 7, pp. 1–30, Dec. 2006.



Julia Moosbauer is currently pursuing the Doctoral degree with the Chair for Statistical Learning and Data Science, Ludwig-Maximilians-Universität München, Munich, Germany.

She is a member of the Munich Center for Machine Learning. Her research focus lies in the interface between automated and explainable machine learning with the goal to increase transparency and trust into automated machine learning systems. She is also interested in hyperparameter optimization algorithms (in particular, Bayesian

optimization), algorithm configuration, multiobjective optimization, and (sequential) experimental design.

Martin Binder is currently pursuing the Doctoral degree with the Chair for Statistical Learning and Data Science, Ludwig-Maximilians-Universität München, Munich, Germany.

He is a member of the Munich Center for Machine Learning. He is mostly working on black-box optimization methods for automatic machine learning and hyperparameter optimization, with a focus on multifidelity optimization. He also works on deep learning, specifically self-supervised learning, for genome sequence classification and analysis.



Lennart Schneider is currently pursuing the Doctoral degree with the Chair of Statistical Learning and Data Science, Ludwig-Maximilians-Universität München, Munich, Germany.

His research focuses on automated machine learning, hyperparameter optimization, multiobjective optimization, and neural architecture search.



Florian Pfisterer received the master's degree in statistics from the Ludwig Maximilian University of Munich, Munich, Germany, in 2018, where he is currently pursuing the Ph.D. degree with the Statistical Learning and Data Science Group, Department of Statistics.

His research interests are in the field of automated machine learning, multiobjective optimization, and algorithmic fairness.



Marc Becker received the master's degree in geoinformatics from Friedrich-Schiller University Jena, Jena, Germany, in 2020.

Since 2020, he has been a Research Engineer with the Ludwig Maximilian University of Munich, Munich, Germany, and a Main Developer of the mlr3 optimization packages.



Michel Lang received the Ph.D. degree in statistics, Dortmund Technical University, Dortmund, Germany, in 2015.

He is a former member of the Munich Center for Machine Learning. He is currently a Scientific Manager of the Research Center Trustworthy Data Science and Security, Dortmund, Germany. He is the author of many popular R packages for machine learning and parallelization, e.g., mlr3 or batchtools. His research areas include machine learning, optimization, and software development.



Lars Kotthoff received the Doctoral degree in computer science from the University of St Andrews, St Andrews, U.K., in 2012.

He is an Assistant Professor of Computer Science with the University of Wyoming, Laramie, WY, USA. His research focuses on data-driven combinatorial optimization, automated machine learning, and applying machine learning in other disciplines, for example, materials science.



Bernd Bischl studied computer science, artificial intelligence, and data sciences in Hamburg, Germany; Edinburgh, U.K.; and Dortmund, Germany. He received the Ph.D. degree in statistics from Dortmund Technical University, Dortmund, in 2013 with a thesis on "Model and Algorithm Selection in Statistical Learning and Optimization."

He holds the Chair of Statistical Learning and Data Science with the Department of Statistics, Ludwig-Maximilians-Universität München, Munich, Germany, and he is the Co-Director of the Munich

Center for Machine Learning, one of Germany's national competence centers for ML. He is a member of ELLIS, and a Faculty Member of ELLIS Munich, an active developer of several R-packages, leads the "mlr" (Machine Learning in R) Engineering Group and is the Co-Founder of the Science Platform "OpenML" for open and reproducible ML. Furthermore, he leads the Munich branch of the Fraunhofer ADA Lovelace Center for Analytics, Data, and Applications, i.e., a new type of research infrastructure to support businesses in Bavaria, especially in the SME sector. His research interests include AutoML, model selection, interpretable ML, as well as the development of statistical software.