

Practitioner Motives to Select Hyperparameter Optimization Methods

Niklas Hasebrook

NIKLAS.HASEBROOK@GMAIL.COM

*Institute of Applied Informatics and Formal Description Methods, Karlsruhe Institute of Technology,
Karlsruhe, Germany*

Felix Morsbach

FELIX.MORSBACH@KIT.EDU

KASTEL Security Research Labs, Karlsruhe Institute of Technology, Karlsruhe, Germany

Niclas Kannengiesser

NICLAS.KANNENGIESSER@KIT.EDU

*Institute of Applied Informatics and Formal Description Methods, Karlsruhe Institute of Technology,
Karlsruhe, Germany*

KASTEL Security Research Labs, Karlsruhe Institute of Technology, Karlsruhe, Germany

Marc Zöller

MARC.ZOELLER@USU.COM

USU Software AG, Karlsruhe, Germany

Jörg Franke

FRANKEJ@CS.UNI-FREIBURG.DE

Machine Learning Lab, University of Freiburg, Freiburg, Germany

Marius Lindauer

M.LINDAUER@AI.UNI-HANNOVER.DE

Institute of Artificial Intelligence, Leibniz University Hannover, Hannover, Germany

Frank Hutter

FH@CS.UNI-FREIBURG.DE

Machine Learning Lab, University of Freiburg, Freiburg, Germany

Ali Sunyaev

SUNYAEV@KIT.EDU

*Institute of Applied Informatics and Formal Description Methods, Karlsruhe Institute of Technology,
Karlsruhe, Germany*

KASTEL Security Research Labs, Karlsruhe Institute of Technology, Karlsruhe, Germany

Abstract

Advanced programmatic hyperparameter optimization (HPO) methods, such as Bayesian optimization, have high sample efficiency in reproducibly finding optimal hyperparameter values of machine learning (ML) models. Yet, ML practitioners often apply less sample-efficient HPO methods, such as grid search, which often results in under-optimized ML models. As a reason for this behavior, we suspect practitioners choose HPO methods based on individual motives, consisting of contextual factors and individual goals. However, practitioners' motives still need to be clarified, hindering the evaluation of HPO methods for achieving specific goals and the user-centered development of HPO tools. To understand practitioners' motives for using specific HPO methods, we used a mixed-methods approach involving 20 semi-structured interviews and a survey study with 71 ML experts to gather

evidence of the external validity of the interview results. By presenting six main goals (e.g., improving model understanding) and 14 contextual factors affecting practitioners’ selection of HPO methods (e.g., available computer resources), our study explains why practitioners use HPO methods that seem inappropriate at first glance. This study lays a foundation for designing user-centered and context-adaptive HPO tools and, thus, linking social and technical research on HPO.

Keywords: Artificial Intelligence (AI), Automated Machine Learning (AutoML), Human-AI Collaboration, Hyperparameter Optimization (HPO), User-centered HPO

1. Introduction

The performance of machine learning (ML) models depends strongly on the choice of hyperparameter values (Bergstra and Bengio, 2012; Melis et al., 2018; Henderson et al., 2018; Chen et al., 2018; Zhang et al., 2021; Kadra et al., 2021). Finding hyperparameter values that maximize the performance of ML models is challenging for ML practitioners (e.g., engineers and scientists from academia and industry), even for highly experienced ones, because search spaces of hyperparameter values can be large and the relationships between ML model performance, hyperparameter values, and dataset are seldom apparent. ML practitioners need to test various hyperparameter values, often in a trial-and-error manner, to find those that contribute to the desired performance of ML models. This process of searching and testing hyperparameter values to meet specific requirements for ML models is called hyperparameter optimization (HPO).

Testing many hyperparameter values manually is often cumbersome, tedious, and error-prone. To help practitioners in automating HPO, several programmatic HPO methods, including grid search, random search, and Bayesian optimization, were developed. Practitioners can find equivalent hyperparameter values using most programmatic HPO methods or manual tuning. Nonetheless, HPO methods differ considerably in the way they search for hyperparameter values, which affects the suitability of the HPO method to satisfy specific requirements, such as replicability and high sample efficiency.

Replicability (i.e., the deterministic execution of HPO) is paramount for ML research to ensure comparability between ML models and reproducibility of prior results (Musgrave et al., 2020). Manual tuning is non-programmatic and influenced by unconscious factors, like personal experiences (Tversky and Kahneman, 1974; Kahneman and Klein, 2009), tacit knowledge (Polanyi and Sen, 2009), and intuition (Helmstadter, 1970), and thus it is not replicable at all. In contrast, even the simplest HPO methods, such as grid search, ensure replicability to a certain extent. When comparing ML models produced by HPO methods with different sample efficiencies, the comparison and resulting conclusions might be skewed, which is particularly a challenge in ML research (McKinney et al., 2020). The sample efficiency of Bayesian optimization is often superior to grid search (Snoek et al., 2012). Despite the apparent superiority of some HPO methods over others, ML practitioners often apply seemingly inferior HPO methods. For example, practitioners mostly prefer to perform manual tuning or choose grid search over the more sample-efficient Bayesian optimization (Bouthillier and Varoquaux, 2020).

The dominant use of such seemingly inferior HPO methods suggests that there are factors that inhibit the use of superior HPO methods or that practitioners have different motives for HPO. Such motives appear to be beyond the fulfillment of publicly-discussed goals of HPO and ML models, such as replicability, high sample efficiency, and optimal ML model performance (Claesen and De Moor, 2015). However, practitioners’ motives (i.e., specific combinations of goals and contextual factors) to use specific HPO methods remain unclear. To support the development of user-centric HPO methods and corresponding software tools, we need to understand the motives of practitioners pursued in HPO.

Research investigated HPO of ML models from a foremost technical perspective (Feurer and Hutter, 2019; Bischl et al., 2023). The development and improvement of programmatic HPO methods used to automate HPO have been driven by powerful mathematical approaches and their manifestations in HPO tools (e.g., Bergstra et al., 2015; Li et al., 2018b). The superiority of HPO methods over others is shown based on traditional performance metrics such as the minimization of the generalization error and sample efficiency (e.g., Zöllner and Huber, 2021; Gijssbers et al., 2022; Lindauer et al., 2022).

Taking a sociotechnical perspective on HPO can direct research in this field. Extant studies focus on three predominant goals. First, by understanding how to engage practitioners in HPO, the expertise and experience of practitioners can be leveraged to improve HPO efficiency (Lee and Macke, 2020; Wang et al., 2021c). Second, by better understanding how to guide practitioners through HPO, practitioners are aimed to be supported in their decisions and action in data science works (e.g., Wang et al., 2019; Crisan and Fiore-Gartland, 2021). Third, explainability and transparency of HPO can help practitioners better understand the execution of programmatic HPO methods (e.g., Drozdal et al., 2020; Zöllner et al., 2022). Although research supports our assumption that practitioners perform HPO to achieve goals beyond traditional technocentric ones, practitioners’ motives for engaging in HPO remain unclear. This unclarity inhibits the user-centered improvement of HPO methods and tools to support the attainment of individual practitioner goals and hinders the targeted analysis of human performance in HPO. In order to overcome this issue, we shed light on the possible answers to the following research question: *Why do ML practitioners choose different HPO methods?*

We applied a two-step research approach consisting of an interview study and a survey study. First, we conducted semi-structured interviews with 20 ML experts to unveil HPO methods applied by practitioners, the goals pursued by practitioners when applying those methods, and the contextual factors that influence practitioners’ decisions for HPO methods to attain their goals. Second, we performed an online survey for seven months with 71 participants to collect evidence for the external validity of the relevance of the HPO methods, goals, and contextual factors identified in the interviews.

This work has three main contributions. First, we present six principal goals (e.g., comply with target audience, increase ML model performance) pursued by practitioners in HPO and fourteen contextual factors (e.g., available compute resources, HPO method traceability) that can influence practitioner decisions for specific HPO methods. We thereby support user-centric research on HPO by providing a better foundation for interactions of practitioners with HPO methods. Researchers can use the set of identified goals to provide HPO methods

tailored to pursuing the different goals while still preserving the benefits of advanced HPO methods. Second, we present a mapping between goals pursued in HPO, HPO methods used to achieve these goals, and contextual factors that influence practitioner decisions for using HPO methods. We thereby deepen the understanding of why practitioners use different HPO methods. Our mapping informs practitioners using HPO methods and can be used to better align the development of HPO methods to practitioner needs. Third, we present an overview of the success perceived by practitioners when using HPO methods in particular contexts (i.e., configurations of contextual factors and goals). This overview can serve as a foundation to advance the level of automation of HPO tools.

The remainder of this work is structured into five sections. First, we describe the current state of research on HPO. Second, we describe the approach we applied to answer our research question. Third, we report our results, including four HPO methods, six goals, and fourteen contextual factors. Fourth, we discuss our principal findings, explain the contributions of this work, describe possible threats to the validity of our results, and outline future research directions. We conclude with our main takeaways.

2. Background

The research field of automated machine learning (AutoML; see, e.g., Hutter et al. (2019)) aims to automate all aspects related to creating ML models, for example, by programmatic HPO methods. AutoML research has taken a foremost technical stance on HPO (e.g., Feurer and Hutter, 2019; Bischl et al., 2023). Typical works (e.g., Bergstra et al., 2011; Jamieson and Talwalkar, 2016; Golovin et al., 2017; Falkner et al., 2018; Li et al., 2018b; He et al., 2018; Wang et al., 2021a) focus on the performance optimization of ML models in terms of smaller generalization errors, smaller ML model size, or lower latency.

In AutoML research, HPO is investigated from a mathematical perspective and is often treated as a black-box optimization problem: Given a problem instance in the form of a dataset and a loss function, a black-box optimizer (i.e., the HPO method) searches for hyperparameter values in a predefined search space optimizing a given metric. Most works focus on maximizing the predictive accuracy of ML models on validation data. Multi-objective optimization methods can be used to specify additional properties of the resulting ML model, such as low model complexity, fast inference and algorithmic fairness (Gardner et al., 2019; Binder et al., 2020; Karl et al., 2022; Dooley et al., 2022).

Grid search was one of the earliest HPO methods that could be executed programmatically. While easy to implement and parallelize, it became increasingly unsuited for modern HPO problems due to the curse of dimensionality (Bergstra and Bengio, 2012). At the beginning of HPO, practitioners were faced with small search spaces that included only a few hyperparameters. Since then, AutoML has moved to ever-increasing search spaces covering the construction of complete ML pipelines. In practice, not all hyperparameters have a similar impact on the final performance of ML models (Bergstra and Bengio, 2012; Rijn and Hutter, 2018). Due to the rigid search strategy, grid search tends to spend a large portion of the optimization budget on exploring irrelevant parts of the search space (Bergstra and Bengio,

2012). This limitation can be circumvented using random search, which simply samples hyperparameter values at random.

To cope with large search spaces, the pure exploration by grid search and random search has been complemented by the exploitation of knowledge of well-performing regions (Snoek et al., 2012). A well-established approach for combining exploitation and exploration is evolutionary optimization. Evolutionary optimization methods, a class of optimizers inspired by natural evolution, often perform well optimizing a black-box function (Olson and Moore, 2019). Alternatively, Bayesian optimization can be used. Bayesian optimization, an optimization method for noisy black-box functions, constructs an internal probabilistic model, mapping hyperparameter values to expected performance, to achieve a good balance of exploration and exploitation (Shahriari et al., 2016; Frazier, 2018; Garnett, 2023). Extant research focuses on increasing the sample efficiency (Bergstra et al., 2011; Hutter et al., 2011; Snoek et al., 2012), reducing the time for evaluating the objective function (Swersky et al., 2014; Domhan et al., 2015; Li et al., 2018b), or transferring knowledge from prior optimization runs on similar problem instances (Dyrmishi et al., 2019; Vanschoren, 2019).

Besides the technical perspective on HPO, studies on human perception of programmatic HPO methods focus on deepening the understanding of practitioner perceptions of the advantages and disadvantages of HPO methods (Gil et al., 2019; Wang et al., 2019, 2021c). Practitioners acknowledge the benefits of automation leading to faster turn-around time for building individual models and, therefore, higher productivity (Wang et al., 2019, 2021b). These automatically created models are often used by practitioners to create a first baseline model for further manual HPO tuning or to gain data insights (Wang et al., 2019, 2021b; Crisan and Fiore-Gartland, 2021). However, a recurring criticism is insufficient confidence in the functioning of programmatic HPO methods (Lee and Macke, 2020; Wang et al., 2019; Drozdal et al., 2020; Khuat et al., 2022). Even though practitioners acknowledge that programmatic HPO provides well-performing ML models (Wang et al., 2021b), they refuse to use them to not be accountable for ML models they do not understand (Drozdal et al., 2020). Missing confidence is mostly associated with the perceived black-box nature of some programmatic HPO methods leading to missing transparency of the method internals. Practitioners rather wished for augmentation of their daily data science work (e.g., through guidance) instead of automating it (Crisan and Fiore-Gartland, 2021).

Research on human-guided HPO focus on approaches to involve humans in programmatic HPO methods to improve HPO with dormant domain expertise (Wang et al., 2019). This requires identifying how and when to engage humans in HPO to achieve the best combination of automation and human knowledge (Crisan and Fiore-Gartland, 2021). Especially the engagement of practitioners in ML model development, including HPO, appears promising for a higher level of automation. For other tasks, including requirement analysis or data acquisition, practitioners prefer a strong human involvement with minimal automation (Wang et al., 2021c). Interactions of practitioners with software tools for programmatic HPO were further structured into different modes of cooperation between practitioners and software tools, ranging from manual tuning to full automation, in the literature (Lee and Macke, 2020; Crisan and Fiore-Gartland, 2021; Wang et al., 2021c).

Even though extant literature offers important insights into the perception of programmatic HPO methods by practitioners and their interaction with it, prior research basically compares manual tuning with programmatic HPO methods. A distinction between the vastly different programmatic HPO methods is not performed, making it very hard to understand the reasons for the selection of HPO methods by practitioners. Furthermore, interactions proposed in the literature are often not validated with actual practitioners (e.g., Lee and Macke, 2020; Gil et al., 2019).

Despite the valuable contributions of extant publications, from a technical and sociotechnical perspective to the knowledge base of HPO, it still remains unclear why practitioners chose specific HPO methods. To learn the motives of practitioners for selecting HPO methods is the goal of this study.

3. Methods

We applied a mixed-methods research approach incorporating two principal steps: First, we conducted semi-structured interviews with ML experts. Second, we performed a survey using an online questionnaire to collect evidence for the external validity of the interviews.

3.1 Semi-structured Interviews with Machine Learning Experts

To identify practitioners’ principal goals pursued in HPO and to understand decisions for specific HPO methods to achieve these goals, we chose an exploratory, qualitative research approach and conducted semi-structured expert interviews.

Data Gathering We contacted potential interviewees among contacts from ongoing research projects, authors of scientific studies, and company contacts. Potential interviewees were selected to have heterogeneous experiences with HPO and ML, ranging from novices to experts and different ML research areas. Among the contacted potential interviewees, 20 agreed to participate in our study. The ML experts were associated with thirteen different organizations and had an average work experience in ML of about five years. For more detailed information about the interviewees, please see Table 1.

Field	Highest Degree of Education	Years of Experience	Research Areas
Academia (14)	Bachelor (2)	< 2 (4)	CV (6)
Industry (6)	Master (16)	2–4 (6)	NLP (3)
	PhD (2)	5–7 (7)	RL (3)
		> 7 (3)	TS (2)

CV: Computer vision, NLP: Natural language processing, RL: Reinforcement learning, TS: Time Series

Table 1: Overview of the demographic data of the 20 interviewees in our study. The numbers in parentheses show the number of interviewees with the respective characteristics. The interviewees could name multiple research areas.

In the interviews, we first introduced interviewees to the background and goals of our study. Second, we asked the ML experts to name HPO methods they used in ML projects. Finally, we gathered insights into why they selected HPO methods for optimization and asked about contextual factors that impacted their HPO method selection. We interviewed the ML experts under consideration of methodological recommendations, for example, to be open-minded and not to bias interviewees in their responses (Gorden, 1975; Louise Barriball and While, 1994; McIntosh and Morse, 2015). The interviews took between 18 and 61 minutes, with an average time of 31 minutes.

Data Analysis Two co-authors independently analyzed the transcripts using thematic analysis (Braun and Clarke, 2006, 2012). The thematic analysis comprises six steps: (1) familiarize yourself with the data, (2) generate initial codes, (3) search for themes, (4) review themes, (5) define and name themes, and (6) produce the report.

After familiarizing ourselves with our transcripts (Step 1), we coded them (Step 2) to identify HPO methods applied by practitioners, to extract practitioners’ goals targeted in HPO, and to reveal contextual factors that impacted the interviewees’ decisions for HPO methods. We incorporated contextual factors to better understand what influenced practitioner decisions for HPO methods (Gigerenzer and Brighton, 2009). Two authors independently read the transcripts, identified quotes relevant to our study, and labeled the quote with a name, a so-called code, that expresses a potentially relevant HPO method, goal, or contextual factor. The first coding iteration revealed 241 preliminary codes. Next, we harmonized these codes into 21 distinct codes so that no different codes had the same semantics. For example, we merged the contextual factors *knowledge about Bayesian optimization* and *knowledge about grid search* into the contextual factor *HPO method comprehension*. We developed preliminary themes (Step 3) to group the harmonized codes. If a code did not suit an existing theme, we created a new theme. For example, we assigned the contextual factor *available computational resources* to the theme *technical environment*, while we created a new theme *own knowledge* for the contextual factor *HPO method comprehension*. Our set of themes was comprised of four themes associated with HPO methods, six themes associated with goals, and three themes associated with contextual factors. In Step 4, we reviewed the preliminary themes within the author team. Subsequently, we developed an intuitive name for each theme and a precise description (Step 5). Finally, we collated the set of 13 themes into three categories: HPO methods, principal goals, and contextual factors (Step 6).

3.2 Online Survey

We conducted a survey study based on an online questionnaire to collect evidence for the external validity of the interview study results and to learn whether ML practitioners perceive that they succeeded in achieving their goals through their decisions to use specific HPO methods.

Questionnaire Structure Our online questionnaire was structured into four sections: *Introduction*, *HPO Methods and Goals*, *Contextual Factor Integration*, and *Demographics*. In the *Introduction* section, we described the motivation for and the structure of the questionnaire. In *Methods and Goals*, we showed participants a matrix that listed all identified

goals and HPO methods. Participants were asked to select all pairs of HPO methods and goals, which reflected that they used a given HPO method to achieve specific goals. In a second question, we asked participants to indicate, for each selected pair, whether or not they felt they had successfully achieved each goal. In the *Contextual Factor Integration* section, we wanted to learn how influential practitioners perceived the contextual factors in their HPO method selections. For each pair of HPO methods and goals previously selected, we asked the participants to express their perceived influence of each contextual factor on the selection of an HPO method to achieve a particular goal on a five-point Likert scale—0 represents very low perceived influence, two corresponds to a neutral response (i.e., the contextual factor was not perceived as influential), and four represents very high perceived influence. In the *Demographics* section, we aimed to learn more about the background and ML experience of our study participants.

Data Gathering To solicit participants for our survey, we contacted ML practitioners via email and promoted our study on social media platforms. Over the course of seven months, a total of 71 participants completed the *HPO Methods and Goals* section, of which 29 participants discontinued the questionnaire after completing the subsequent *Contextual Factor Integration* section. 31 participants completed the entire questionnaire.

The 31 participants who completed the questionnaire came from seven countries. Most participants worked in large organizations with more than 500 employees, including automotive companies, universities, and companies specializing in IT support and services. Table 2 shows more demographic details about the participants who completed the questionnaire.

Field	Highest Degree of Education	Years of Experience	Research Areas
Academia (15)	High School (1)	< 2 (5)	CV (10)
Industry (16)	Bachelor (1)	2–4 (14)	TD (16)
	Master (14)	5–7 (7)	NLP (10)
	Diploma (2)	8–10 (2)	TS (5)
	PhD (13)	> 10 (3)	Other (8)

CV: Computer vision, TD: Tabular data, NLP: Natural language processing, TS: Time series analysis

Table 2: Overview of the demographic data of the 31 participants that completed our survey study. The numbers in parentheses show the number of interviewees with the respective characteristics. The participants could choose multiple research areas.

Data Analysis By analyzing the survey responses, we sought to learn how frequently practitioners tend to choose which HPO methods to pursue specific goals, given which contextual factors. We extracted the number of identical responses and interrelated them.

Several survey respondents aborted the questionnaire before completing it. We filtered out incomplete responses by identifying the last completed study section. This enabled us to incorporate incomplete survey responses into our analysis. By including incomplete

responses, we were able to increase the number of responses to our questionnaire by 40 participants who would otherwise have been excluded from the study. After the analysis, we created visualizations that depict answer frequencies and their relationships.

4. Results

The study participants, including the interviewees and the survey participants, applied four principal HPO methods to achieve six goals under consideration of fourteen contextual factors (see Sections 4.1.1–4.1.3). Our study participants pursued their goals with varying self-reported success rates (see Subsection 4.2).

4.1 Hyperparameter Optimization Practices

In the following, we first briefly introduce the used HPO methods. Second, we report the goals pursued by the study participants and which of the four HPO methods practitioners used to achieve which goals. Third, we introduce fourteen contextual factors and how they can influence practitioners in their decisions for HPO methods to achieve specific goals.

4.1.1 HYPERPARAMETER OPTIMIZATION METHODS USED BY PRACTITIONERS

The participants in our interview study primarily applied four HPO methods: manual tuning, grid search, random search, and Bayesian optimization.

Manual Tuning Manual tuning refers to a set of HPO methods where a practitioner decides hyperparameter values based on personal explicit and implicit knowledge and external influences (e.g., results from literature). The dependence of manual tuning on the practitioner experience and even unconscious rationals for decisions make manual methods very individual to practitioners, make the explication of applied methods difficult, and, thus, are hardly replicability (Musgrave et al., 2020). Usually, only intermediate data (e.g., used hyperparameter values) can be used to replicate manual tuning, while reasons for the selection of these values remain unclear. Because formalization of HPO methods within manual tuning is difficult, the number of strategies for manual tuning applied by practitioners is unknown. In addition, the difficulty in explicating HPO methods used in manual tuning makes it hard to evaluate their sample efficiency. Therefore, we are not aware of any hard evidence for manual tuning to outperform advanced methods, such as Bayesian optimization. However, extant publications show that advanced HPO methods can outperform methods in manual tuning in certain use cases (Feurer et al., 2016; Melis et al., 2018; Chen et al., 2018; Zhang et al., 2021; Kadra et al., 2021).

Grid Search Grid search refers to the process of evaluating the Cartesian product of a finite set of values for each hyperparameter. Every possible combination of hyperparameter values included in the defined subset of the search space is evaluated (Montgomery, 2017) and thus it does not scale well with the number of hyperparameters. Grid search allows replicating an experiment because a deterministic procedure selects hyperparameter values to be evaluated. For replication, the originally applied search space and discretization strategy must be known. The sample efficiency of grid search tends to be lower compared to random

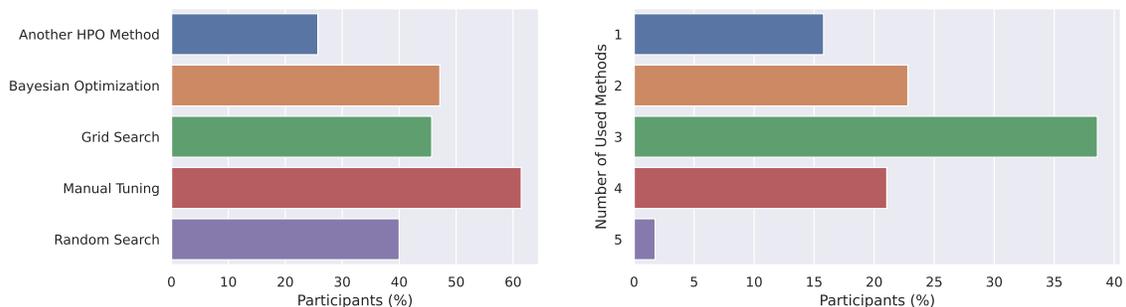
search and Bayesian optimization (Snoek et al., 2012; Eggenberger et al., 2013; Turner et al., 2021), in particular, because it cannot make use of the low effective dimensionality of HPO problems (Bergstra et al., 2011).

Random Search Random search refers to the process of sampling random hyperparameter values from a defined search space until a budget is exhausted (Bergstra and Bengio, 2012). Random search can be replicated if the used search space, the randomness generator, and the corresponding seed are known. Random search reaches better sample efficiency in high-dimensional search spaces than grid search if some hyperparameters have a larger influence on the performance of the ML model than others (Bergstra and Bengio, 2012).

Bayesian Optimization Bayesian optimization refers to the process of using a sequential approach based on a surrogate model to find appropriate hyperparameter values for an ML model in a defined search space (e.g., Brochu et al., 2010; Shahriari et al., 2016; Feurer and Hutter, 2019; Garnett, 2023). HPO based on Bayesian optimization can be replicated with fixed random seeds if the search space, the acquisition function, and the surrogate model, including the hyperparameters of the surrogate model, are known.

Survey Results Most survey participants used manual tuning, followed by Bayesian optimization, grid search, and random search. 25% of the survey participants also used other HPO methods that are out of the scope of this work (see Figure 1a). Most survey participants stated to have used at least three HPO methods in their past ML projects; about 2% have even used at least five HPO methods. Only roughly 15% of the survey participants have used just a single HPO method (see Figure 1b).

Even though literature indicates that Bayesian optimization yields better results than grid search and random search in shorter time (e.g., Snoek et al., 2012; Turner et al., 2021), practitioners tend to use seemingly inferior HPO methods. Practitioners appear to not only aim at finding hyperparameter values for optimal ML model performance but pursue manifold goals.



(a) Percentage of survey participants that already used the according HPO method. (b) Number of different HPO methods used by participants.

Figure 1: Overview of HPO methods used by 71 study participants.

4.1.2 GOALS OF PRACTITIONERS PURSUED WITH DIFFERENT HPO METHODS

We identified six goals that the participants pursued through HPO (see Table 3). In the following, we first introduce each goal based on our interview results. Then, we describe the results of our survey study.

Code	Description
Comply with Target Audience	The state where the applied HPO method and the resulting ML model fulfill the expectations of addressees
Decrease Necessary Computations	The state where an ML model is trained with an HPO method that requires less computational resources than other methods but still is sufficiently useful for a given purpose
Decrease Practitioner Effort	The state in which a practitioner applies an HPO method for training an ML model that requires fewer resources compared to other HPO methods (e.g., time for learning a new HPO method, for implementing corresponding software tools)
Increase Model Comprehension	The state where a practitioner is able to predict changes in an ML model’s behavior caused by altering hyperparameter values based on an understanding of the ML model’s inner workings
Increase Model Performance	The state where a refined version of an ML model outperforms its original version in terms of a specified metric
Satisfy Requirements	The state where the development and training of an ML model satisfies social and technical demands imposed by stakeholders

Table 3: Overview of goals practitioners pursue in HPO

Comply with Target Audience The goal *comply with target audience* refers to the alignment of individuals’ behaviors with the behavior expected by the target audience. Three interviewees stated to have decided on HPO methods in order to comply with the expectations of their target audiences regarding applied HPO methods and the resulting ML model. Two other interviewees, both from academia, considered Bayesian optimization uncommon in their research communities and saw the need to explain it in scientific papers about their ML model. This would require additional explanations of Bayesian optimization, even though the authors assumed the exact HPO method was not relevant to their scientific work. Therefore, they decided to use grid search as they assumed this HPO method to be well-known in their research communities.

Decrease Necessary Computations Extensive searches for optimal hyperparameter values in large search spaces usually require a vast amount of computing resources. A decrease in necessary computation for HPO is achieved when an HPO method is applied that requires fewer compute resources than other methods but is still sufficiently useful.

“This whole method was already super, super expensive [...] and if you would make again hyperparameter optimization, then it becomes even more expensive.” (Interviewee #8)

Decrease Practitioner Effort Practitioners choose HPO methods to reduce their overhead, for example, in terms of the additional time required to understand the HPO method or to integrate the HPO method into workflows. To decrease their efforts in HPO, the interviewees applied grid search and manual tuning. In particular, practitioners stated to have applied manual tuning to avoid efforts related to the integration of libraries for programmatic HPO methods into cluster infrastructures.

“HPO is time-consuming sometimes, because it requires some extra lines of code to wrap all your models with this HPO method and then set up the scripts to run them on a cluster.” (Interviewee #5)

Increase Model Comprehension Increasing model comprehension refers to reaching the state where a practitioner can predict changes in an ML model’s behavior caused by tuning hyperparameter values based on an understanding of the ML model’s inner workings. To increase their model comprehension, the interviewees reported having applied manual tuning. The interviewees claimed that manual tuning improves their understanding of hyperparameter influences on ML models because they can formulate hypotheses about hyperparameter influences and evaluate them immediately. The interviewees explained to iteratively improve their model comprehension by tuning hyperparameter values, observing the influences of these values on their ML models, and testing their hypotheses.

Increase Model Performance ML model performance is increased when a refined version of the ML model outperforms its original version in terms of a specified metric. The interviewees chose manual tuning, grid search, random search, and Bayesian optimization HPO methods to achieve this goal for example, to prototype novel ML models.

“If the only concern is to find the best model possible and no one asks how I got there, and I do not have a lot of time, I probably would use a random search.” (Interviewee #6)

Satisfy Requirements The goal to satisfy requirements refers to reaching the state in which the development and training of an ML model fulfill social and technical constraints imposed by stakeholders (e.g., business clients, ethics commissions) and the environment (e.g., available compute resources). Ten interviewees described that their decisions for HPO methods were influenced by the goal of fulfilling such requirements. For example, one interviewee reported preferring manual tuning to meet hard-to-formalize requirements, such as a smooth behavior of the model output.

Survey Results Our survey results indicate that all goals extracted from the interview study are also pursued by the survey participants (see Figure 2). More than 97% of the survey participants pursued the goal *increase model performance* and 77% aimed to achieve *decrease*

necessary computations or *decrease practitioner effort* (76%). 69% of the practitioners aimed to *increase model comprehension*. The least pursued goals are *satisfy requirements* (63%) and *comply with target audience* (58%).

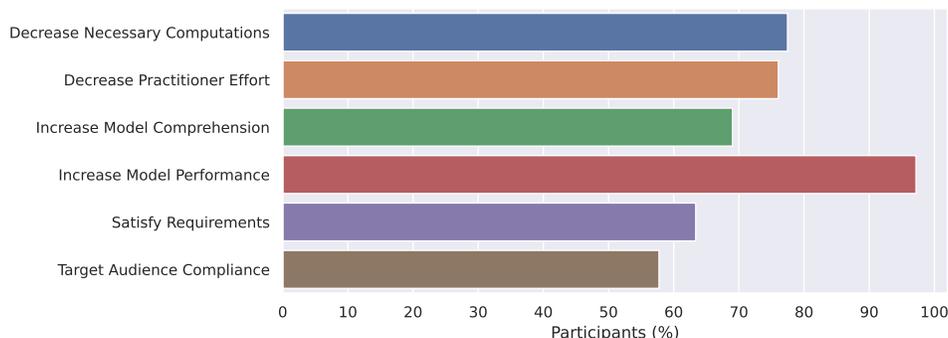


Figure 2: Relative frequency of pursued goals (all survey participants included that answered principal section two of the questionnaire).

Figure 3 shows how often the survey participants used HPO methods to reach a specific goal. 70% of the survey participants tried to decrease necessary computations by using Bayesian optimization and 60% by other HPO methods. About 50% of the participants tried to decrease the necessary computations by applying manual tuning. Random search and grid search were least often used, with 40% and 30%, respectively. Decreasing practitioner effort was only of interest for less than 50% of the participants, with very similar responses for all HPO methods (yet a notable exception of manual tuning with only 35%). The survey participants mainly used manual tuning to increase their comprehension of the ML models. Interestingly, participants also tried to use Bayesian optimization 2–3 times more often than grid or random search to achieve this goal despite its black-box nature. Increasing model performance is the most important goal for HPO. About 70% of the participants pursued this goal independently of the HPO method. Only random search was less often used to achieve this goal (< 60%). Requirements were mostly aimed to be satisfied using another HPO method, manual tuning, or Bayesian optimization (33%). Grid search (19%) and random search (11%) were rarely selected to achieve this goal. *Comply with target audience* was the least relevant goal for the participants. ~40% of the participants used another HPO method for this goal. There are no apparent differences between manual tuning and Bayesian optimization visible, with 33% of the participants using the respective method. Only random search and manual tuning were very seldom used to achieve this goal (~15%).

Apparently, participants use different HPO methods even though they aim to achieve the same goal and a clear mapping of HPO methods to goals is not possible. For example, one interviewee chose manual tuning to increase ML model comprehension while another interviewee preferred Bayesian optimization.

“*Because especially when entering new areas, we would like to understand step by step what is working and what is not.*” (Interviewee #9)

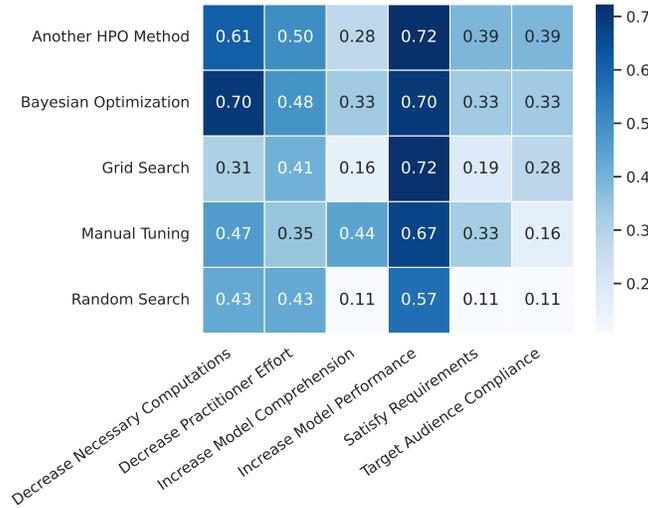


Figure 3: Frequency of goal and HPO method combinations. Per cell, all presented values are normalized to the number of participants having applied the respective HPO method (all survey participants included that answered principal section two of the questionnaire).

“I almost always select Bayesian optimization to get an idea in which region I find the [hyper-] parameters.” (Interviewee #15)

Due to the ambiguities of responses, the selection of HPO methods appears not to be simply explainable based on the practitioners’ goals. Contextual factors, which influence practitioners in their decision for HPO methods to reach goals, seem to be important to understand the motives of practitioners to use specific HPO methods.

4.1.3 CONTEXTUAL FACTORS THAT INFLUENCE HPO METHOD SELECTIONS

We identified fourteen contextual factors that can influence practitioner decisions for using HPO methods to achieve specific goals. These contextual factors can be grouped into three themes: *own knowledge*, *social environment*, and *technical environment*.

Own Knowledge Practitioner decisions for HPO methods depend on their own knowledge, which refers to the internal knowledge of a practitioner about HPO and ML models that guides them. We identified three contextual factors related to own knowledge: *HPO method comprehension*, *ML model comprehension*, and *personal experiences* (see Table 4).

HPO method comprehension refers to the degree to which practitioners understand how HPO methods work and how to apply them. Practitioners tend to neglect HPO methods they do not sufficiently understand. Two interviewees stated to have disregarded Bayesian optimization because they feel to have not sufficiently understood its inner workings. Another interviewee perceived grid search to be faster to implement and easier to use compared to Bayesian optimization, because using the latter would have required the interviewees to learn

Code	Description
HPO Method Comprehension	The self-perceived level of knowledge a practitioner has about the inner workings of an HPO method
ML Model Comprehension	The self-perceived level of understanding about the inner workings of an ML model with which a practitioner is able to explain changes in the behavior of the model caused by altering hyperparameter values
Personal Experience	The available internal knowledge that has been generated by past activities (e.g., personal best practices to solve a specific problem type)

Table 4: Overview of principal contextual factors that can influence practitioner decisions for HPO methods related to own knowledge.

an HPO method they were not experienced with. Another interviewee perceived random search as uncontrolled, which caused them to decide against it. Two interviewees decided to use grid search because they perceived grid search as easy to understand and implement.

ML model comprehension refers to a practitioner’s ability to explain changes in an ML model’s behavior caused when altering hyperparameter values based on an understanding of the inner workings of the ML model. The perceived degree of ML model comprehension plays an important role. Interviewees who perceived their ML model comprehension as high stated to have chosen manual tuning. Due to their deep ML model comprehension, these interviewees claimed that they are able to find appropriate sets of hyperparameter values without the need for extensive HPO. The interviewees perceived programmatic HPO methods as not taking advantage of the known effects of hyperparameters:

“Relationships between hyperparameters are often deducible, but optimizers [here: HPO methods] usually do not support functionalities for this.” (Interviewee #1)

Interviewees who perceived their ML model comprehension as low tended to use random search or Bayesian optimization. Low model comprehension made it difficult for interviewees to predict the challenges they will encounter in HPO. To better react to unforeseen challenges, practitioners choose manual tuning. For example, manual tuning can facilitate spotting and correcting mistakes when errors occur during the development of a novel model type because feedback loops are faster compared to those of programmatic HPO methods:

“Because we altered the standard architecture as a whole, we were not really sure what problems we will face. So that was one of the reasons to stick with manual tuning.” (Interviewee #3)

Personal experiences refers to the available internal knowledge that a practitioner generated through past activities (e.g., personal best practices for solving a specific type of problem). Practitioners tend to use HPO methods with which they have positive experiences:

“I have also made good experiences with it [here: Bayesian optimization] in a previous paper ”
 (Interviewee #2)

Social Environment The choice for an HPO method is also influenced by the social environment of ML practitioners, especially by four contextual factors (see Table 5): *acceptance of advanced HPO methods, literature, shared opinions, and tension for resources.*

Code	Description
Acceptance of Advanced HPO Methods	The extent to which advanced HPO methods (e.g., Bayesian optimization) are valued by a target group
Literature	The knowledge acquired on the basis of published text documents (e.g., articles, blog entries, and papers)
Shared Opinions	The knowledge acquired on the basis of advice by peers (e.g., colleagues)
Tension for Resources	The degree to which limited compute resources cause conflicts between practitioners regarding the allocation of those resources

Table 5: Overview of principal contextual factors that can influence practitioner decisions for HPO methods related to the social environment.

Acceptance of advanced methods refers to the extent to which advanced HPO methods such as Bayesian optimization are valued by a target group. Low acceptance of advanced HPO methods in a community targeted by a practitioner can make them choose manual tuning or avoid extensive HPO entirely. For example, an academic stated that they perceived the use of advanced HPO methods and extensive HPO as not being valued by their community. According to the interviewee, their community encourages the use of pre-trained ML models in combination with manual fine-tuning to avoid extensive HPO. Although the interviewee perceived Bayesian optimization as more suitable, they felt discouraged by the attitude of their community and applied manual tuning instead.

Shared opinions refers to the knowledge acquired on the basis of advice by peers (e.g., colleagues). The interviewees explained to have chosen HPO methods that are considered as commonly used in their labs or by their peers. In various communities, different HPO methods are applied so frequently that their use becomes habitual. For example, manual tuning was commonly used in one research group, while Bayesian optimization was considered the primarily applied method in another one. The interviewees associated with those communities applied the respectively manifested HPO methods. This indicates that the immediate social environment has a noticeable influence on practitioners’ HPO method choices.

Literature refers to the knowledge acquired on the basis of published text documents (e.g., articles, blog entries, papers). Practitioners are guided in their choice of HPO methods by recommendations from the literature on ML models similar to their own. All practitioners

that primarily based their decisions on literature, chose Bayesian optimization because it attests Bayesian optimization to a high sample efficiency (e.g., Turner et al., 2021).

Tension for shared resources refers to the degree to which limited compute resources cause conflicts between practitioners. The availability of only shared resources can cause tensions among colleagues, for example, when practitioners need to compete for computing resources to perform HPO. For example, such tensions caused one scientist in academia to choose manual tuning in order to avoid arguing with colleagues over computing resources.

Technical Environment Contextual factors associated with the technical environment refer to technical boundaries, such as limited computational resources, that guide a practitioner in selecting an HPO method. The interviewees stated seven contextual factors associated with the technical environment (see Table 6): *available compute resources, cost of objective function, HPO method traceability, HPO setup readiness, parallelization possibilities, search space size, and usability of HPO tools*.

Code	Description
Available Compute Resources	The amount of compute resources available for HPO
Cost of Objective Function	The amount of compute resources required to evaluate a single point within the hyperparameter space
HPO Method Traceability	The extent to which the sequence of sample points can be backtraced or predicted
HPO Setup Readiness	The degree to which HPO libraries and test environments are ready to use (e.g. preinstalled HPO libraries on the cluster)
Parallelization Possibilities	The degree to which multiple independent ML models can be simultaneously evaluated
Search Space Size	The number of possible hyperparameter value combinations
Usability of HPO Tools	The perceived ease with which practitioners can achieve their goal by using an HPO method and corresponding implementations

Table 6: Overview of principal contextual factors that can influence practitioner decisions for HPO methods related to the technical environment

Available compute resources refers to the amount of compute resources available for HPO. Practitioners choose manual tuning when faced with limited available compute resources. They perceive that in combination with a high level of model comprehension, they can outperform programmatic HPO methods. If the available compute resources are too scarce, the exploration of large search spaces is not possible. Practitioners need to reduce the search space, decrease the number of necessary function evaluations, or decrease computational cost per function evaluation (e.g., by low-fidelity approximations) to still perform HPO. To decrease the number of necessary function evaluations, three scientists stated to have used manual tuning as they were able to predict the impact of specific hyperparameter values on

model performance. Given limited available compute resources and a degree of ML model comprehension, the scientists perceived manual tuning as superior compared to Bayesian optimization and random search.

Two interviewees chose HPO methods depending on the *cost of the objective function* they sought to optimize (i.e., training of an ML model). The cost of the objective function refers to the amount of compute resources required to evaluate a single point within the hyperparameter space. Similar to limited compute resources, the interviewees chose manual tuning when faced with expensive objective functions. When the interviewees perceived their level of model comprehension as high, they found manual tuning more efficient.

HPO method traceability refers to the extent to which a sequence of sample points can be backtracked or predicted by the practitioner, which requires that the selection of samples by the HPO method is comprehensible for and reproducible by practitioners.

HPO setup readiness refers to the degree to which HPO libraries and test environments are ready to use (e.g., preinstalled HPO libraries on the cluster). Some of our study participants stated to not be willing to set up new HPO tools but rather use already set up tooling, regardless of the quality produced by the corresponding HPO method.

Parallelization possibilities of HPO methods refer to the degree to which independent ML models can be simultaneously evaluated. Limited parallelization possibilities can be caused, for example, by software license limitations. Two interviewees chose Bayesian optimization if parallelization of HPO was not possible due to a misconception of the sequential proceeding in Bayesian optimization. Another interviewee stated that they chose Bayesian optimization if their objective function is expensive and HPO parallelization is not possible.

Usability of HPO tools refers to the perceived ease with which practitioners can achieve their goal by using an HPO method and corresponding implementations. Within the scope of usability, practitioners demanded more automation of cumbersome tasks in HPO such as infrastructure orchestration:

“*What beats everything for me is that I have a dashboard that’s somewhere in the cloud that orchestrates my various agents, where I can sort of say online, ‘Start another agent on this machine,’ or that on the machine I just have to say, ‘Start another agent on this sweep here,’ and I don’t have to worry about the agents talking to each other or having a shared database running on some cluster. This functionality, it overrides everything. If I had some mega highly optimized Bayesian optimization tool that didn’t have that functionality, I wouldn’t use it.*” (Interviewee #14)

Survey Results Our survey study results show that each contextual factor was considered by at least 60% of the participants (see Figure 4). 90% of the survey participants considered the decision factors *personal experience*, *HPO setup readiness* and *search space size* in their selections of HPO methods. The least considered contextual factors are *acceptance of advanced methods*, *tension for resources*, and *parallelization possibilities*, as they were only relevant for less than 70% of the survey participants in their past ML projects. The remaining contextual factors, covering all three themes with, for example, *shared opinions*,

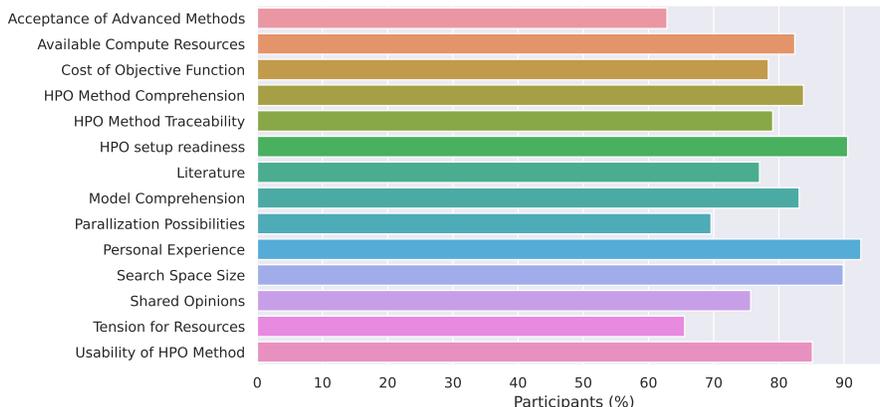


Figure 4: Percentage of participants that incorporated the individual contextual factors (all survey participants included that principal section three of the questionnaire).

model comprehension, and *HPO method traceability*, have been considered by 75–85% of the survey participants.

The identified contextual factors are of different self-perceived relevance for the selection of HPO methods (see Figure 5). The self-perceived relevance of contextual factors interdepends with the consideration of contextual factors. *Usability of HPO tools* and *search space size* are the most relevant contextual factors, closely followed by *personal experience* and *HPO setup readiness*. Further rather relevant context factors are *available compute resources*, *HPO method* and *model comprehension*, and the *cost of the objective function*. All contextual factors associated with the social environment are less relevant for the survey participants, with *tension for resources* being considered the least.

Figure 6 contains the self-perceived relevance of contextual factors for each of the evaluated HPO methods. The survey participants considered the *search space size*, *acceptance of advanced methods*, *available compute resources*, and *cost of the objective function* mostly when selecting Bayesian optimization. The selection of grid search was mostly influenced by *usability of HPO tools*, *HPO setup readiness*, and *search space size* with similar results for random search. Finally, survey participants considered their *personal experience*, *model comprehension*, and *HPO setup readiness* most relevant when selecting manual tuning.

4.2 Perceived Success of Using HPO Methods to Achieve Specific Goals

Practitioners appear to have individual motives to use specific HPO methods. Yet, some decisions for HPO methods may not produce the attempted results. To distinguish between the successful and unsuccessful experiences of practitioners in using HPO to reach their goals, we asked the study participants to what extent they perceive to have attained which goals with the selected HPO methods.

Figure 7 shows the success rate for each goal perceived by the survey participants. Roughly 75% of the participants responded to have successfully *increased model performance*, reached

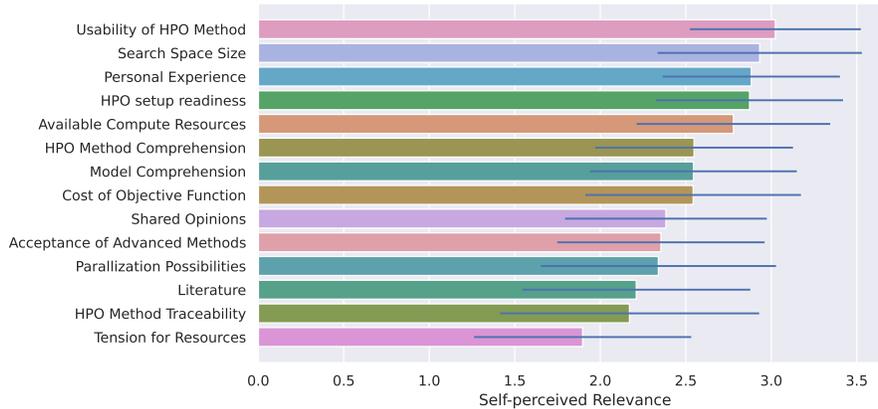


Figure 5: Overview of the average self-perceived relevance of contextual factors. Results are reported on a scale from 0 (very low) to 5 (very high) (all survey participants included that answered principal section three of the questionnaire). On average, no contextual factor was rated with a higher relevance than 3, which is why we masked out the values 4 and 5. Blue lines indicate error bars of one standard deviation.

complied with target audience, or satisfied requirements. 70% of the participants state that they were able to *increase their model comprehension* and 65% achieved the goal *decrease practitioners effort*. Only 55% of the participants perceived themselves as successful in *decreasing necessary computations*.

These numbers strongly vary when considering separate combinations of goal and HPO method (see Figure 8). The survey participants perceived themselves as not successful when they tried to *decrease necessary computations* using manual tuning, grid search, or random search. They perceived themselves as rather successful in reaching this goal when using Bayesian optimization or another HPO method. *Decreasing practitioner effort* was best achieved using grid search or random search according to the survey participants. Bayesian optimization and other HPO methods were perceived as less effective to decrease effort. A potential explanation could be that those HPO methods require some effort to set up in the first place. The participants perceived manual tuning as ineffective in decreasing their efforts. In contrast, participants perceived manual tuning as very helpful to *increase model comprehension*. Even though Bayesian optimization is considered a black-box optimization technique (Frazier, 2018), it was also perceived as suitable to increase model comprehension. Random search, grid search, and another HPO method were perceived as unsuitable for increasing model comprehension. Most participants did not perceive noteworthy differences between the effectiveness of HPO methods to successfully *increase ML model performance* and to *satisfy requirements*. Random search, grid search, and other HPO methods were successfully used by survey participants to achieve *comply with target audience*. Bayesian optimization and manual tuning were applied with lower success rates for this goal.

PRACTITIONER MOTIVES TO SELECT HPO METHODS

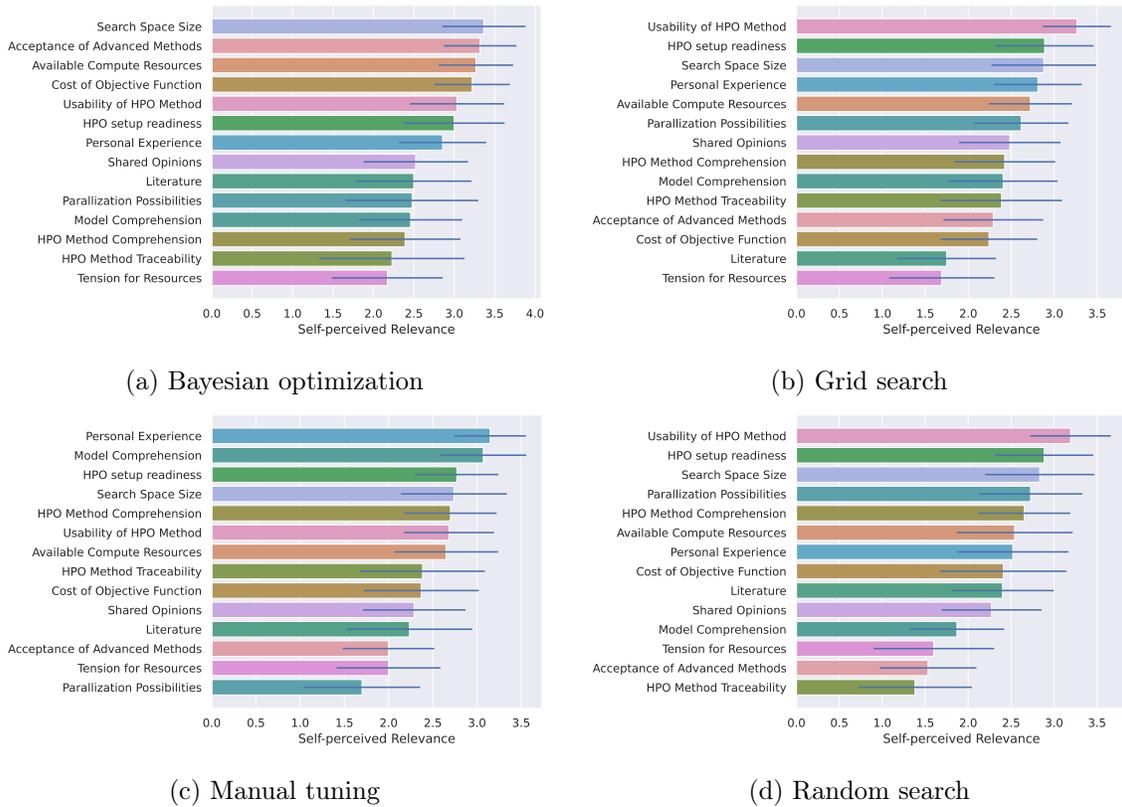


Figure 6: Overview of the average self-perceived relevance of contextual factors per HPO method. Results are reported on a scale from 0 (very low) to 5 (very high) (all survey participants included that answered principal section three of the questionnaire). Blue lines indicate error bars of one standard deviation.

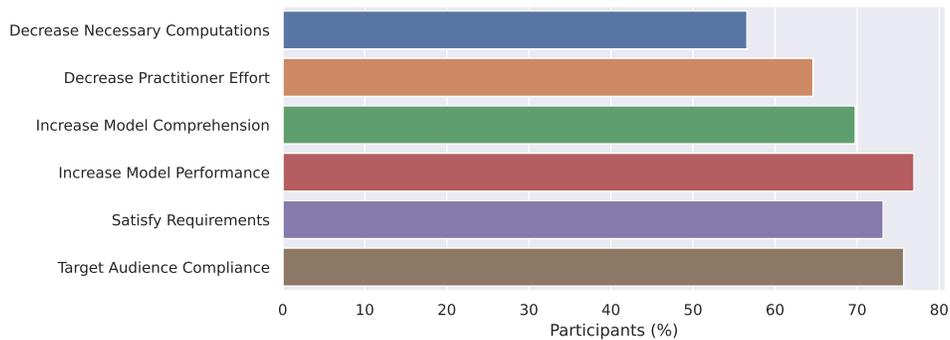


Figure 7: Self-reported success rate per goal (all survey participants included that answered principal section two of the questionnaire).

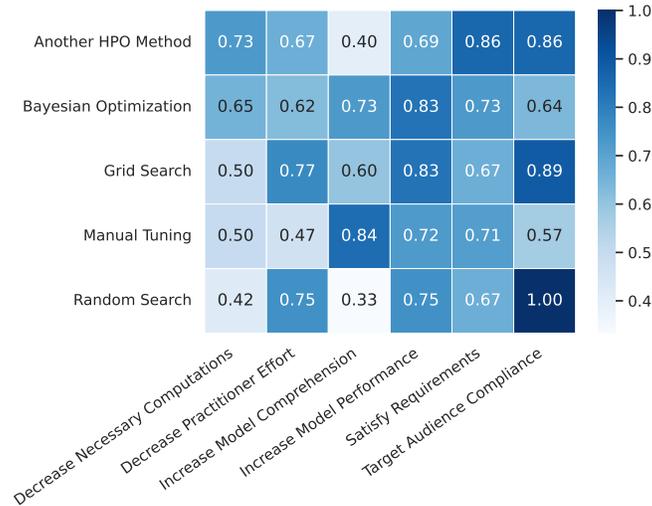


Figure 8: Study participants’ self-reported success rates per goal-method combination (based on responses of all survey participants that completed principal section two of the questionnaire).

5. Discussion

In this section, we summarize and discuss our principal findings based on the results presented in the previous section. We describe how our results contribute to HPO practice and research, describe future research directions, and discuss the limitations of our study.

5.1 Principal Findings

We conducted 20 semi-structured interviewees and an online survey with 71 participants to understand practitioner motives for choosing HPO methods. Our study reveals that practitioner motives comprise six goals (e.g., decrease practitioner effort, increase model comprehension), which practitioners pursue in the HPO of ML models, and fourteen contextual factors (e.g., HPO method comprehension, personal experience), which influence practitioners in their decisions for HPO methods to reach those goals.

The multitude of goals pursued by practitioners confirms our assumption that they have strong motives for HPO beyond improving ML model performance. Besides improving ML model performance in the first place, practitioners are second most interested in decreasing the required amount of necessary computations and their personal efforts to perform HPO. This perfectly aligns with the intended goals of HPO.

To decrease necessary computations, our study participants predominantly used Bayesian optimization. The frequent use of Bayesian optimization to decrease necessary computations suggests that practitioner perceptions of the benefits of Bayesian optimization are coherent with its benefits empirically shown in prior research (Turner et al., 2021).

Despite Bayesian optimization’s high sample efficiency, some practitioners prefer to use manual tuning to decrease the number of necessary computations. In particular, if practitioners assume that their ML model comprehension is high, they expect to outperform Bayesian optimization. Yet, it is difficult to compare manual tuning to programmatic HPO methods due to its reliance on a mixture of explicit and implicit knowledge that often cannot be fully extracted from observations of practitioner actions. One way to leverage this model comprehension could be to integrate user priors on the location of well-performing hyperparameter configurations into Bayesian Optimization (Souza et al., 2021; Hvarfner et al., 2022; Mallik et al., 2022).

Our study results show that various practitioners are interested in understanding their subject of work, including HPO methods and ML models, prior to using them. To better understand ML models, most practitioners opted for manual tuning instead of programmatic HPO methods and chose HPO methods they understood over methods they would have to study first. Participants perceived programmatic HPO tools as unsuited to advance their ML model comprehension. Even though many software packages for advanced HPO methods exist that can be used out of the box without understanding internals, practitioners apparently are reluctant to use methods they do not understand. In summary, practitioners tend to rely on their own knowledge rather than giving up control to insufficiently understood ML models and HPO tools.

To help practitioners to improve ML model comprehension, multiple software tools were designed. Such tools, predominantly HPO tools, mainly focus on providing measures for the impacts of hyperparameter values on the final ML model performance (e.g., Hutter et al., 2014; Biedenkapp et al., 2017; Moosbauer et al., 2021; Segel et al., 2023). With a focus on HPO methods, visual analytics aims to facilitate understanding the internal behaviors of HPO methods by visualizations for practitioners (e.g., Golovin et al., 2017; Biedenkapp et al., 2018; Zöllner et al., 2022; Sass et al., 2022).

Despite the existence of such tools to support practitioners in understanding ML models and HPO methods, practitioners still prefer to use manual tuning, which may have different reasons. The first reason may be that practitioners are unaware of HPO tools that can support ML model comprehension. A second reason may be that HPO tools do not fulfill the information needs of practitioners to increase their ML model comprehension because HPO tools mainly focus on the performance of ML models, which is, as shown in this study, only one of the manifold goals pursued by practitioners in HPO. A third reason may be that HPO tools themselves are hard to comprehend for practitioners (e.g., because those HPO tools implement complex HPO methods), which leads practitioners to prefer HPO methods they are familiar with.

Our study shows that the comprehensibility of HPO methods and the usability of HPO tools are paramount to practitioners. Practitioners tend to avoid using HPO methods that are difficult to integrate into workflows or require prior training. According to our study participants, these challenges occur, especially when using complex HPO methods. To make achievements of technocentric research on HPO actionable in practice, it is important to consider social aspects when developing HPO tools, such as seamless integration of HPO tools into workflows, prior experiences of practitioners who should use HPO tools, and under-

standing practitioners’ motivations for HPO. We propose three sociotechnical improvements for programmatic HPO methods based on our findings in the following that should be pursued in research in addition to further technical improvements with a focus on performance metrics.

Aid ML Model Comprehension To aid practitioners in gaining ML model comprehension, HPO libraries should generate reports about the behavior of different ML models, for example, the importance of individual hyperparameters. For the generation of such reports, many methods are already available, such as functional ANOVA (Hutter et al., 2014), ablation (Biedenkapp et al., 2017), local parameter importance (Biedenkapp et al., 2018), partial dependence plots (Moosbauer et al., 2021), and symbolic regressions (Segel et al., 2023). Alternatively, HPO tools could provide additional insights about model behavior. Especially for more complex search spaces used for building complete ML pipelines, information about the transformation of input data can provide additional model comprehension (Zöller et al., 2022). Including such reports in HPO libraries can facilitate leveraging the benefits of advanced HPO methods (e.g., high sample efficiency) while still helping practitioners to *increase ML model comprehension*.

Explain HPO Internals HPO tools should provide more support explaining their internal behavior to make them more comprehensible for practitioners. An easy approach would be simple visualizations of what hyperparameter values were actually evaluated using parallel coordinates plots (Golovin et al., 2017). More sophisticated approaches could provide information about the internal proceeding of their optimizers, for example, the surrogate model in Bayesian optimization (Biedenkapp et al., 2018). These measures could help educate practitioners about HPO methods and increase practitioners confidence in the functioning of programmatic HPO tools.

Integrate Practitioners’ ML Model Comprehension into HPO To increase the efficiency of HPO methods, they should allow the incorporation of comprehension of practitioners about ML models. Practitioners should be enabled to input their knowledge about behaviors of ML models into HPO libraries prior to HPO on a case-by-case basis. For example, practitioners could specify their perceived hyperparameter importance or influences between hyperparameters. Furthermore, practitioner knowledge could be directly incorporated into the search strategy of HPO methods. Promising work in this direction includes various methods for integrating prior knowledge into Bayesian optimization. This can be achieved by directly specifying priors about the location of the optimum (Li et al., 2020; Ramachandran et al., 2020; Souza et al., 2021; Hvarfner et al., 2022), or structural priors, for example, in the form of log-transformations of hyperparameters (Hutter et al., 2011), monotonicity constraints (Li et al., 2018a), or hyperparameter warping (Snoek et al., 2014).

5.2 Contributions to Practice and Research

Our main contributions are three-fold. First, we present six principal goals (e.g., comply with target audience, increase ML model performance) pursued by practitioners in HPO and fourteen contextual factors (e.g., available compute resources, HPO method traceability)

that can influence practitioner decisions for specific HPO methods. We thereby support user-centric research in HPO by offering a foundation to better understand practitioner motives for HPO. Researchers can use the set of identified goals to provide HPO methods tailored to pursuing specific goals while still preserving the benefits of advanced HPO methods. Moreover, research on human-in-the-loop in (Auto)ML can use our results to better describe contexts for information needs of practitioners (e.g., to increase the transparency of HPO tools) and to purposefully engage practitioners in programmatic HPO depending on goals and context factors. For example, practitioners may become more engaged in HPO when they aim to improve their ML model comprehension.

Second, we present a mapping between goals, HPO methods, and contextual factors that influence practitioner decisions for using HPO methods. We thereby deepen the understanding of why practitioners use different HPO methods. Our mapping informs researchers and developers of programmatic HPO methods and can be used to better align the development of new as well as the adoption of existing HPO methods to the needs of practitioners. This mapping revealed potential input parameters and their interrelationships that can be used for the further automation of HPO tools. For example, dedicated HPO tools can be developed for specific contextual factors and goals best suited to meet practitioner motives.

Third, we present an overview of the success perceived by practitioners when using HPO methods in particular contexts (i.e., configurations of contextual factors and goals). This overview can serve as a foundation to advance the level of automation of HPO tools. Moreover, we support research on user-centric HPO tools by grounding the usefulness of HPO methods with practitioners’ reality highlighting possibilities for HPO tool advancement. Practitioners’ self-perceived success in attaining goals under consideration of respective contextual factors may serve as a new lightweight benchmark metric that helps the sociotechnical improvement of HPO methods.

5.3 Limitations and Future Work

We performed semi-structured interviews, a qualitative and explorative research approach. Interviews mainly rely on the interviewees’ knowledge, perceptions, and capabilities to verbalize responses to our questions. The interviewed ML experts may bias our results despite our efforts to reduce such biases (e.g., by not asking leading questions or asking the interviewees for clarifications of statements). We aimed to minimize biases in the analysis of the interviews by having two scientists independently code the transcripts of the interviews and then discuss their codes to agree on a shared understanding. However, despite these efforts, we cannot guarantee to have eliminated biases in our results. Moreover, our results may not be comprehensive as ML practitioners that did not participate in our interview study may have additional goals and use HPO methods not mentioned by our interviewees. Additional interviews or focus group workshops should be conducted to increase the comprehensiveness of the findings presented in this work.

25% of the survey participants responded to have used other HPO methods than Bayesian optimization, manual tuning, grid search, and random search. Other HPO methods include evolutionary optimization, gradient-based optimization, or population-based approaches. Such other HPO methods were not mentioned frequently by the interviewees during the

initial interviews and, therefore, not explicitly included in the online survey. As we did not ask to record which other HPO methods participants used, potential insights for the other approaches could not be gained. The investigation of practitioner motives for choosing HPO methods not in the scope of this work offers an avenue for future research.

Future investigations of human decision-making in ML are a promising research direction that can help to improve AutoML by incorporating human knowledge. The interviewed practitioners reported actions they applied in HPO, which are agnostic to the selected HPO methods, such as choosing a set of hyperparameters to tune and defining corresponding search ranges. With similar reasoning, the interviewees applied similar actions for choosing HPO methods, hyperparameters, and hyperparameter values and stated they had attained their goals by taking those actions. Such a similar proceeding of practitioners in HPO allows for the assumption that best practices for actions taken during HPO exist. Since our interviewees mainly stated that they largely unconsciously compared HPO methods but achieved satisfactory outcomes, the identification of heuristics applied in their decision-making (Tversky and Kahneman, 1974; Gigerenzer and Brighton, 2009; Kahneman and Klein, 2009) in HPO appears to be of great potential to advance automation of HPO tasks in AutoML. By identifying the heuristics of practitioners in HPO, a better understanding of how practitioners use HPO methods can be reached to improve AutoML tools by an automated selection of the best suitable HPO methods for a given problem instance based on a specific goal and set of contextual factors.

6. Conclusion

This study sheds light on the ML practitioners’ motives for choosing different HPO methods. While programmatic HPO methods, such as Bayesian optimization, achieve high efficiency of the HPO process; practitioners sometimes opt for efficiency-wise inferior methods like manual tuning or grid search. To understand practitioners’ motives for their choices, we employed a two-step research approach consisting of semi-structured interviews and an online survey. Through thematic analysis of the interview transcripts and survey responses, we identified six principal goals pursued by practitioners in HPO and fourteen contextual factors that influence their decisions regarding HPO methods. The identified goals encompassed various aspects such as decreased practitioner effort, decreased necessary computations, or increased model comprehension.

The findings of this study can be leveraged to improve HPO practices and guide the development of user-centric HPO methods and software tools. By considering practitioners’ goals and contextual factors, researchers can refine existing HPO methods and tools and design new approaches that better meet the diverse motives of practitioners. Understanding the motives behind ML practitioners’ choices of different HPO methods is crucial to advancing HPO and developing effective and user-friendly HPO tools that cater to practitioners’ specific goals and contexts. By bridging the gap between technology advancements and practitioner needs, this research contributes to enhancing HPO practices. It promotes the broader adoption of efficient and reliable HPO methods in ML.

Our study calls for further research on HPO, particularly in exploring the engagement of practitioners, decision support systems for HPO, and enhancing transparency and explainability of programmatic HPO methods. We will build on the findings presented in this work and seek to identify human heuristics applied in HPO. After identifying human heuristics (e.g., Godbole et al., 2023), we aim to implement them in algorithms for AutoML and evaluate these algorithms in comparison to the performance of human decision-making.

Acknowledgments

We would like to thank all study participants for their time and valuable contributions, which formed the basis for this work. This work was supported by KASTEL Security Research Labs. Frank Hutter and Marius Lindauer acknowledge funding by the European Union (via ERC Consolidator Grant DeepLearning 2.0, grant no. 101045765, and ERC Starting Grant “ixAutoML”, grant no. 101041029, respectively). However, the views and opinions expressed are those of the author(s) only and do not necessarily reflect those of the European Union or the European Research Council. Neither the European Union nor the granting authority can be held responsible for them.



References

- James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *Journal of Machine Learning research*, 13(10), 2012.
- James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-parameter optimization. In *Advances in Neural Information Processing Systems*, volume 24, pages 2546–2554. Curran Associates, Inc., 2011.
- James Bergstra, Brent Komer, Chris Eliasmith, Dan Yamins, and David D. Cox. Hyperopt: a python library for model selection and hyperparameter optimization. *Computational Science & Discovery*, 8(1):014008, 2015.
- André Biedenkapp, Marius Lindauer, Katharina Eggensperger, Chris Fawcett, Holger Hoos, and Frank Hutter. Efficient parameter importance analysis via ablation with surrogates. In S. Singh and S. Markovitch, editors, *Proceedings of the 31st AAAI Conference on Artificial Intelligence*, volume 31, 2017.
- André Biedenkapp, Joshua Marben, Marius Lindauer, and Frank Hutter. Cave: Configuration assessment, visualization and evaluation. In *International Conference on Learning and Intelligent Optimization*, pages 115–130. Springer, 2018.
- Martin Binder, Julia Moosbauer, Janek Thomas, and Bernd Bischl. Multi-objective hyper-parameter tuning and feature selection using filter ensembles. In *Proceedings of the 2020*

- Genetic and Evolutionary Computation Conference, GECCO '20*, page 471–479. Association for Computing Machinery, 2020.
- Bernd Bischl, Martin Binder, Michel Lang, Tobias Pielok, Jakob Richter, Stefan Coors, Janek Thomas, Theresa Ullmann, Marc Becker, Anne-Laure Boulesteix, Difan Deng, and Marius Lindauer. Hyperparameter optimization: Foundations, algorithms, best practices, and open challenges. *WIREs Data Mining and Knowledge Discovery*, 13(2):e1484, 2023.
- Xavier Bouthillier and Gaël Varoquaux. Survey of machine-learning experimental methods at NeurIPS2019 and ICLR2020. Research report, Inria Saclay Ile de France, 2020.
- Virginia Braun and Victoria Clarke. Using thematic analysis in psychology. *Qualitative Research in Psychology*, 3(2):77–101, 2006.
- Virginia Braun and Victoria Clarke. *Thematic analysis*, page 57–71. American Psychological Association, 2012.
- Eric Brochu, Vlad M. Cora, and Nando De Freitas. A tutorial on bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning. *arXiv preprint arXiv:1012.2599*, 2010.
- Yutian Chen, Aja Huang, Ziyu Wang, Ioannis Antonoglou, Julian Schrittwieser, David Silver, and Nando de Freitas. Bayesian optimization in alphago. *arXiv preprint arXiv:1812.06855*, 2018.
- Marc Claesen and Bart De Moor. Hyperparameter search in machine learning. *arXiv preprint arXiv:1502.02127*, 2015.
- Anamaria Crisan and Brittany Fiore-Gartland. Fits and starts: Enterprise use of automl and the role of humans in the loop. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems, CHI '21*, New York, NY, USA, 2021. Association for Computing Machinery.
- T. Domhan, J. Springenberg, and F. Hutter. Speeding up automatic Hyperparameter Optimization of deep neural networks by extrapolation of learning curves. In Q. Yang and M. Wooldridge, editors, *Proceedings of the 24th International Joint Conference on Artificial Intelligence (IJCAI'15)*, 2015.
- Samuel Dooley, Rhea Sanjay Sukthanker, John P Dickerson, Colin White, Frank Hutter, and Micah Goldblum. On the importance of architectures and hyperparameters for fairness in face recognition. In *Workshop on Trustworthy and Socially Responsible Machine Learning, NeurIPS 2022*, 2022.
- Jaimie Drozdal, Justin Weisz, Dakuo Wang, Gaurav Dass, Bingsheng Yao, Changruo Zhao, Michael Muller, Lin Ju, and Hui Su. Trust in automl: Exploring information needs for establishing trust in automated machine learning systems. In *Proceedings of the 25th International Conference on Intelligent User Interfaces, IUI '20*, page 297–307, New York, NY, USA, 2020. Association for Computing Machinery.

- Salijona Dyrnishi, Radwa Elshawi, and Sherif Sakr. A decision support framework for automl systems: A meta-learning approach. In *2019 International Conference on Data Mining Workshops (ICDMW)*, pages 97–106, 2019.
- Katharina Eggenberger, Matthias Feurer, Frank Hutter, James Bergstra, Jasper Snoek, Holger Hoos, and Kevin Leyton-Brown. Towards an empirical foundation for assessing bayesian optimization of hyperparameters. In *NIPS workshop on Bayesian Optimization in Theory and Practice*, volume 10, 2013.
- Stefan Falkner, Aaron Klein, and Frank Hutter. Bohb: Robust and efficient hyperparameter optimization at scale. In *Proceedings of the 35th International Conference on Machine Learning*, pages 1437–1446. PMLR, 2018.
- Matthias Feurer and Frank Hutter. *Hyperparameter Optimization*, pages 3–33. Springer, Cham, 2019.
- Matthias Feurer, Matthias Klein, and Frank Hutter. Winning the automl challenge with auto-sklearn, 2016. URL <https://www.kdnuggets.com/2016/08/winning-automl-challenge-auto-sklearn.html>.
- Peter I. Frazier. A tutorial on bayesian optimization. *arXiv preprint arXiv: 1807.02811*, pages 1–22, 2018.
- Steven Gardner, Oleg Golovidov, Joshua Griffin, Patrick Koch, Wayne Thompson, Brett Wujek, and Yan Xu. Constrained multi-objective optimization for automated machine learning. In *2019 IEEE International Conference on Data Science and Advanced Analytics (DSAA)*, pages 364–373, 2019.
- Roman Garnett. *Introduction*, page 1–14. Cambridge University Press, 2023. doi: 10.1017/9781108348973.002.
- Gerd Gigerenzer and Henry Brighton. Homo heuristicus: Why biased minds make better inferences. *Topics in Cognitive Science*, 1(1):107–143, 2009.
- Pieter Gijsbers, Marcos L. P. Bueno, Stefan Coors, Erin LeDell, Sébastien Poirier, Janek Thomas, Bernd Bischl, and Joaquin Vanschoren. Amlb: an automl benchmark. *arXiv preprint arXiv:2207.12560*, 2022.
- Yolanda Gil, James Honaker, Shikhar Gupta, Yibo Ma, Vito D’Orazio, Daniel Garijo, Shruti Gadewar, Qifan Yang, and Neda Jahanshad. Towards human-guided machine learning. In *Proceedings of the 24th International Conference on Intelligent User Interfaces, IUI ’19*, page 614–624, New York, NY, USA, 2019. Association for Computing Machinery.
- V. Godbole, G. E. Dahl, J. Gilmer, C. J. Shallue, and Z. Nado. Deep learning tuning playbook, 2023. URL http://github.com/google-research/tuning_playbook. Version 1.
- Daniel Golovin, Benjamin Solnik, Subhodeep Moitra, Greg Kochanski, John Karro, and D. Sculley. Google vizier: A service for black-box optimization. In *Proceedings of the*

- 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '17*, page 1487–1495, New York, NY, USA, 2017.
- Raymond L. Gorden. *Interviewing: strategy, techniques, and tactics*. The Dorsey series in sociology. Dorsey Press, rev. ed edition, 1975.
- Yihui He, Ji Lin, Zhijian Liu, Hanrui Wang, Li-Jia Li, and Song Han. Amc: Automl for model compression and acceleration on mobile devices. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 784–800, 2018.
- Gerald Carl Helmstadter. *Research concepts in human behavior: Education, psychology, sociology*. Appleton-Century-Crofts, 1970.
- Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger. Deep reinforcement learning that matters. In *AAAI'18/IAAI'18/EAAI'18*, volume 32, 2018.
- Frank Hutter, Holger Hoos, and Kevin Leyton-Brown. Sequential model-based optimization for general algorithm configuration. In *International conference on learning and intelligent optimization*, pages 507–523. Springer, 2011.
- Frank Hutter, Holger Hoos, and Kevin Leyton-Brown. An efficient approach for assessing hyperparameter importance. In Eric P. Xing and Tony Jebara, editors, *Proceedings of the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pages 754–762, Beijing, China, 22–24 Jun 2014. PMLR. URL <https://proceedings.mlr.press/v32/hutter14.html>.
- Frank Hutter, Lars Kotthoff, and Joaquin Vanschoren, editors. *Automated Machine Learning - Methods, Systems, Challenges*. Springer, 2019.
- Carl Hvarfner, Danny Stoll, Artyr Souza, Marius Lindauer, Frank Hutter, and Luigi Nardi. pibo: Augmenting acquisition functions with user beliefs for bayesian optimization. In *International Conference on Learning Representations*, pages 1–15, 2022.
- Kevin Jamieson and Amee Talwalkar. Non-stochastic best arm identification and hyperparameter optimization. In *Artificial Intelligence and Statistics*, pages 240–248, 2016.
- Arlind Kadra, Marius Lindauer, Frank Hutter, and Josif Grabocka. Well-tuned simple nets excel on tabular datasets. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 23928–23941, 2021.
- Daniel Kahneman and Gary Klein. Conditions for intuitive expertise: a failure to disagree. *American psychologist*, 64(6):515, 2009.
- Florian Karl, Tobias Pielok, Julia Moosbauer, Florian Pfisterer, Stefan Coors, Martin Binder, Lennart Schneider, Janek Thomas, Jakob Richter, Michel Lang, et al. Multi-objective hyperparameter optimization—an overview. *arXiv preprint arXiv:2206.07438*, 2022.

- Thanh Tung Khuat, David Jacob Kedziora, and Bogdan Gabrys. The roles and modes of human interactions with automated machine learning systems. *arXiv preprint arXiv:2205.04139*, 2022.
- Doris Jung-Lin Lee and Stephen Macke. A human-in-the-loop perspective on automl: Milestones and the road ahead. *IEEE Data Engineering Bulletin*, 2020.
- Cheng Li, Santu Rana, Sunil Gupta, Vu Nguyen, Svetha Venkatesh, Alessandra Sutti, David Rubin, Teo Slezak, Murray Height, Mazher Mohammed, and Ian Gibson. Accelerating experimental design by incorporating experimenter hunches. In *IEEE International Conference on Data Mining*, pages 257–266. IEEE, 2018a.
- Cheng Li, Sunil Gupta, Santu Rana, Vu Nguyen, Antonio Robles-Kelly, and Svetha Venkatesh. Incorporating expert prior knowledge into experimental design via posterior sampling. *arXiv preprint arXiv:2002.11256*, 2020.
- Lisha Li, Kevin G. Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. Hyperband: A novel bandit-based approach to hyperparameter optimization. *Journal of Machine Learning Research*, 18:1–52, 2018b.
- Marius Lindauer, Katharina Eggenberger, Matthias Feurer, André Biedenkapp, Difan Deng, Carolin Benjamins, René Sass, and Frank Hutter. SMAC3: A versatile bayesian optimization package for hyperparameter optimization. *Journal of Machine Learning Research (JMLR)*, 23(54):1–9, 2022.
- K. Louise Barriball and Alison While. Collecting data using a semi-structured interview: a discussion paper. *Journal of Advanced Nursing*, 19(2):328–335, 1994.
- Neeratoy Mallik, Carl Hvarfner, Danny Stoll, Maciej Janowski, Eddie Bergman, Marius Lindauer, Luigi Nardi, and Frank Hutter. Priorband: Hyperband + human expert knowledge. In *Sixth Workshop on Meta-Learning at the Conference on Neural Information Processing Systems*, 2022.
- Michele J. McIntosh and Janice M. Morse. Situating and constructing diversity in semi-structured interviews. *Global Qualitative Nursing Research*, 2:2333393615597674, 2015.
- Scott Mayer McKinney, Marcin Sieniek, Varun Godbole, Jonathan Godwin, Natasha Antropova, Hutan Ashrafian, Trevor Back, Mary Chesus, Greg S. Corrado, Ara Darzi, Mozziyar Etemadi, Florencia Garcia-Vicente, Fiona J. Gilbert, Mark Halling-Brown, Demis Hassabis, Sunny Jansen, Alan Karthikesalingam, Christopher J. Kelly, Dominic King, Joseph R. Ledsam, David Melnick, Hormuz Mostofi, Lily Peng, Joshua Jay Reichler, Bernardino Romera-Paredes, Richard Sidebottom, Mustafa Suleyman, Daniel Tse, Kenneth C. Young, Jeffrey De Fauw, and Shravya Shetty. International evaluation of an ai system for breast cancer screening. *Nature*, 577(7788):89–94, 2020.
- Gábor Melis, Chris Dyer, and Phil Blunsom. On the state of the art of evaluation in neural language models. In *International Conference on Learning Representations*, 2018.
- Douglas C. Montgomery. *Design and analysis of experiments*. Wiley, 2017.

- Julia Moosbauer, Julia Herbinger, Giuseppe Casalicchio, Marius Lindauer, and Bernd Bischl. Explaining hyperparameter optimization via partial dependence plots. In *Advances in Neural Information Processing Systems*, pages 2280–2291, 2021.
- Kevin Musgrave, Serge Belongie, and Ser-Nam Lim. A metric learning reality check. In *European Conference on Computer Vision*, pages 681–699. Springer, 2020.
- Randal S Olson and Jason H Moore. Tpot: A tree-based pipeline optimization tool for automating machine learning. *Automated Machine Learning*, page 151, 2019.
- Michael Polanyi and Amartya Sen. *The tacit dimension*. University of Chicago press, 2009.
- Anil Ramachandran, Sunil Gupta, Santu Rana, Cheng Li, and Svetha Venkatesh. Incorporating expert prior in bayesian optimisation via space warping. *Knowledge-Based Systems*, 195:105663, 2020.
- Jan N. Van Rijn and Frank Hutter. Hyperparameter importance across datasets. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 2367–2376. Association for Computing Machinery, 2018.
- René Sass, Eddie Bergman, André Biedenkapp, Frank Hutter, and Marius Lindauer. Deepcave: An interactive analysis tool for automated machine learning. In *ICML Workshop on Adaptive Experimental Design and Active Learning in the Real World (ReALML)*, 2022. Workshop on Adaptive Experimental Design and Active Learning in the Real World (ReALML@ICML’22).
- Sarah Segel, Helena Graf, Alexander Tornede, Bernd Bischl, and Marius Lindauer. Symbolic explanations for hyperparameter optimization. In *AutoML Conference 2023*. PMLR, 2023.
- Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P. Adams, and Nando de Freitas. Taking the human out of the loop: A review of bayesian optimization. *Proceedings of the IEEE*, 104(1):148–175, 2016.
- Jasper Snoek, Hugo Larochelle, and Ryan P. Adams. Practical bayesian optimization of machine learning algorithms. *Advances in neural information processing systems*, 25, 2012.
- Jasper Snoek, Kevin Swersky, Richard Zemel, and Ryan P. Adams. Input warping for bayesian optimization of non-stationary functions. In *Proceedings of the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pages 1674–1682, 2014.
- Arthur Souza, Luigi Nardi, Leonardo Oliveira, Kunle Olukotun, Marius Lindauer, and Frank Hutter. Bayesian optimization with a prior for the optimum. In *Machine Learning and Knowledge Discovery in Databases. Research Track: European Conference, ECML PKDD 2021, Bilbao, Spain, September 13–17, 2021, Proceedings, Part III 21*, volume 12977 of *Lecture Notes in Computer Science*, pages 265–296. Springer, 2021.
- Kevin Swersky, Jasper Snoek, and Ryan Prescott Adams. Freeze-thaw bayesian optimization. *arXiv preprint arXiv:1406.3896*, 2014.

- Ryan Turner, David Eriksson, Michael McCourt, Juha Kiili, Eero Laaksonen, Zhen Xu, and Isabelle Guyon. Bayesian optimization is superior to random search for machine learning hyperparameter tuning: Analysis of the black-box optimization challenge 2020. In Hugo Jair Escalante and Katja Hofmann, editors, *Proceedings of the NeurIPS 2020 Competition and Demonstration Track*, volume 133, pages 3–26. PMLR, 2021.
- Amos Tversky and Daniel Kahneman. Judgment under uncertainty: Heuristics and biases: Biases in judgments reveal some heuristics of thinking under uncertainty. *Science*, 185(4157):1124–1131, 1974.
- Joaquin Vanschoren. Meta-learning: A survey. In Frank Hutter, Lars Kotthoff, and Joaquin Vanschoren, editors, *Automatic Machine Learning: Methods, Systems, Challenges*, pages 35–61. Springer, 2019.
- Chi Wang, Qingyun Wu, Markus Weimer, and Erkang Zhu. Flaml: A fast and lightweight automl library. In *Machine Learning and Systems*, pages 434–447, 2021a.
- Dakuo Wang, Justin D. Weisz, Michael Muller, Parikshit Ram, Werner Geyer, Casey Dugan, Yla Tausczik, Horst Samulowitz, and Alexander Gray. Human-ai collaboration in data science: Exploring data scientists’ perceptions of automated ai. *Proc. ACM Hum.-Comput. Interact.*, 3(CSCW), nov 2019.
- Dakuo Wang, Josh Andres, Justin D. Weisz, Erick Oduor, and Casey Dugan. Autods: Towards human-centered automation of data science. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, CHI ’21, New York, NY, USA, 2021b. Association for Computing Machinery.
- Dakuo Wang, Q. Vera Liao, Yunfeng Zhang, Udayan Khurana, Horst Samulowitz, Soya Park, Michael Muller, and Lisa Amini. How much automation does a data scientist want? *arXiv preprint arXiv: 2101.03970*, pages 1–31, 2021c.
- Baohe Zhang, Raghu Rajan, Luis Pineda, Nathan Lambert, André Biedenkapp, Kurtland Chua, Frank Hutter, and Roberto Calandra. On the importance of hyperparameter optimization for model-based reinforcement learning. In *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, pages 4015–4023. PMLR, 2021.
- Marc-André Zöller and Marco F. Huber. Benchmark and survey of automated machine learning frameworks. *Journal of artificial intelligence research*, 70:409–472, 2021.
- Marc-André Zöller, Waldemar Titov, Thomas Schlegel, and Marco F. Huber. Xautoml: A visual analytics tool for understanding and validating automated machine learning. *arXiv preprint arXiv: 2202.11954*, 2022.