HYPERPARAMETER OPTIMIZATION VIA INTERACTING WITH PROBABILISTIC CIRCUITS

Anonymous authors

Paper under double-blind review

ABSTRACT

Despite the growing interest in designing truly interactive hyperparameter optimization (HPO) methods, to date, only a few allow to include human feedback. However, these methods add friction to the interactive process, rigidly requiring to fully specify the user input as prior distribution ex ante and often imposing additional constraints on the optimization framework. This hinders the flexible incorporation of expertise and valuable knowledge of domain experts, which might provide partial feedback at any time during optimization. To overcome these limitations, we introduce a novel Bayesian optimization approach leveraging tractable probabilistic models named probabilistic circuits (PCs) as surrogate model. PCs encode a tractable joint distribution over the hybrid hyperparameter space and evaluation scores, and enable exact conditional inference and sampling, allowing users to provide valuable insights interactively and generate configurations adhering to their feedback. We demonstrate the benefits of the resulting interactive HPO through an extensive empirical evaluation of diverse benchmarks, including the challenging setting of neural architecture search.

025 026 027

004

010 011

012

013

014

015

016

017

018

019

021

023

1 INTRODUCTION

028 Hyperparameters crucially influence the performance of machine learning (ML) algorithms and must 029 be set carefully to fully unleash the algorithm's potential (Bergstra & Bengio, 2012; Hutter et al., 2013; Probst et al., 2019). Manually finding good hyperparameters is a tedious and costly task. Thus, various 031 approaches have been proposed to automatize this process which is referred to as *hyperparameter* optimization (HPO) (Bischl et al., 2023). Generally, HPO is framed as optimizing an expensive black-033 box function since the true functional form of the objective is commonly unknown, and the evaluation 034 of hyperparameter configurations is costly, as it requires training ML models several times. Given the rise of deep learning, there is also a growing interest in optimizing the hyperparameters defining the 035 architecture of neural models, commonly referred to as neural architecture search (NAS). The main 036 goal in NAS and HPO is to efficiently traverse the search space to find good configurations quickly 037 and avoid unpromising regions. Bayesian optimization (BO) methods have proven to be sample efficient and converge on good configurations quickly. They exploit knowledge about previously evaluated configurations by learning a surrogate model at each iteration to approximate the optimized 040 unknown objective function (Hutter et al., 2011; Falkner et al., 2018). A selection policy employs the 041 learned surrogate to determine the next configuration to be evaluated (Garnett, 2023; Wang et al., 042 2022). Selection policies aim to balance the exploitation of promising regions of the search space 043 with the exploration of undiscovered ones. Prominent policies optimize an acquisition function 044 assessing the utility of each configuration in terms of the mentioned trade-off (Hutter et al., 2011; Hvarfner et al., 2022) or employ sampling strategies that frame the selection of the next configuration as a multi-armed bandit problem and maximize a reward provided by the surrogate (Shahriari et al., 046 2016; Wang et al., 2022). 047

Although the recent advancements in HPO and NAS could facilitate the design and optimization of ML models for non-experts, in most of the cases, hyperparameters are still tuned manually (Bouthillier & Varoquaux, 2020) and the majority of cutting-edge neural architectures, e.g., transformers (Vaswani et al., 2017), are derived by hand. Given that many ML practitioners perform hyperparameter tuning purely based on their knowledge, experience, and intuition, integrating this valuable knowledge to guide HPO algorithms *during* optimization is of high value since it can substantially foster the search and mitigate its cost. For example, in Fig. 1 (Left), three hyperparameters of the



Figure 1: Interactive Bayesian Hyperparameter Optimization. (Left) We devise interactive Bayesian HPO by employing PCs as surrogate models encoding a joint distribution over hyperparameters and evaluation scores. They allow users to directly condition the surrogate on their beliefs while new candidates are generated via tractable sampling. Thus, user knowledge (priors or point-wise values) is reflected accurately. (Right) Accurately reflecting user beliefs is crucial for interactive HPO to fully leverage user knowledge. IBO-HPC (our method) precisely reflects the user prior provided 072 over the hyperparameter R (image resolution), while πBO and BOPrO fail to do so. 073

074 JAHS benchmark (Bansal et al., 2022) are optimized (depth multiplier N, width multiplier W, and 075 resolution R of a deep CNN). During an HPO run, a user might realize that values around N = 16076 and W = 3 yield high-performing models. Hence, a user can guide the HPO algorithm with the 077 obtained knowledge (here, N = 16 and W = 3) without restarting the optimization from scratch. This can considerably increase the convergence speed and quality of the final solution by focusing the optimization on remaining hyperparameters (here R, details in App. A). Recent works by Souza 079 et al. (2021) and Hvarfner et al. (2022) allow users to infuse knowledge into a BO framework via user-defined prior distributions. To integrate user feedback, both approaches employ weighting 081 schemes that alter the behavior of the acquisition function to follow the user-defined priors. Given a configuration, Souza et al. (2021) weight the surrogate's prediction with the prior, while, Hvarfner 083 et al. (2022) directly reshape the acquisition function by weighting it with the prior. Although these 084 approaches are valid and principled ways to guide an HPO task, their weighting schemes limit the set 085 of compatible acquisition functions and, therefore, the selection policies. Furthermore, they assume user knowledge to be available only *ex ante*, i.e., before the optimization, hindering users to adjust 087 their beliefs flexibly, at *any* iteration during the optimization. Moreover, these approaches might 088 not reflect user knowledge precisely as defined in the prior due to the non-linear integration of the priors in the acquisition function. For example, in Fig. 1 (Right), the configurations selected by 089 both BOPrO (Souza et al., 2021) and π BO (Hvarfner et al., 2022) during the first 20 iterations of 090 optimization remarkably deviate from the given user prior. 091

092 To overcome the above limitations and to integrate user feedback in HPO and NAS more flexibly. making the optimization truly interactive, we introduce INTERACTIVE BAYESIAN OPTIMIZATION 094 VIA HYPERPARAMETER PROBABILISTIC CIRCUITS (IBO-HPC).¹ This novel BO method provides an elegant and flexible mechanism to incorporate user knowledge at any time during optimization 095 without relying on weighting schemes reshaping an acquisition function. Instead, we derive a selection 096 policy that accurately reflects user beliefs by allowing users to directly *condition* the surrogate on those beliefs. Conditions can be specific values or distributions representing users' uncertainties over 098 values of a subset of hyperparameters. Besides this natural incorporation of user knowledge, we introduce a novel, purely data-driven selection policy which suggests new configuration candidates 100 leveraging conditional sampling and avoiding an additional inner loop optimization of an acquisition 101 function. We achieve these benefits by employing as a surrogate a fairly recent tractable probabilistic 102 model called probabilistic circuits (PCs) (Choi et al., 2020). In contrast with other common choices 103 in HPO and NAS such as Gaussian processes (GPs) (Rasmussen & Williams, 2006) and random 104 forests (RFs) (Breiman, 2001), PCs compactly model a joint distribution. PCs allow the tractable

105

068

069

¹⁰⁶ ¹We make our code available at https://anonymous.4open.science/r/jahs_sand-40B3/ 107 and provide all logs and data at https://hessenbox.tu-darmstadt.de/getlink/ fiL7ifX7hq1TfA7MnxndkELS/logs_v2.zip.

and exact computation of arbitrary marginal and conditional distributions, and provide tractable
 (conditional) sampling, rendering them a natural choice for interactive HPO. To ensure robustness
 against potentially misleading user input, we devise a decay mechanism to decrease the influence of
 user knowledge over time. In the following, we refer to these PCs that encode joint distributions over
 hyperparameters and evaluations as *hyperparameter probabilistic circuits* (HPC).

We make the following contributions: (1) We introduce a novel HPO method named IBO-HPC that enables direct incorporation of user knowledge into the selection policy at any time s.t. the user knowledge is precisely reflected in the selection policy as specified by the user. (2) We formally define a notion of interactive policy in HPO, and show that IBO-HPC conforms to this notion and is guaranteed to reflect user knowledge as provided in the optimization process. (3) We provide an extensive empirical evaluation of IBO-HPC showing that it is competitive with strong HPO and NAS baselines without user interaction and outperforms them when leveraging user knowledge.

120 121

2 RELATED WORK

122 123

Bayesian optimization (BO) is an approach to optimize an unknown and expensive black-box function and thus is a natural choice for tackling HPO tasks. BO methods sequentially update, based on observations, a surrogate model that approximates the unknown objective function and captures dependencies between configurations from the search space and the evaluation function's values. In a selection policy, the surrogate is used to choose only promising configurations evaluated next (Hutter et al., 2011; Snoek et al., 2012; Mockus, 1975; Shahriari et al., 2016). Common choices for surrogate models are Gaussian processes (GPs) (Rasmussen & Williams, 2006) or random forests (RFs) (Breiman, 2001).

131 Due to the rise of deep learning, neural architecture search (NAS) has become increasingly relevant. 132 Besides BO, different approaches based on local search (Den Ottelander et al., 2021), reinforcement 133 learning (Pham et al., 2018; Zoph & Le, 2017), and gradient descent (Liu et al., 2019) have been 134 proposed. See White et al. (2023) for a comprehensive survey on NAS. To foster the development 135 and fair comparison of HPO and NAS algorithms, different benchmarks have been proposed. These 136 benchmarks define a search space over architectures and hyperparameters and train all candidates 137 to provide several quantities such as validation and test accuracy. Ying et al. (2019) and Dong & Yang (2020) introduce such benchmarks for NAS (NAS-Bench-101/201) while Bansal et al. (2022) 138 introduces a benchmark for joint optimization of hyperparameters and neural architectures (JAHS). 139

Leveraging previous HPO runs to increase the efficiency of HPO on related problems has been considered in several works under the umbrella of hyperparameter transfer learning (HTL). Prominent approaches perform HTL by projecting objective responses of all runs to a common response surface, via learning a common representation of HPO tasks or by pruning the search space based on previous tasks (Yogatama & Mann, 2014; Wistuba et al., 2015; Perrone et al., 2018; Vanschoren, 2018; Salinas et al., 2020; Horváth et al., 2021).

146 With the goal of making HPO more trustworthy, several works have introduced methods to provide 147 explanations in HPO (Hutter et al., 2014; Moosbauer et al., 2021; Watanabe et al., 2023; Segel et al., 2023), while involving feedback during search has received little attention. Hvarfner et al. (2022) 148 allow users to provide prior beliefs via prior distributions over the search space. The provided user 149 prior is used to reshape the acquisition function when selecting new configurations, thus, favoring 150 configurations that received a high likelihood in the prior. Mallik et al. (2023) propose a similar 151 mechanism to incorporate user knowledge in multi-fidelity optimization. However, such approaches 152 assume an *ex ante* full specification of the priors, often with additional constraints such as requiring 153 invertible priors (Ramachandran et al., 2020) or a specific acquisition function (Souza et al., 2021). 154 As illustrated in Sec. 1, the mentioned approaches have several drawbacks in reflecting given user 155 knowledge well when selecting new configurations. We believe that truly interactive HPO methods 156 should not assume ex ante specification of user knowledge and should reflect user knowledge properly.

157 158

3 INTERACTIVE HYPERPARAMETER OPTIMIZATION

159 160

We now present IBO-HPC, which provides a flexible way to interact with the optimiz

161 We now present IBO-HPC, which provides a flexible way to interact with the optimization process via probabilistic queries. This way, user knowledge is decoupled from shaping acquisition functions

while reflecting given knowledge in the selection policy as specified by the user. In this section, we
briefly revise BO as a general framework for HPO before introducing our definition of a *feedback adhering interactive policy*. Next, we introduce HPCs and present the concrete instantiation of our
interactive optimization method. For completeness, we start with a general formal definition of the
HPO problem (Kohavi & John, 1995; Hutter et al., 2019).

167 Definition 1 (Hyperparameter optimization). Given hyperparameters $\mathcal{H} = \{H_1, \ldots, H_n\}$ with **168** associated domains $\mathbf{H}_1, \ldots, \mathbf{H}_n$, and a set of problem instances \mathcal{X} , we define a search space **169** $\Theta = \mathbf{H}_1 \times \cdots \times \mathbf{H}_n$. For a given problem instance $\mathbf{x} \in \mathcal{X}$ and evaluation function $f : \Theta \times \mathcal{X} \to \mathbb{R}$, hyperparameter optimization aims to solve $\theta^* = \arg \min_{\theta \in \Theta} f(\theta; \mathbf{x})$.

171 172

193

3.1 INTERACTIVITY IN BAYESIAN HPO

173 **Bayesian Optimization.** BO aims to optimize a black-box objective function $f: \Theta \to \mathbb{R}$ which is 174 costly to evaluate, i.e., to find the input $\hat{\theta}^* \in \arg \min_{\theta \in \Theta} f(\theta)$ (Shahriari et al., 2016). BO typically 175 tackles such problem in sequential steps, leveraging two key ingredients: a probabilistic surrogate 176 model and a selection policy determining the next θ' to be evaluated. Given a set \mathcal{D}_n of observations 177 that correspond to the configurations with associated evaluations $(\theta_j, f(\theta_j))_{j=1...n}$, the surrogate 178 $s \in S$ aims to induce a distribution over functions $p(f|\mathcal{D}_n)$. The selection policy uses s to select the 179 next $\theta' \in \Theta$ s.t. it achieves a good exploration-exploitation trade-off. Prominent selection policies 180 optimize an acquisition function $a: \Theta \times S \to \mathbb{R}$, such as expected improvement (Jones et al., 1998), 181 that estimates the utility of an evaluation at an arbitrary point $\theta \in \Theta$ under a surrogate $s \in S$. The 182 configuration θ' , evaluated next, is selected by optimizing a. A popular alternative to obtain the next configuration θ' is Thompson sampling (Wang et al., 2022). The obtained tuple (θ' , $f(\theta')$) is added 183 to \mathcal{D}_n and used to update the surrogate model for the next iteration. This process is repeated until 184 convergence. 185

An interactive BO method should be capable of incorporating, at any time, the knowledge provided
 by users; also, the selection policy should reflect the provided knowledge as specified by the user.
 Consequently, we formalize the concept of an interactive selection policy, or interactive policy for
 short, that adheres to these requirements.

Definition 2 (Interactive Policy). An interactive policy is a function $p : S \times \Theta \times K \to \mathcal{P}(\Theta)$ mapping from the set of surrogates S and search space Θ to the set of all distributions $\mathcal{P}(\Theta)$ over the search space Θ while accepting user knowledge $\mathcal{K} \in \mathcal{K}$.

Note that the collection of user knowledge \mathcal{K} and the set of surrogates \mathcal{S} are left unspecified to keep Def. 2 rather general. Since Def. 2 can be trivially achieved, e.g. by ignoring user knowledge, we introduce *feedback adhering interactive policies* that guarantee that (i) user knowledge affects the policy's outcome and (ii) user knowledge is reflected in the interactive policy as specified.

Definition 3 (Feedback Adhering Interactive Policy). Given user knowledge $\mathcal{K} \in \mathcal{K}$ and surrogate s_t $\in S$ at iteration t, an interactive policy p is called efficacious if $p(\Theta, s_t, \mathcal{K}) \neq p(\Theta, s_t, \emptyset)$ where \emptyset indicates that p is applied without user knowledge. If further \mathcal{K} is provided as a distribution $q(\hat{\mathcal{H}})$ over $\hat{\mathcal{H}} \subset \mathcal{H}$, we call p feedback adhering if it is efficacious and $\int_{\mathcal{H} \setminus \hat{\mathcal{H}}} p(\Theta, s_t, \mathcal{K}) = q(\hat{\mathcal{H}})$

holds, i.e., the distribution over $\hat{\mathcal{H}}$ induced by the selection policy equals the prior $q(\hat{\mathcal{H}})$ in the next iteration.

In Def. 3, the first condition ensures that the user knowledge has an effect on the sampling policy. The second condition ensures that in the first iteration, after that a user provides a distribution over a subset of hyperparameters, the values sampled for the specified hyperparameters follow exactly the distribution q given by the user. Note that user knowledge could also be misleading; thus, Def. 3 does not require user knowledge to have exclusively positive effects. Equipped with Def. 2 and 3, we now introduce IBO-HPC that adheres to both definitions.

- 3.2 INTERACTIVE BAYESIAN OPTIMIZATION WITH HYPERPARAMETER PROBABILISTIC CIRCUITS
- In this section, we introduce an interactive Bayesian optimization method that fulfills Def. 3. It
 employs hyperparameter probabilistic circuits (HPCs) as a surrogate model and a selection policy
 leveraging the flexible inference and sampling of HPCs, avoiding an additional inner loop optimization

216 of an acquisition function. Before delving into our method, we first provide preliminaries on 217 probabilistic circuits.

4

6

7

8

10

11

12

13

15

16

17

18

218 Hyperparameter Probabilistic Circuits 219 (HPCs). Motivated by the lack of truly inter-220 active Bayesian HPO methods, we seek a policy 221 that enables flexible interactions with the opti-222 mization procedure by providing an arbitrary amount of knowledge about hyperparameters 224 at any time *during* the optimization while 225 reflecting user beliefs as specified. Probabilistic 226 circuits (Choi et al., 2020) are computation graphs that compactly represent multivariate 227 distributions. PCs can answer a wide range of 228 probabilistic queries in a tractable fashion and 229 (conditionally) generate new samples. These 230 features make them a good candidate for our 231 purpose of building a policy adhering to Def. 3. 232

More formally, a PC is a computational graph 233 encoding a distribution over a set of random vari-234 ables **X**. It is defined as a tuple (\mathcal{G}, ϕ) where 235 $\mathcal{G} = (V, E)$ is a rooted, directed acyclic graph 236 and ϕ : $V \to 2^{\mathbf{X}}$ is the scope function as-237 signing a subset of random variables to each 238 node in \mathcal{G} . For each internal node N of \mathcal{G} , the 239 scope is defined as the union of scopes of its 240 children, i.e. $\phi(N) = \bigcup_{N' \in ch(N)} \phi(N')$. Each 241 leaf node L computes a distribution/density over 242 its scope $\phi(L)$. All internal nodes of \mathcal{G} are ei-243 ther a sum node S or a product node P where 14 each sum node computes a convex combination 244 of its children, i.e., $S = \sum_{N \in ch(S)} w_{S,N}N$, and 245 each product computes a product of its children, 246 i.e., $P = \prod_{N \in ch(P)} N$. We assume *smooth* and 247 decomposable PCs (see App. C for details); 248 19 thus, our method can exploit tractable inference, 20 249 sampling, and conditioning of PCs. For a more 21 250 detailed description of PCs, refer to App. C; for 251 an overview, see Fig. 1 (Left). 252

Algorithm 1: Interactive BO with HPCs (IBO-HPC). Our interactive BO method allows for flexible incorporation of user knowledge at any iteration via conditional sampling enabling true interaction with users.

Data: Search space Θ over $\mathcal{H} = \{H_1, \ldots, H_n\},$ problem instance $\mathbf{x} \in \mathcal{X}$, prior distribution $p(\mathbf{\Theta})$, objective $f: \boldsymbol{\Theta} \times \boldsymbol{\mathcal{X}} \to \mathbb{R}$, user prior $q(\mathcal{H})$ is optional and can be provided at any time, decay γ 1 $\mathcal{D} = \emptyset;$ **2** for $i \in \{1, ..., J\}$ do $\boldsymbol{\theta} \sim p(\boldsymbol{\Theta});$ $\mathcal{D} = \mathcal{D} \cup \{(\boldsymbol{\theta}, f(\boldsymbol{\theta}, \mathbf{x}))\};\$ 5 while not converged do every L-th iteration, fit HPC s on \mathcal{D} ; $f^* = \max_f \mathcal{D};$ $b \sim \text{Ber}(\rho)$: if prior $q(\mathcal{H})$ given and b = 1 then sample N conditions $\mathbf{h} \sim q(\mathcal{H})$; $\mathbf{C} = \emptyset;$ for condition h_i in h do sample *B* configurations
$$\begin{split} \boldsymbol{\theta}_{1,...,B}^{\prime} &\sim s(\boldsymbol{\mathcal{H}} \setminus \hat{\boldsymbol{\mathcal{H}}} | \hat{\boldsymbol{\mathcal{H}}}, f^{*}); \\ \boldsymbol{\theta}_{i}^{*} &= \arg \max_{\boldsymbol{\theta}^{\prime} \in \boldsymbol{\theta}_{1,...,B}^{\prime}} s(\boldsymbol{\theta}^{\prime} | f^{*}); \end{split}$$
 $\mathbf{C} = \mathbf{C} \cup \boldsymbol{\theta}_i^*;$ $\theta^* \sim \mathcal{U}(\mathbf{C});$ else $\boldsymbol{\theta}^* \sim s(\boldsymbol{\mathcal{H}}|f^*);$ $\mathcal{D} = \mathcal{D} \cup (\theta', f(\theta', \mathbf{x}));$ $\rho = \gamma \cdot \rho;$ present evaluations \mathcal{D} ;

253 We jointly model the hyperparameters and evaluation scores with PCs; thus, we refer to these surrogates as hyperparameter probabilistic circuits (HPC). Given the hybrid (discrete and continuous) 254 nature of hyperparameter search spaces, in this work, we focus on a type of PCs tailored for hybrid 255 domains named mixed sum-product networks (MSPNs). An MSPN is a decomposable and smooth 256 PC with piecewise polynomial leaves. These properties types allow MSPNs to represent valid 257 distributions over hybrid domains (i.e. discrete and continuous variables) (Molina et al., 2018). 258

Method. We now describe our method shown in Algorithm 1. Since our surrogate is a density 259 estimator, we start off by sampling J hyperparameter configurations from a prior distribution p, e.g., 260 a uniform distribution, and evaluate them by querying the objective function f (Line 1-5). The 261 function f yields a performance score of a model trained to solve a given problem instance x with 262 the sampled configuration θ . After evaluating each sampled θ we obtain a set \mathcal{D} of pairs $(\theta, f_{\theta}(\mathbf{x}))$ 263 and fit a HPC s estimating the joint distribution $p(\mathcal{H}, F)$, where \mathcal{H} is the set of hyperparameters and 264 F is a random variable representing the evaluation score (Line 7). IBO-HPC proceeds by selecting a 265 configuration θ that gets evaluated next, i.e., our *feedback adhering interactive policy* is applied. Our 266 policy exploits the flexible and exact inference of HPCs to derive arbitrary conditional distributions 267 according to the partial evidence at hand (Peharz et al., 2015). We target the configurations that are 268 likely to achieve a better evaluation score. Thus, a posterior distribution over the hyperparameter space is derived by conditioning on the best score $f^* = \max_f \mathcal{D}$ observed so far alongside with

276 277

278

279

280

281 282

317

(optional) user knowledge \mathcal{K} . For now, \mathcal{K} is assumed to be given in the form of conditions such as $\hat{\mathcal{H}} = \hat{\mathbf{h}}$ where $\hat{\mathcal{H}} \subset \mathcal{H}$ is a subset of hyperparameters being set to $\hat{\mathbf{h}}$. With Bayes rule and tractable marginal inference and sampling of HPCs, we obtain the conditional distribution and use it to sample a new configuration from promising regions in the search space.

$$p(\mathcal{H} \setminus \hat{\mathcal{H}} | \hat{\mathcal{H}}, F = f^*) = s(\mathcal{H} \setminus \hat{\mathcal{H}} | \hat{\mathcal{H}}, F = f^*)$$

$$\boldsymbol{\theta} \sim p(\mathcal{H} \setminus \hat{\mathcal{H}} | \hat{\mathcal{H}}, F = f^*)$$
(1)

Since users might be uncertain about hyperparameter values, defining a prior $q(\hat{\mathcal{H}})$ over $\hat{\mathcal{H}}$ might be more reasonable than setting a fixed value for certain hyperparameters. The prior $q(\hat{\mathcal{H}})$ is interpreted as a distribution over conditions of the form $\hat{\mathcal{H}} = \hat{\mathbf{h}}$ where $\hat{\mathbf{h}} \sim q(\hat{\mathcal{H}})$. This induces a different weighted version of the distribution given in Eq. 1.

$$p(\mathcal{H} \setminus \hat{\mathcal{H}} | \hat{\mathcal{H}}, F = f^*) \cdot q(\hat{\mathcal{H}}) = s(\mathcal{H} \setminus \hat{\mathcal{H}} | \hat{\mathcal{H}}, F = f^*) \cdot q(\hat{\mathcal{H}})$$
(2)

Since user intuition can be wrong, we allow IBO-HPC to recover from sub-optimal user knowledge by deciding whether or not to use the provided \mathcal{K} based on a Bernoulli distribution with success probability ρ . To ensure that IBO-HPC gradually recovers when misleading \mathcal{K} is provided, we decrease the likelihood of using \mathcal{K} in each iteration after \mathcal{K} was supplied via a decay factor γ . When user knowledge is provided at iteration T, the distribution over configurations after T + t iterations reads: $\gamma^t \rho \cdot s(\mathcal{H} \setminus \hat{\mathcal{H}} | \hat{\mathcal{H}} | F - f^*) \cdot g(\hat{\mathcal{H}}) + (1 - \gamma^t \rho) \cdot s(\mathcal{H} | F - f^*)$ (3)

$$\gamma^t \rho \cdot s(\mathcal{H} \setminus \hat{\mathcal{H}} | \hat{\mathcal{H}}, F = f^*) \cdot q(\hat{\mathcal{H}}) + (1 - \gamma^t \rho) \cdot s(\mathcal{H} | F = f^*)$$
(3)

290 Note that fusing the distribution in Eq. 2 with the HPC to allow exact inference and conditioning 291 is non-trivial since the prior q is defined over an arbitrary subset and no further assumptions about 292 q are made. Thus we approximate Eq. 2 by sampling N times from $q(\hat{\mathcal{H}})$ and use Eq. 1 to obtain 293 N conditional distributions respecting the user prior $q(\hat{\mathcal{H}})$. To select a promising configuration for 294 evaluation from the approximated distribution in Eq. 2 while still achieving exploration, we sample 295 B configurations from all N conditionals. For each conditional, the configuration maximizing the 296 likelihood $s(\mathcal{H}|F = f^*)$ is selected to reduce the candidate set to configurations likely to achieve 297 a high evaluation score. This leaves us with N configurations from which we sample uniformly to select the configuration evaluated next (Line 10-16). We found that setting B = 1 works surprisingly 298 well. A discussion about the quality of our approximation is given in App. B.4. The surrogate is 299 kept fixed for L optimization rounds before retraining it. This fosters exploration by leveraging 300 uncertainty encoded in the (conditional) distribution of the surrogate. An iteration is concluded by 301 updating the set of evaluations \mathcal{D} that can be presented to the users (Line 20-21). The algorithm runs 302 until convergence or another condition for termination, e.g., a time budget limit is encountered. 303

Remark 1. Although different, our sampling policy shares similarities with Thompson Sampling (TS): TS samples function values from the posterior and selects the next configuration based on the obtained maximum of the function. Instead of sampling the function value from the posterior, we use the maximum obtained so far and use it to sample the next configuration.

Theoretical Properties. After presenting IBO-HPC as an instance of an IBO algorithm, we now show that the policy applied to select the next configuration is a feedback adhering interactive policy according to Def. 3. Since modeling dependencies among hyperparameters is non-trivial for users, we only require users to provide a product distribution as prior knowledge.

Proposition 1 (IBO-HPC Policy is feedback adhering interactive). Given a search space Θ over hyperparameters \mathcal{H} , an HPC $s \in \mathcal{S}$, user knowledge $\mathcal{K} \in \mathcal{K}$ in form of a prior q over $\hat{\mathcal{H}} \subset \mathcal{H}$ s.t. the marginal distribution over $\hat{\mathcal{H}}$ of s conditioned on f^* is different than $q(\hat{\mathcal{H}})$, i.e., $\int_{\mathcal{H}\setminus\hat{\mathcal{H}}} s(\hat{\mathcal{H}}|F =$ $f^*) \neq q(\hat{\mathcal{H}})$, the selection policy of IBO-HPC is feedback adhering interactive. The proof is provided in the App. B.1.

318 Besides being feedback adhering, IBO-HPC is a global optimizer for black-box optimization.

Proposition 2 (IBO-HPC is a global optimizer). *IBO-HPC minimizes simple regret, which is defined as* $r = f(\mathbf{h}) - f(\mathbf{h}^*)$ *for a hyperparameter configuration* $\mathbf{h} \in \Theta$ *and global optimum* \mathbf{h}^* . *A proof is given in B.2.*

To analyze the convergence rate of IBO-HPC, we involve the expected improvement (EI) with a surrogate PC s. This gives us the following proposition.



Figure 2: **IBO-HPC outperforms state of the art.** For 5/5 tasks across three challenging benchmarks, IBO-HPC is competitive with strong baselines when no user knowledge is provided. When beneficial user beliefs (\triangle) are provided, either after 5 iterations (—) or after 15 iterations (—), it outperforms all competitors w.r.t. convergence and solution quality on 4/5 tasks.

Proposition 3 (Convergence of IBO-HPC). Assume a non-noisy differentiable L-Lipschitz continuous function $f : \mathbb{R}^d \to \mathbb{R}$ with global optimum $\mathbf{h}^* \in \mathbb{R}^d$ that is convex within a ball $B_r(\mathbf{h}^*) = {\mathbf{h} \in \mathbb{R}^d : ||\mathbf{h} - \mathbf{h}^*|| < r}$. Further, assume we have given a dataset $\mathcal{D} = {(\mathbf{h}_1, y_1), \ldots, (\mathbf{h}_n, y_n)}$ where all $\mathbf{h}_i \in B_r$ and $y_i = f(\mathbf{h}_i)$ and a decomposable, complete PC s over $\mathcal{H} \cup {F}$ where the support of $\mathcal{H} = B_r$ and the support of $F = \mathbb{R}$. Assume s locally maximizes the likelihood over \mathcal{D} and that all leaves are Gaussians. Then, the convergence rate of IBO-HPC is lower bounded by the expected improvement (EI) in each iteration, that is

$$\sum_{k=1}^{\tau_s} w_i \cdot \Big(\prod_{j=1}^d \operatorname{erf}\Big(\frac{\mathbf{h}_{tj}^* - \boldsymbol{\mu}_{ij}}{\Sigma_{i_{jj}}\sqrt{2}}\Big) - \prod_{j=1}^d \operatorname{erf}\Big(\frac{\mathbf{h}_j^* - \boldsymbol{\mu}_{ij}}{\Sigma_{i_{jj}}\sqrt{2}}\Big) + L\epsilon_i\Big).$$
(4)

Here, τ_s is the number of induced trees of s (see Def. 4 in App. B.3), $\epsilon_i = ||\boldsymbol{\mu}_i + \alpha_i \cdot diag(\Sigma_i) - \mathbf{h}^*||$, each $\boldsymbol{\mu}_i$ is the mean vector of a d-dimensional multivariate Gaussian defined by the *i*-th induced tree, Σ_i is the corresponding correlation matrix and \mathbf{h}_t^* is the best performing configuration until iteration t. A proof is given in App. B.3.

Intuitively, the EI (and thus the convergence rate) is determined by (1) the probability of sampling a configuration in a region bringing us closer to \mathbf{h}^* and (2) how much we expect to move towards the optimum \mathbf{h}^* if (1) occurs. While (1) is lower bounded by considering how close the mixture means are to the best obtained configuration \mathbf{h}_t^* at iteration t, and how far off the mixture means are from the optimum \mathbf{h}^* , (2) is lower bounded for each mixture component by ϵ_i and the Lipschitz constant L.

370 371 372

373

350

351

352

353

354

355

362

363

364 365

366

367

368

369

4 EXPERIMENTAL EVALUATION

We now provide an extensive empirical evaluation of IBO-HPC and aim to answer the following
research questions: (Q1) Can IBO-PC compete with prominent baseline HPO algorithms? (Q2)
How does the performance of IBO-HPC, provided with user knowledge at various points during
optimization, compare to existing approaches incorporating user knowledge ex ante? (Q3) Is IBO-HPC capable of reliably recovering from misleading user interactions?

378 Experimental Setup. We compare IBO-HPC against seven diverse competitors: local search 379 (LS) (White et al., 2020), BO with random forest (RF) surrogate (Head et al.), and SMAC (Hutter 380 et al., 2011) as standard HPO methods. Furthermore, we used random search (RS) (Bergstra & 381 Bengio, 2012) with user priors, BOPrO (Souza et al., 2021), π BO (Hvarfner et al., 2022) and 382 Priorband (Mallik et al., 2023) which allow ex ante incorporation of user knowledge. Since Priorband is a multi-fidelity method, we reserved the number of epochs as the fidelity in each benchmark. For all non-multi-fidelity methods, set it to the highest possible value. For our evaluation, we employ 384 four real-world benchmarks, i.e., NAS-Bench-101 (Ying et al., 2019) and NAS-Bench-201 (Dong & 385 Yang, 2020) as tabular NAS benchmarks, JAHS (Bansal et al., 2022) as a surrogate benchmark for 386 joint architecture and hyperparameter search, as well as HPO-B (Arango et al., 2021) containing a 387 large collection of challenging OpenML tasks. We evaluate IBO-HPC on seven diverse tasks and five 388 different search spaces covering continuous, discrete, and mixed spaces. Our evaluation considers the 389 optimization of neural architectures and tree-based models on image and tabular data. For NAS (and 390 NAS+HPO), we consider CIFAR-10 (JAHS, NAS-Bench-101, NAS-Bench-201), Fashion-MNIST 391 (JAHS), and Colorectal Histology (JAHS). To make the optimization problem more challenging, we 392 extended the search space definition of JAHS (see App. D). This is possible since JAHS is a surrogate 393 benchmark. Additionally, we consider optimizing an XGBoost and random forest classifier on the credit-g dataset, taken from the HPO-B benchmark. As usual, we optimize the validation accuracy as 394 our objective. All algorithms were repeated on 500 different seeds for a maximum of 2000 iterations 395 (100 iterations for HPO-B). We report the mean test error against computational cost and provide 396 standard error to quantify uncertainty. The computational costs are reported as the accumulated 397 wall-clock time of training and evaluation of each sampled configuration, where the benchmarks 398 provide wall-clock times for each configuration. All experiments were run on DGX-A100 machines. 399

User Interactions. For the experiments, beneficial and misleading user interactions have been defined 400 as user priors for each benchmark. To define priors, we randomly sampled 10k configurations and 401 kept the best/worst performing ones, denoted as h^+ and h^- , respectively. To demonstrate that user 402 priors over a few hyperparameters are enough to improve the performance of IBO-HPC significantly, 403 we defined beneficial interactions by selecting a small subset of hyperparameters $\mathcal{H} \subset \mathcal{H}$. Then, 404 we defined a prior over each $H \in \mathcal{H}$ s.t. the probability of sampling the value of H given in \mathbf{h}^+ , 405 denoted by $\mathbf{h}^+[H]$, is 1000 times higher than sampling a different value than $\mathbf{h}^+[H]$. For misleading 406 interaction, \mathcal{H} was chosen to be large to demonstrate that IBO-HPC recovers even if a large amount 407 of misleading information is provided. We then defined priors over $\hat{\mathcal{H}}$ as for beneficial interactions; 408 however, this time the probability to sample $\mathbf{h}^{-}[H]$ is 1000 times higher than for other values for 409 each $H \in \mathcal{H}$. The priors were chosen to be rather strong since, as emphasized in Sec. 1, the stronger 410 the prior, the better πBO and BOPrO reflect user knowledge in their selection policy. Striving for a 411 fair comparison, we opted for such strong priors. Further, we aimed to show that IBO-HPC reliably 412 recovers from receiving large amounts of strongly misleading knowledge. Sometimes, it is easier for 413 users to specify a concrete value for certain hyperparameters instead of defining a distribution. Thus, 414 we also conducted experiments with priors defined as a point mass. See App. D for further details. 415

416 417

418 419

4.1 (Q1) IBO-HPC IS COMPETITIVE IN HPO & NAS

To demonstrate the effectiveness of IBO-HPC, we ran IBO-HPC on all tasks without user interaction. 420 We compared its performance against two strong BO baselines and LS. Fig. 2 shows that the 421 performance of IBO-HPC without user interaction is competitive to or outperforms BO baselines in 422 4/5 tasks across the NAS-101/201 and JAHS benchmarks. We obtained similar results on HPO-B 423 (see App. D.3). These results show that IBO-HPC performs well in complex and realistic settings. 424 Also, it underlines that HPCs accurately capture characteristics of the objective function and that our 425 sampling-based selection policy reliably identifies good configurations. Besides the quality of the 426 final solution, we also observe that IBO-HPC converges at rates similar to those of the baselines. The 427 only exception, where SMAC and LS achieve better results than IBO-HPC, is NAS-101. While LS is 428 known to be a strong baseline for NAS (White et al., 2020; Den Ottelander et al., 2021), we posit that SMAC's more complex selection policy is particularly effective in handling sparse search spaces like 429 NAS-101, which comes with a much sparser representation than NAS-201 and JAHS. To summarize, 430 we state that IBO-HPC is competitive with existing strong BO baselines when no interaction takes 431 place, and answer (Q1) affirmatively.

454

455

456

457 458 459

460



Figure 3: **IBO-HPC recovers from misleading interactions.** IBO-HPC automatically recovers (—) from misleading feedback provided as point values at the 5th iteration of the search (1st \triangle). Also, when providing harmful and beneficial beliefs alternatively (\triangle / \triangle), IBO-HPC (—) catches up with or outperforms π BO (…) and BOPrO (…) in 4/5 cases.

4.2 (Q2, Q3) INTERACTIVE AND RESILIENT HPO & NAS WITH IBO-HPC

We now demonstrate that IBO-HPC successfully handles various kinds of user knowledge (point values and distributions), analyze the benefits of user knowledge w.r.t. convergence speed, and demonstrate IBO-HPC's recovery from misleading beliefs. Therefore, different beneficial and misleading user beliefs about hyperparameters for all benchmarks were defined (details in App. D).

465 **Beneficial Interactions.** Fig. 2 shows a clear positive effect of providing beneficial user beliefs 466 (i.e. distribution or fixing hyperparameters) to IBO-HPC across all tasks. This holds for very early interactions (after 5 iterations; — and —) and later interactions (after 15 iterations; —). Remarkably, 467 we observed a clear benefit in terms of convergence speed and improvement in solution quality, 468 especially for more complex search spaces. Besides outperforming strong HPO baselines incapable 469 of incorporating user knowledge, IBO-HPC is competitive to or outperforms πBO , Priorband and 470 BOPrO, which assume user priors to given ex ante. These results show that IBO-HPC's policy reflects 471 user beliefs accurately and leverages provided knowledge effectively by sampling configurations 472 that are promising according to both, user belief and the surrogate HPC. App. D provides additional 473 experiments on HPO-B showing similar results. 474

Recovery and Multiple Interactions. User beliefs could also be misleading for the optimization 475 process; thus, an interactive HPO algorithm should recover from such misleading interactions and 476 allow users to correct their initial beliefs. We demonstrate that IBO-HPC recovers from misleading 477 user knowledge by deliberately providing IBO-HPC with known sub-optimal values for a large number 478 of hyperparameters to ensure a significant negative effect on the optimization process (see App. D for 479 details). Fig. 3 shows that IBO-HPC (—) recovers similarly well or better as πBO and BOPrO from 480 misleading interactions. In most cases, IBO-HPC catches up with standard HPO competitor methods. 481 This confirms that IBO-HPC's recovery mechanism works reliably and that misleading user beliefs 482 do not deteriorate IBO-HPC's performance in the long run. Again, on NAS-101, IBO-HPC is less 483 effective, which we attribute to the sparse nature of NAS-101 (see Sec. 4.1). Since users might revise their beliefs when no improvement is obtained, we demonstrate that IBO-HPC successfully handles 484 multiple, contradictory interactions. Therefore, we first provided IBO-HPC with the same misleading 485 beliefs as before at an early stage (after 5 iterations), followed by an alternation of beneficial and



511

486



Figure 4: IBO-HPC achieves considerable runtime improvements. (a) IBO-HPC is more efficient than SMAC in 4/5 cases (averaged over 20 runs). With the number of hyperparameters increasing, the gap between IBO-HPC and SMAC in terms of computational efficiency is larger. Runtimes of 2000 iterations normalized between [0, 1] are reported per benchmark (highest obtained runtime for a given benchmark is 1). (b) Beneficial interactions lead to significant speed-ups, from 2 to $10 \times$.

harmful beliefs every 10 iterations. As expected, the misleading interactions decelerate IBO-HPC, 508 and the recovery mechanism is triggered. In contrast, with beneficial interactions, IBO-HPC quickly 509 catches up with the competitors or even outperforms them, confirming that IBO-HPC leverages 510 valuable feedback in critical conditions (see Fig. 3 (—) and App. D).

Speed-up. Both, the runtime of the optimization loop (fitting surrogate and suggesting the next 512 configuration) and convergence speed are crucial for efficient HPO. We, therefore, analyze the average 513 runtime of IBO-HPC's optimization loop and the increase in convergence speed when valuable user 514 knowledge is provided to IBO-HPC as a distribution. In Fig. 4 (a) we compare the runtime of the 515 optimization loop of SMAC and IBO-HPC, averaged over 20 runs. IBO-HPC is considerably faster 516 than SMAC in 4/5 cases, especially in larger search spaces. We attribute this to the efficiency of PCs 517 and of our selection policy (i.e., conditional sampling). In Fig. 4 (b), we ran IBO-HPC without user 518 interaction and obtained the wall-clock time needed for the best evaluation result (denoted as t_w). 519 Then, we ran IBO-HPC with beneficial user knowledge and measured the estimated wall-clock time 520 until IBO-HPC found an equally well or better-performing configuration (denoted as t_i). Fig. 4 (b) reports the relative performance speedup $\frac{t_w}{t_i}$ for all 500 runs. A median speed-up of 2 to 10× with 521 useful user interactions, clearly demonstrates IBO-HPC's increase in convergence speed while saving 522 resources. Since IBO-HPC effectively incorporates various user interactions, leading to remarkable 523 speedups, and provides a reliable recovery mechanism, we answer (Q2) and (Q3) positively. 524

- 525
- 526 527

5 CONCLUSION

- 528 529
- 530

We introduced a novel definition of interactive BO policies and an interactive BO method named 531 IBO-HPC that leverages the flexible inference of probabilistic circuits to flexibly incorporate user beliefs. With no user knowledge, IBO-HPC is competitive with strong baselines and it outperforms 532 competitors when knowledge is available. Also, it reliably recovers from misleading user beliefs and 533 converges significantly faster when provided with valuable user knowledge, thus, saving resources. 534

Limitations & Future Work. Whereas IBO-HPC enables flexible interactive BO, the necessity to 536 retrain the surrogate model at each iteration remains. Thus, prospective directions could explore 537 methods of continual learning (Mundt et al., 2023) to increase the overall efficiency and knowledge reuse over different HPO and NAS settings. Moreover, to model hybrid domains, we relied on 538 HPCs employing piecewise polynomials that might not be sufficient to model complex distributions. Therefore, more sophisticated alternatives could further improve performance.

540 REFERENCES

547

554

558

565

566

569

580

581

582

Sebastian Pineda Arango, Hadi S. Jomaa, Martin Wistuba, and Josif Grabocka. Hpo-b: A large-scale
 reproducible benchmark for black-box hpo based on openml, 2021.

- Archit Bansal, Danny Stoll, Maciej Janowski, Arber Zela, and Frank Hutter. Jahs-bench-201: A
 foundation for research on joint architecture and hyperparameter search. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *Journal of Machine Learning Research*, 13(2), 2012.
- 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. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 13(2), 2023.
- Xavier Bouthillier and Gaël Varoquaux. Survey of machine-learning experimental methods at
 NeurIPS2019 and ICLR2020. Research report, Inria Saclay Ile de France, January 2020.
- ⁵⁷ Leo Breiman. Random forests. *Machine learning*, 45:5–32, 2001.
- 559 YooJung Choi, Antonio Vergari, and Guy Van den Broeck. Probabilistic circuits: A unifying 560 framework for tractable probabilistic models. Technical report, UCLA, 2020.
- Tom Den Ottelander, Arkadiy Dushatskiy, Marco Virgolin, and Peter A. N. Bosman. Local search is a remarkably strong baseline for neural architecture search. In *Evolutionary Multi-Criterion Optimization: 11th International Conference*, 2021.
 - Xuanyi Dong and Yi Yang. Nas-bench-201: Extending the scope of reproducible neural architecture search. In *International Conference on Learning Representations (ICLR)*, 2020.
- Stefan Falkner, Aaron Klein, and Frank Hutter. BOHB: Robust and efficient hyperparameter opti mization at scale. In *International Conference on Machine Learning (ICML)*, 2018.
- 570 Roman Garnett. *Bayesian Optimization*. Cambridge University Press, 2023.
- Robert Gens and Pedro Domingos. Learning the structure of sum-product networks. In *Proceedings* of the 30th International Conference on Machine Learning, volume 28, pp. 873–880. PMLR, 2013.
- Tim Head, MechCoder Gilles, Louppe Iaroslav, Shcherbatyi fcharras, Zé Vinícius, cmmalone, Christopher Schröder, nel215, Nuno Campos, Todd Young, Stefano Cereda, Thomas Fan, Justus Schwabedal, Mikhail Hvass-Labs, Pak SoManyUsernamesTaken, Fred Callaway, Loïc Estève, Lilian Besson, Peter M. Landwehr, Pavel Komarov, Mehdi Cherti, Kejia (KJ) Shi, Karlson Pfannschmidt, Fabian Linzberger, Christophe Cauet, Anna Gut, Andreas Mueller, and Alexander Fabisch. skopt. URL https://github.com/scikit-optimize/scikit-optimize.
 - Samuel Horváth, Aaron Klein, Peter Richtarik, and Cedric Archambeau. Hyperparameter transfer learning with adaptive complexity. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 1378–1386, 2021.
- Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. Sequential model-based optimization for general algorithm configuration. In *Learning and Intelligent Optimization: 5th International Conference*, 2011.
- Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. Identifying key algorithm parameters
 and instance features using forward selection. In *Learning and Intelligent Optimization: 7th International Conference*, 2013.
- Frank Hutter, Holger Hoos, and Kevin Leyton-Brown. An efficient approach for assessing hyperparameter importance. *International Conference on Machine Learning (ICML)*, 2014.
- 593 Frank Hutter, Lars Kotthoff, and Joaquin Vanschoren. *Automated Machine Learning: Methods, Systems, Challenges.* Springer Nature, 2019.

594 595 596	Carl Hvarfner, Danny Stoll, Artur Souza, Marius Lindauer, Frank Hutter, and Luigi Nardi. π bo: Augmenting acquisition functions with user beliefs for bayesian optimization. In <i>International</i> <i>Conference on Learning Representations (ICLR)</i> , 2022.				
598 599	Donald R Jones, Matthias Schonlau, and William J Welch. Efficient global optimization of expensive black-box functions. <i>Journal of Global Optimization</i> , 13:455–492, 1998.				
600 601 602	Ron Kohavi and George H. John. Automatic parameter selection by minimizing estimated error. In <i>International Conference on Machine Learning (ICML)</i> , 1995.				
603 604	Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. In <i>International Conference on Learning Representations (ICLR)</i> , 2019.				
605 606 607 608	Neeratyoy Mallik, Eddie Bergman, Carl Hvarfner, Danny Stoll, Maciej Janowski, Marius Lindauer, Luigi Nardi, and Frank Hutter. Priorband: Practical hyperparameter optimization in the age of deep learning. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> , 2023.				
609 610	J. Mockus. On the bayes methods for seeking the extremal point. <i>IFAC Proceedings Volumes</i> , 8(1, Part 1):428–431, 1975.				
611 612 613 614	Alejandro Molina, Antonio Vergari, Nicola Di Mauro, Sriraam Natarajan, Floriana Esposito, and Kristian Kersting. Mixed sum-product networks: A deep architecture for hybrid domains. In AAAI Conference on Artificial Intelligence, 2018.				
615 616 617	Julia Moosbauer, Julia Herbinger, Giuseppe Casalicchio, Marius Lindauer, and Bernd Bischl. Explain ing hyperparameter optimization via partial dependence plots. In <i>Advances in Neural Information</i> <i>Processing Systems (NeurIPS)</i> , 2021.				
618 619 620	Martin Mundt, Yongwon Hong, Iuliia Pliushch, and Visvanathan Ramesh. A wholistic view of continual learning with deep neural networks: Forgotten lessons and the bridge to active and open world learning. <i>Neural Networks</i> , 160:306–336, 2023.				
622 623 624	Robert Peharz, Sebastian Tschiatschek, Franz Pernkopf, and Pedro M. Domingos. On theoretical properties of sum-product networks. In <i>International Conference on Artificial Intelligence and Statistics (AISTATS)</i> , 2015.				
625 626 627	Valerio Perrone, Rodolphe Jenatton, Matthias W Seeger, and Cedric Archambeau. Scalable hyper- parameter transfer learning. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , volume 31, 2018.				
628 629 630	Hieu Pham, Melody Guan, Barret Zoph, Quoc Le, and Jeff Dean. Efficient neural architecture search via parameters sharing. In <i>International Conference on Machine Learning (ICML)</i> , 2018.				
631 632 633 634	Philipp Probst, Anne-Laure Boulesteix, and Bernd Bischl. Tunability: Importance of hyperparameters of machine learning algorithms. <i>The Journal of Machine Learning Research</i> , 20(1):1934–1965, 2019.				
635 636	Anil Ramachandran, Sunil Gupta, Santu Rana, Cheng Li, and Svetha Venkatesh. Incorporating expert prior in bayesian optimisation via space warping. <i>Knowledge-Based Systems</i> , 195:105663, 2020.				
637 638 639	Carl Rasmussen and Christopher Williams. <i>Gaussian Processes for Machine Learning</i> . MIT Pre 2006.				
640 641	David Salinas, Huibin Shen, and Valerio Perrone. A quantile-based approach for hyperparameter transfer learning. In <i>International Conference on Machine Learning (ICML)</i> , pp. 8438–8448, 2020.				
642 643 644	Sarah Segel, Helena Graf, Alexander Tornede, Bernd Bischl, and Marius Lindauer. Symbolic explanations for hyperparameter optimization. In <i>International Conference on Automated Machine Learning</i> , 2023.				
646 647	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. <i>Proceedings of the IEEE</i> , 104(1): 148–175, 2016.				

- Jasper Snoek, Hugo Larochelle, and Ryan P. Adams. Practical bayesian optimization of machine learning algorithms. In *Advances in Neural Information Processing Systems (NIPS)*, 2012.
- Artur L. F. Souza, Luigi Nardi, Leonardo B. Oliveira, Kunle Olukotun, Marius Lindauer, and Frank
 Hutter. Bayesian optimization with a prior for the optimum. In *European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML-PKDD)*,
 2021.
- ⁶⁵⁵ Joaquin Vanschoren. Meta-learning: A survey, 2018.

662

668

684

696 697

- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems (NIPS)*, 2017.
- Xilu Wang, Yaochu Jin, Sebastian Schmitt, and Markus Olhofer. Recent advances in bayesian
 optimization, 2022.
- Shuhei Watanabe, Archit Bansal, and Frank Hutter. PED-ANOVA: efficiently quantifying hyper parameter importance in arbitrary subspaces. In *International Joint Conference on Artificial Intelligence (IJCAI)*, 2023.
- Colin White, Sam Nolen, and Yash Savani. Exploring the loss landscape in neural architecture search.
 In *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2020.
- Colin White, Mahmoud Safari, Rhea Sukthanker, Binxin Ru, Thomas Elsken, Arber Zela, Debadeepta
 Dey, and Frank Hutter. Neural architecture search: Insights from 1000 papers. *arXiv preprint arXiv:2301.08727*, 2023.
- Martin Wistuba, Nicolas Schilling, and Lars Schmidt-Thieme. Sequential model-free hyperparameter
 tuning. In 2015 IEEE International Conference on Data Mining, pp. 1033–1038, 2015.
- 674
 675
 676
 676
 677
 677
 678
 679
 679
 679
 670
 670
 670
 671
 672
 674
 673
 674
 674
 674
 674
 674
 674
 675
 676
 677
 677
 678
 678
 679
 679
 679
 679
 679
 670
 670
 670
 671
 672
 674
 675
 675
 676
 677
 676
 677
 678
 678
 678
 679
 679
 679
 679
 679
 679
 679
 670
 670
 670
 670
 671
 672
 672
 673
 674
 674
 675
 675
 676
 676
 677
 678
 678
 678
 678
 678
 678
 679
 679
 679
 679
 679
 679
 679
 670
 670
 670
 670
 670
 671
 672
 672
 673
 674
 674
 675
 675
 676
 677
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
 678
- Dani Yogatama and Gideon S. Mann. Efficient transfer learning method for automatic hyperparameter
 tuning. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2014.
- Han Zhao, Tameem Adel, Geoff Gordon, and Brandon Amos. Collapsed variational inference for sum-product networks. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48, pp. 1310–1318. PMLR, 2016.
 - Barret Zoph and Quoc V. Le. Neural architecture search with reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2017.

702 A MOTIVATION & REAL-WORLD EXAMPLE

704 Reflecting user knowledge accurately is crucial for interactive HPO methods to fully benefit from 705 human knowledge and improve trustworthiness. Existing weighting scheme based methods like πBO 706 and BOPrO fail to reflect user priors accurately in their selection policy as it can be seen in Fig. 5 (a). Here, we show a 1d-example of a Branin function with an optimum around x = 0.5. The user prior 708 (in red) is placed at x = 0.3. Both π BO and BOPrO fail to select the next configuration from the high-density region of the prior; thus, the user prior is not incorporated in the selection process as 709 a user would expect. We followed the recommendation of (Hvarfner et al., 2022) and set $\beta = \frac{T}{10}$ 710 where we ran πBO for T = 10 iterations. Our method, IBO-HPC, solves this issue, which we now 711 demonstrate based on a real-world example. 712

713 To this end, we ran πBO , BOPrO – both of which leverage a weighting scheme to incorporate user 714 priors –, and IBO-HPC for T = 100 iterations on the CIFAR-10 task of the JAHS benchmark (Bansal 715 et al., 2022). Thus, we set the decay parameter of πBO to 10. We specified a Gaussian prior distribution with $\mu = 1$ and $\sigma = 0.3$ (Fig. 5 (b), purple) over the hyperparameter RESOLUTION 716 (R) that controls the down-/up-sampling rate of an image fed into a neural network. The rest of the 717 hyperparameters for this specific task (i.e. the network architecture and all other hyperparameters; 718 see App. D for details) were optimized by π BO, BOPrO and IBO-HPC without any user knowledge. 719 All methods received the same user prior (π BO and BOPrO from the beginning of the optimization; 720 IBO-HPC after 5 iterations). From the iteration the user prior was provided on, we then considered 721 the values chosen for RESOLUTION by πBO , BOPrO, and IBO-HPC for the next 20 iterations and 722 estimated a density of selected values for R (see Fig. 5 (b)). We chose 20 as the horizon under 723 consideration because for higher β , the prior is weighted down later in πBO (see (Hvarfner et al., 724 2022), Alg. 1) and BOPrO (see (Souza et al., 2021) Eq. 4). In the JAHS setup with T = 100 and 725 $\beta = 10$, the prior is weighted down after the 10th iteration in πBO and BOPrO. In the 20th iteration, π BO and BOPrO exponentially weigh down the prior with exponent 0.5. The density value of the 726 mode of our prior is then $1.26^{0.5} \approx 1.12$. For IBO-HPC, we chose the decay $\gamma = 0.995$; hence, 727 after 20 iterations, we get $1.26 \cdot \gamma^{20} \approx 1.14$ for the mode of the prior. Thus, we weigh down the 728 prior by approximately the same factor in πBO , BOPro, and IBO-HPC, ensuring a fair comparison. 729 We obtained that neither the choices for R by πBO (green dashed line) nor the choices of BOPrO 730 (red dashed line) reflect the user prior as specified. While πBO 's choices of RESOLUTION are 731 biased towards smaller values, BOPrO does not reflect the user's uncertainty well in its choices of 732 RESOLUTION. In contrast, IBO-HPC (blue solid line) precisely reflects the user prior as specified (up 733 to random variations due to sampling). 734



Figure 5: **IBO-HPC reflects user priors as specified.** In contrast to other weighting scheme based methods like π BO and BOPrO, IBO-HPC reflects the user prior as specified in its selection policy.

B PROOFS

750

751 752 753

754 755

In this section we provide the proof of Proposition 1 of the main paper.

756 B.1 IBO-HPC'S POLICY IS FEEDBACK ADHERING INTERACTIVE

Proposition 1 (IBO-HPC Policy is feedback adhering interactive). Given a search space Θ over hyperparameters \mathcal{H} , an HPC $s \in \mathcal{S}$, user knowledge $\mathcal{K} \in \mathcal{K}$ in form of a prior q over $\hat{\mathcal{H}} \subset \mathcal{H}$ s.t. $\int_{\mathcal{H} \setminus \hat{\mathcal{H}}} s(\mathcal{H} | F = f^*) \neq q(\hat{\mathcal{H}})$, the selection policy of IBO-HPC is feedback adhering interactive.

Proof. We have to show that the policy of IBO-HPC is feedback adhering, i.e. it conforms with Def. 3: The distribution over the configuration space used to obtain new configurations is different if user knowledge is provided from the distribution used if no user knowledge is provided (policy is efficacious) and the provided user knowledge is represented during configuration selection as specified (feedback adhering).

We first show that the selection policy of IBO-HPC is efficacious.

IBO-HPC selection policy is efficacious. Since the decay mechanism allowing IBO-HPC to recover from misleading knowledge can be treated as a constant in each iteration, it is enough if $s(\mathcal{H} \setminus \hat{\mathcal{H}} | \hat{\mathcal{H}} = \hat{\mathbf{h}}, F = f^*) \cdot q(\hat{\mathcal{H}} = \hat{\mathbf{h}}) \neq s(\mathcal{H} \setminus \hat{\mathcal{H}} | \hat{\mathcal{H}} = \emptyset, F = f^*) \cdot q(\hat{\mathcal{H}} = \emptyset)$ holds for any surrogate *s* representing a joint distribution over search space \mathcal{H} and prior *q* over $\hat{\mathcal{H}} \subset \mathcal{H}$ to make the policy efficacious. Note that we assume that \mathcal{K} is given in form of a prior $q(\hat{\mathcal{H}})$ over $\hat{\mathcal{H}}$ as before. Since $\emptyset \notin \hat{\mathcal{H}}$ is assumed, our policy ignores any prior if no user knowledge is provided. Thus, in this case, the policy samples from the distribution

$$s(\mathcal{H}|F = f^*) = s(\mathcal{H} \setminus \hat{\mathcal{H}}|\hat{\mathcal{H}}, F = f^*) \cdot \int_{\mathcal{H} \setminus \hat{\mathcal{H}}} s(\mathcal{H}|F = f^*)$$
(5)

⁷⁷⁹ Since $s(\mathcal{H} \setminus \hat{\mathcal{H}} | \hat{\mathcal{H}}, F = f^*)$ is the same, regardless of whether user knowledge is given or not, user knowledge will lead to a different distribution if $\int_{\mathcal{H} \setminus \hat{\mathcal{H}}} s(\mathcal{H} | F = f^*) \neq q(\hat{\mathcal{H}})$ holds. Since Prop. 1 demands that this is the case, our policy is efficacious according to Def. 2.

783 We can now proceed and show feedback adherence of the IBO-HPC selection policy.

784 **IBO-HPC selection policy is feedback adhering.** The proof that our policy is feedback adhering 785 directly follows by design: If a user prior $q(\hat{\mathcal{H}})$ is given, Eq. 3 is approximated by sampling N 786 conditions $\mathbf{h}'_{1,\dots,N} \sim q(\hat{\mathcal{H}})$ and computing N conditionals $s(\mathcal{H} \setminus \hat{\mathcal{H}} | \hat{\mathcal{H}} = h'_1, F = f^*), \dots, s(\mathcal{H} \setminus \hat{\mathcal{H}})$ 787 $\hat{\mathcal{H}}|\hat{\mathcal{H}} = h'_N, F = f^*)$. We can approximate $q(\hat{\mathcal{H}})$ arbitrarily close with $N \to \infty$. To select 788 the next configuration, we sample B configurations from each of the N conditionals and select the 789 configuration maximizing $s(\mathcal{H}|F=f^*)$ for each conditional. This leaves us with N candidates. Note 790 that at this point, the hyperparameters $\hat{\mathcal{H}}$ still follow $q(\hat{\mathcal{H}})$ with $N \to \infty$ as the conditions of $s(\mathcal{H} \setminus \mathcal{H})$ 791 $\hat{\mathcal{H}}|\hat{\mathcal{H}}=h'_1,F=f^*),\ldots,s(\mathcal{H}\setminus\hat{\mathcal{H}}|\hat{\mathcal{H}}=h'_N,F=f^*)$ remain fixed and only hyperparameters 792 793 of the set $\mathcal{H} \setminus \mathcal{H}$ can vary/are sampled. Thus, maximizing the likelihood $s(\mathcal{H}|F = f^*)$ is only 794 done w.r.t. hyperparameters in $\mathcal{H} \setminus \mathcal{H}$. This implies that sampling hyperparameters $\mathcal{H} \setminus \mathcal{H}$ can be biased while sampling from $q(\hat{\mathcal{H}})$ is unaffected because the conditions $\mathbf{h}'_{1,\dots,N}$ are sampled first 796 in i.i.d. fashion. Our policy selects the configuration evaluated next by uniformly sampling from the remaining N candidates. Since uniformly sampling L times from a set of N samples from a 797 distribution q results in approximating q arbitrarily close for $N \to \infty$ and $L \to \infty$, we conclude that 798 user priors are exactly reflected as specified in our selection policy. This concludes our proof that the 799 selection policy of IBO-HPC is efficacious and feedback adhering. 800

801 802

807

777

778

B.2 IBO-HPC MINIMIZES SIMPLE REGRET

803 804 We introduce the following proposition:

Proposition 4 (IBO-HPC minimizes Simple Regret). *IBO-HPC minimizes simple regret, which is defined as* $r = f(\mathbf{h}) - f(\mathbf{h}^*)$ *for a hyperparameter configuration* $\mathbf{h} \in \Theta$ *and global optimum* \mathbf{h}^* .

Proof. Assume that w > 0 holds for each weight w of a PC s, that each leaf node of s is a distribution p s.t. p(x) > 0 for some x and assume f is not noisy. Then, the PC fulfills the positivity assumption, i.e. $s(\mathcal{H} = \mathbf{h}, F = f(\mathbf{h})) > 0$. It follows that $s(\mathcal{H} = \mathbf{h}|F = f^*) > 0$ for any f^* and any $\mathbf{h} \in \Theta$.

Thus, with iterations $T \to \infty$, the probability of sampling the global optimum \mathbf{h}^* in one of the iterations gets 1, and thus $r = f(\mathbf{h}^*) - f(\mathbf{h}^*) = 0$.

813 B.3 CONVERGENCE SPEED OF IBO-HPC

In this section, we analyze the convergence speed of IBO-HPC at each iteration. Therefore, let us state a well-known result of the PC literature on which our analysis is based.

Definition 4. Induced Trees (Zhao et al., 2016). Given a complete and decomposable PC s over $\mathcal{H} = \{H_1, \ldots, H_n\}, \mathcal{T} = (\mathcal{T}_V, \mathcal{T}_E)$ is called an induced tree PC from s if

- 1. $N \in \mathcal{T}_V$ where N is the root of s.
- 2. for all sum nodes $S \in T_V$, exactly one child of S in s is in T_V , and the corresponding edge is in T_E .
- 3. for all product node $P \in T_V$, all children of P in s are in T_V , and the corresponding edges in T_E .
- 827 We can use Def. 4 to represent decomposable and complete PCs as mixtures (Zhao et al., 2016).

828 **Proposition 5** (Induced Tree Representation). Let τ_s be the total number of induced trees in s. Then 829 the output at the root of s can be written as $\sum_{t=1}^{\tau_s} \prod_{(k,j) \in \mathcal{T}_{tE}} w_{kj} \prod_{i=1}^n p_t(H_i = h_i)$, where \mathcal{T}_t is 830 the t-th unique induced tree of s and $p_t(H_i)$ is a univariate distribution over H_i in \mathcal{T}_t as a leaf node. 831

832 With this, we are ready to analyze the convergence speed of IBO-HPC in each iteration. Assume a non-noisy differentiable L-Lipschitz continuous function $f: \mathbb{R}^d \to \mathbb{R}$ with global optimum $\mathbf{h}^* \in \mathbb{R}^d$ 833 that is convex within a ball $B_r(\mathbf{h}^*) = {\mathbf{h} \in \mathbb{R}^d : ||\mathbf{h} - \mathbf{h}^*|| < r}$. Further, assume we have given 834 a dataset $\mathcal{D} = \{(\mathbf{h}_1, y_1), \dots, (\mathbf{h}_n, y_n)\}$ where all $\mathbf{h}_i \in B_r$ and $y_i = f(\mathbf{h}_i)$ and a decomposable, 835 complete PC s over $\mathcal{H} \cup \{F\}$ where the support of $\mathcal{H} = B_r$ and the support of $F = \mathbb{R}$. Assume 836 s locally maximizes the likelihood over \mathcal{D} and that all leaves are Gaussians. Note that LearnSPN 837 yields decomposable and complete PCs that locally maximize the likelihood of the given data (Gens 838 & Domingos, 2013). 839

We analyze the convergence properties of our algorithm by examining the expected improvement (EI) in each iteration. Therefore, denote the best score obtained until iteration t as y_t^* and its corresponding configuration as \mathbf{h}_t^* . For better readability, we write $s(\mathcal{H} = \mathbf{h}|F = y_t^*)$ as $s(\mathbf{h}|y_t^*)$ from now on. Then, the expected improvement of IBO-HPC within B_r is given by

847

815

816

817

818

819

820 821

822

823

824

825 826

 $\int_{\mathbf{h}\in B_r} s(\mathbf{h}|y_t^*) \cdot \mathbb{I}[f(\mathbf{h}) < y_t^*] \cdot f(\mathbf{h})$ (6)

$$= \int_{\mathbf{h}^*}^{\mathbf{h}^*_t} s(\mathbf{h}|y_t^*) \cdot f(\mathbf{h}).$$
(7)

848 849 850

858 859

861

Here, w.l.o.g. we assume that $\mathbf{h}_k^* < \mathbf{h}_{tk}^*$ for all dimensions $k \in \{1, \dots, d\}$ and call I the indicator function. Using Prop. 5, the fact that the first product of the induced tree representation of a PC *s* acts as an edge selector, the fact that the conditional of a PC is a PC again, and the Gaussian leaf parameterization of *s*, we can write *s* as a Gaussian Mixture, i.e., $s(\mathbf{h}|y_t^*) = \sum_{i=1}^{\tau_s} w_i \phi(\mathbf{h}; \boldsymbol{\mu}_i, \Sigma_i)$. Here, ϕ is the density of the Gaussian distribution parameterized by mean $\boldsymbol{\mu}$ and covariance matrix Σ and corresponds to the second product in the induced tree representation of *s*. Thus, Eq. 6 can be rewritten as

$$\sum_{i=1}^{\tau_s} w_i \int_{\mathbf{h}^*}^{\mathbf{h}_t^*} \phi(\mathbf{h}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \cdot f(\mathbf{h}).$$
(8)

B62 Due to the *L*-Lipschitz assumption, $||f(\mathbf{h}) - f(\mathbf{h}')|| \le L \cdot ||\mathbf{h} - \mathbf{h}'||$ holds for all $\mathbf{h}, \mathbf{h}' \in B_r$. Hence, we can use a Taylor approximation and write $f(\mathbf{h}) \approx f(\mathbf{h}^*) + \nabla f(\mathbf{h}^*) \cdot ||\mathbf{h} - \mathbf{h}^*||$ which is upper bounded by $f(\mathbf{h}^*) + L||\mathbf{h} - \mathbf{h}^*||$. Then, we can write an upper bound of EI as

$$\sum_{i=1}^{\tau_s} w_i \int_{\mathbf{h}^*}^{\mathbf{h}^*_t} \phi(\mathbf{h}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \cdot (f(\mathbf{h}^*) + L||\mathbf{h} - \mathbf{h}^*||)$$
(9)

$$=\sum_{i=1}^{\tau_s} w_i \left(\int_{\mathbf{h}^*}^{\mathbf{h}^*_t} \phi(\mathbf{h}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \cdot f(\mathbf{h}^*) + \int_{\mathbf{h}^*}^{\mathbf{h}^*_t} \phi(\mathbf{h}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \cdot L ||\mathbf{h} - \mathbf{h}^*||) \right)$$
(10)

$$=\sum_{i=1}^{\tau_s} w_i \cdot \left(f(\mathbf{h}^*) \cdot \int_{\mathbf{h}^*}^{\mathbf{h}^*_t} \phi(\mathbf{h}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) + \mathbb{E}_{\phi_i}[g_{\mathbf{h}^*}(\mathbf{h})] \right).$$
(11)

In the last step, we defined $g_{\mathbf{h}^*}(\mathbf{h}) := L||\mathbf{h} - \mathbf{h}^*||$. Note that we take the expectation w.r.t. the truncated normal distribution because we consider the interval $[\mathbf{h}^*, \mathbf{h}_t^*]$. Also note that $f(\mathbf{h}^*)$ is constant. Thus, we can omit it for the sake of convergence analysis. Since $g_{\mathbf{h}^*}$ is linear, we can use the linearity of the expectation and write

$$\sum_{i=1}^{\tau_s} w_i \cdot \left(\int_{\mathbf{h}^*}^{\mathbf{h}^*_t} \phi(\mathbf{h}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) + g_{\mathbf{h}^*}(\mathbb{E}_{\phi_i}[\mathbf{h}]) \right)$$
(12)

$$i=1$$
 au_s

$$=\sum_{i=1}^{\tau_s} w_i \cdot \left(\left(\Phi(\mathbf{h}_t^*; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) - \Phi(\mathbf{h}^*; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \right) + L || \mathbb{E}_{\phi_i}[\mathbf{h}] - \mathbf{h}^* || \right),$$
(13)

where $\Phi(\mathbf{h}; \boldsymbol{\mu}, \Sigma)$ is the cumulative distribution function of multivariate Gaussian. Since the expectations $\mathbb{E}_{\phi_i}[\mathbf{h}]$ are taken over the truncated normal, they can be lower bounded by $\boldsymbol{\mu} + \alpha \cdot \operatorname{diag}(\Sigma)$. Thus, we have to set a series of α_i where each $\alpha_i = \min(\mathbf{h}_t^* - \boldsymbol{\mu}_i, \mathbf{h}^* - \boldsymbol{\mu}_i)$. Then, we can write

$$\sum_{i=1}^{\tau_s} w_i \cdot \left(\left(\Phi(\mathbf{h}_t^*; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) - \Phi(\mathbf{h}^*; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \right) + L || (\boldsymbol{\mu}_i + \alpha_i \cdot \operatorname{diag}(\boldsymbol{\Sigma}_i)) - \mathbf{h}^* || \right).$$
(14)

Setting $\epsilon_i = ||\boldsymbol{\mu}_i + \alpha_i \cdot \operatorname{diag}(\Sigma_i) - \mathbf{h}^*||$ and splitting the sum yields

 $\sum_{i=1}^{\tau_s} w_i \cdot \Phi(\mathbf{h}_t^*; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) - \sum_{i=1}^{\tau_s} w_i \cdot \Phi(\mathbf{h}^*; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) + \sum_{i=1}^{\tau_s} w_i L \epsilon_i.$ (15)

Using that the cumulative multivariate Gaussian $\Phi(\mathbf{h}_t^*; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ can be lower bounded by $\prod_{j=1}^d \Phi(\mathbf{h}_{ti}^*; \boldsymbol{\mu}_{ij}, \boldsymbol{\Sigma}_{ijj})$, we can lower-bound the entire equation, giving us

$$\sum_{i=1}^{\tau_s} w_i \cdot \prod_{j=1}^d \Phi(\mathbf{h}_{tj}^*; \boldsymbol{\mu}_{ij}, \boldsymbol{\Sigma}_{i_{jj}}) - \sum_{i=1}^{\tau_s} w_i \cdot \prod_{j=1}^d \Phi(\mathbf{h}_j^*; \boldsymbol{\mu}_{ij}, \boldsymbol{\Sigma}_{i_{jj}}) + \sum_{i=1}^{\tau_s} w_i L\epsilon_i.$$
(16)

Since
$$\Phi(\frac{x-\mu}{\sigma}) = \frac{1}{2} \left(1 + \operatorname{erf}\left(\frac{x-\mu}{\sigma\sqrt{2}}\right) \right)$$
 holds, we rewrite

$$\sum_{i=1}^{\tau_s} w_i \cdot \prod_{j=1}^d \operatorname{erf}\left(\frac{\mathbf{h}_{ij}^* - \boldsymbol{\mu}_{ij}}{\Sigma_{i_{jj}}\sqrt{2}}\right) - \sum_{i=1}^{\tau_s} w_i \cdot \prod_{j=1}^d \operatorname{erf}\left(\frac{\mathbf{h}_j^* - \boldsymbol{\mu}_{ij}}{\Sigma_{i_{jj}}\sqrt{2}}\right) + \sum_{i=1}^{\tau_s} w_i L\epsilon_i$$
(17)

 $=\sum_{i=1}^{\tau_s} w_i \cdot \Big(\prod_{j=1}^d \operatorname{erf}\Big(\frac{\mathbf{h}_{tj}^* - \boldsymbol{\mu}_{ij}}{\Sigma_{i_{jj}}\sqrt{2}}\Big) - \prod_{j=1}^d \operatorname{erf}\Big(\frac{\mathbf{h}_j^* - \boldsymbol{\mu}_{ij}}{\Sigma_{i_{jj}}\sqrt{2}}\Big) + L\epsilon_i\Big).$ (18)

917 Note that we dropped constants and scaling by $\frac{1}{2}$ of the error function as it does not affect the overall result.

918 Intuitively spoken, the EI is lower bounded by the cumulative probability mass (given by error 919 function erf) within the region defined by the largest discrepancy between minimal error w.r.t. to the 920 observed data (i.e., bad convergence when s overfits) and the maximal error w.r.t. h^* (i.e., \mathcal{D} does 921 not contain points close to the optimum), multiplied by a linear approximation of the objective f922 between the best observed configuration \mathbf{h}_t^* and \mathbf{h}^* .

923 Note that this result does not incorporate user knowledge. The analysis of the effect of user knowledge 924 is straightforward. If helpful user knowledge is given, this can be seen as shifting at least one 925 dimension j of at least one mean vector μ_k by some δ towards \mathbf{h}^* , i.e., $\mu_{*k} = \mu_k + (0, \dots, \delta, \dots, 0)$. 926 Then, assuming all Σ_i stay as above, 927

928 929

930 931 932

938

939

940 941

942

$$\sum_{i=1}^{\tau_s} w_i \cdot \Big(\prod_{j=1}^d \operatorname{erf}\Big(\frac{\mathbf{h}_{tj}^* - \boldsymbol{\mu}_{ij}}{\Sigma_{i_{jj}}\sqrt{2}}\Big) - \prod_{j=1}^d \operatorname{erf}\Big(\frac{\mathbf{h}_j^* - \boldsymbol{\mu}_{ij}}{\Sigma_{i_{jj}}\sqrt{2}}\Big) + L\epsilon_i\Big)$$

$$\leq \sum_{i=1, i \neq k}^{\tau_s} w_i \cdot \Big(\prod_{j=1}^a \operatorname{erf}\Bigl(\frac{\mathbf{h}_{tj}^* - \boldsymbol{\mu}_{ij}}{\Sigma_{i_{jj}}\sqrt{2}}\Bigr) - \prod_{j=1}^a \operatorname{erf}\Bigl(\frac{\mathbf{h}_j^* - \boldsymbol{\mu}_{ij}}{\Sigma_{i_{jj}}\sqrt{2}}\Bigr) + L\epsilon_i \Big)$$

$$i=1, i\neq k$$
 $j=1$

$$+ w_k \cdot \Big(\prod_{j=1}^a \operatorname{erf}\Big(\frac{\mathbf{h}_{tj}^* - \boldsymbol{\mu}_{kj}}{\Sigma_{k_{jj}}\sqrt{2}}\Big) - \prod_{j=1}^a \operatorname{erf}\Big(\frac{\mathbf{h}_j^* - \boldsymbol{\mu}_{kj}}{\Sigma_{k_{jj}}\sqrt{2}}\Big) + L\epsilon_k\Big).$$

This is easy to see since the distribution we sample configurations from is shifted towards the global optimum h^* , thus increasing the probability of sampling a configuration closer to h^* , ultimately leading to faster convergence.

B.4 ACCURACY OF IBO-HPC'S SELECTION POLICY

943 Here, we briefly discuss the accuracy of IBO-HPC's policy in selecting new configurations for 944 evaluation based on the obtained data (see Eq. 2). Note that the sampling from the distribution 945 provided in Eq. 2 is accurate if (1) the s represents the data \mathcal{D} accurately and (2) sampling from s and 946 the prior q is unbiased (i.e., samples are drawn according to the underlying distribution). Let us start 947 with (1). Since we employ LearnSPN (Gens & Domingos, 2013) to obtain s (a PC in form of SPN), s will locally maximize the log-likelihood of the training data (i.e., the configuration-evaluation pairs 948 obtained). This means that there is no other SPN in the space of the learnable SPNs via LearnSPN 949 that achieves a better log-likelihood given the data.² Hence, as long as the ground truth distribution p950 (or a good approximation of it) is representable by an SPN, we can recover p with arbitrarily small 951 error with iterations $T \to \infty$. 952

Looking at (2), we sample from two distributions when selecting a new configuration. First, we 953 sample from the prior q, then from the conditional $s(\mathcal{H} \setminus \hat{\mathcal{H}} | \hat{\mathcal{H}} = \mathbf{h}, F = f^*)$ where $\mathbf{h} \sim q$. 954 Assuming q is a tractable distribution (e.g., a parametric one such as an isotropic Gaussian), sampling 955 is immediate and not biased (i.e., performed via simple transformations such as the Box-Muller 956 transform). Note that the assumption on q being a tractable (and relatively simple) distribution can 957 be made safely since providing highly complex distributions as user knowledge is hard to do for 958 most users. When considering sampling from the conditional $s(\mathcal{H} \setminus \hat{\mathcal{H}} | \hat{\mathcal{H}} = \mathbf{h}, F = f^*)$, it should 959 be noted that this conditional is a valid PC again (specifically, a PC in the form of an SPN when 960 obtained with LearnSPN). The model is unchanged and only evaluated differently, i.e., by providing 961 the partial evidence at leaves and evaluating the model bottom-up first (see Choi et al. (2020)). Then, 962 PC sampling is performed top-down by sampling from the simple categorical variables represented 963 by the sum nodes and then from the selected univariate leaves. Thus, the process is tractable (linear 964 in the circuit size) and not biased by further operations or assumptions (Choi et al., 2020). Thus, we 965 conclude that the approximation in Eq. 2 is accurate in the limit $N, T \rightarrow \infty$.

966 967 968

969

971

C PROBABILISTIC CIRCUITS

Since probabilistic circuits (PCs) are a key component of our method, we provide more details on 970 these models in the following. Let us first start with a rigorous definition of PCs.

²Assuming an oracle for the variable splitting. See Proposition 1 in Gens & Domingos (2013).

972 **Definition 5.** A probabilistic cricuit (PC) is a computational graph encoding a distribution over a 973 set of random variables **X**. It is defined as a tuple (\mathcal{G}, ϕ) where $\mathcal{G} = (V, E)$ is a rooted, directed 974 acyclic graph and $\phi: V \to 2^{\mathbf{X}}$ is the scope function assigning a subset of random variables to each 975 node in \mathcal{G} . For each internal node N of \mathcal{G} , the scope is defined as the union of scopes of its children, *i.e.* $\phi(N) = \bigcup_{N' \in ch(N)}$. Each leaf node L computes a distribution/density over its scope $\phi(L)$. All 976 internal nodes of \mathcal{G} are either a sum node S or a product node P where each sum node computes a 977 convex combination of its children, i.e., $S = \sum_{N \in ch(S)} w_{S,N}N$, and each product computes a product 978 of its children, i.e., $P = \prod_{N \in ch(P)} N$. 979

980

With this definition at hand, we describe the tractable key operations of PCs relevant to our method in more detail.

Inference. Inference in PCs is a bottom-up procedure. To compute the probability of given evidence 984 $\mathbf{X} = \mathbf{x}$, the densities of the leaf nodes are evaluated first. This yields a density value for each leaf. 985 The leaf densities are then propagated bottom-up by computing all product/sum nodes. Eventually, 986 the root node holds the probability/density of x. Note that typically, multiple leaf nodes correspond to 987 the same random variable. Thus, if the children of a sum node have the same scope, we can interpret 988 sum nodes as mixture models. Conversely, if the children of a product node have *non*-overlapping 989 scopes, a product node can be interpreted as a product distribution of two (independent) random 990 variables. We call these two properties smoothness and decomposability. More formally, smoothness 991 means that for each sum node $S \in V$ it holds that $\phi(N) = \phi(N')$ for $N, N' \in ch(S)$. Decomposability 992 means that for each product node $P \in V$ it holds that $\phi(N) \cap \phi(N') = \emptyset$ for $N, N' \in ch(P), N \neq N'$. 993 Hence, PCs can be interpreted as hierarchical mixture models.

Marginalization. Decomposability implies that marginalization is tractable in PCs and can be done in linear time of the circuit size. This is because integrals that can be rewritten by nesting single-dimensional integrals can be computed only in terms of leaf integrals, which are assumed to be tractable as they follow certain distributions (e.g., Gaussian). Computing such nested integrals only in terms of leaf integrals is possible because single-dimensional integrals commute with the sum operation and affect only a single child of product nodes. For more details on the computational implications of decomposability, refer to (Peharz et al., 2015).

Practically, there are two ways to marginalize certain variables from the scope of a PC. One approach is structure-preserving, and marginalization is achieved by setting all leaves corresponding to the set of random variables that are supposed to be marginalized to 1. The second approach constructs a new PC representing the marginal distribution, i.e. the structure of the PC is changed. The second approach is beneficial if samples should be drawn from the marginalized PC because the sampling procedure remains the same, i.e. the PC is adopted to obtain the marginal distribution, not vice versa.

Conditioning. Computing a conditional distribution $p(\mathbf{X}_1|\mathbf{X}_2) = \frac{p(\mathbf{X})}{\int_{\mathbf{X}_2} p(\mathbf{X})}$ where $\mathbf{X}_1 \cup \mathbf{X}_2 = \mathbf{X}$ and $\mathbf{X}_1 \cap \mathbf{X}_2 = \emptyset$ is achieved by combining marginalization (denominator) and inference (numerator). Since inference is tractable for PCs in general and marginalization is tractable for decomposable PCs, conditioning is also tractable.

Sampling. Sampling in PCs is a top-down procedure and recursively samples a sub-tree, starting at the root. Each sum node S holds a parameter vector \mathbf{w} s.t. $\sum_{i=0}^{|\operatorname{ch}(S)|} \mathbf{w}_i = 1$. Based on the distribution induced by \mathbf{w} , one of the children of S is sampled as a sub-tree. By decomposability, the scope of the children of a product node is non-overlapping; thus, sampling from a product node corresponds to sampling from all its child nodes. If a leaf node is reached, a sample is obtained from the distribution at that leaf.

1018 Learning. Learning PCs consists of two steps: Identify the structure of the PC and learn the parame-1019 ters of the PC. A common approach to learning both the structure and parameters is LearnSPN (Gens 1020 & Domingos, 2013). We employ LearnSPN to learn the PC after obtaining new data. The basic 1021 idea of LearnSPN is to split the data by alternating clustering (i.e., split the data along the sample dimension) and independence tests (i.e., split the data along the features dimension). In other words, 1023 the data matrix is split by rows (samples) and columns (features). Usually, rows are clustered when the independence test fails in splitting the features. Clusters correspond to sum nodes in the learned 1024 PC, while product nodes correspond to successfully passed independence tests (assessing that two 1025 subsets of features are statistically independent). The parameters (i.e., weights of sum nodes) are

set proportional to the cluster sizes of clusters represented by the child nodes of a sum node. Leaf parameters are commonly defined via maximum likelihood estimation.

D EXPERIMENTAL DETAILS

1029

1030 1031

Here we present additional details of our empirical evaluation.

1033 1034 D.1 SEARCH SPACE EXTENSION OF JAHS

To make the HPO problem on JAHS more challenging, we decided to extend the search space slightly as JAHS – as a surrogate benchmark – allows us to query hyperparameter values which were not tested explicitly in the benchmark. We defined three search spaces for JAHS which are presented in the following table.

1040		S1	S2	S 3
1041	Activation	[Mish, ReLU, Hardswish]	[Mish, ReLU, Hardswish]	[Mish, ReLU, Hardswish]
1042	Learning Rate	[1e-3, 1e0]	[1e-3, 1e0]	[1e-3, 1e0]
1043	Weight Decay	[1e-5, 1e-2]	[1e-5, 1e-2]	[1e-5, 1e-2]
1044	Trivial Argument	[True, False]	[True, False]	[True, False]
1045	Opl	0-6	0-6	0-6
1045	Op2 Op3	0-0	0-6	0-6
1046	Op3 Op4	0-6	0-6	0-6
1047	Op4 Op5	0-6	0-6	0-6
1048	Op6	0-6	0-6	0-6
1049	N	1-15	1-11	1-5
1050	W	1-31	1-23	1-16
1051	Epoch	1-200	1-200	1-200
1052	Resolution	0-1	0-1	0-1
1053				
1054	Table 1: JAHS Se	arch Space. We define the	ree versions of the JAHS	search space, ranging from
1055	simpler to harder s	paces.		
1056				
1057				
1057				
1050				
1059				
1060				
1061				
1062				
1063				
1064				
1065				
1066				
1067				
1068				
1069				
1070				
1071				
1072				
1072				
1073				
1074				
10/5				
1076				
1077				
1078				
1079				

D.2 INTERACTIONS

¹⁰⁸² Here we provide the interactions used for our experiments.

JAHS The following JSON code shows the interactions performed in our JAHS experiments. The first interaction is a misleading interaction, followed by a beneficial interaction and a no interaction (for recovery).

```
1087
      [
1088
           {
1089
               "type": "bad",
               "intervention": {"Activation": 1, "LearningRate":
1090
      0.8201676371308472, "N": 15,
1091
               "Op1": 3, "Op2": 4, "Op3": 1, "Op4": 2, "Resolution":
1092
      0.5096959403985494,
1093
               "TrivialAugment": 0, "W": 14,
1094
                "WeightDecay": 0.002697686639935806, "epoch": 10},
1095
               "iteration": 5
1096
           },
1097
           {
               "type": "good",
1099
               "intervention":
                                {"N": 3, "W": 16, "Resolution": 1},
1100
               "iteration": 15
1101
           },
1102
           {
               "type": "good",
1103
               "intervention": null,
1104
               "iteration": 20
1105
           },
1106
           {
1107
               "type": "good",
1108
               "kind": "dist",
1109
               "intervention": {"N": {"dist": "cat", "parameters":
1110
               [1, 1, 1, 1e4, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1]},
1111
               "W": {"dist": "cat", "parameters":
               1112
               "Resolution": {"dist": "uniform", "parameters": [0.98,
1113
      1.02]}},
1114
               "iteration": 5
1115
           }
1116
      ]
1117
      NAS-Bench-101 The following JSON code shows the interactions performed in our experiments on
1118
      NAS-Bench-101. The first interaction is a misleading interaction, followed by a beneficial interaction
1119
      and a no interaction (for recovery).
1120
1121
      [
1122
           {
1123
               "type": "bad",
1124
               "kind": "point",
1125
               "intervention": [0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1,
1126
      0,
         0, 1, 1, 1, 0, 1],
1127
               "iteration": 5
1128
           },
1129
           {
               "type": "good",
1130
1131
               "kind": "point",
               "intervention": [1, 0, 1, 0, 1, 1, 1, 0, 0, 0, 0, 0, 1, 0,
1132
      0, 0, 1, 0, 1, 0, 1],
1133
               "iteration": 12
```

```
},
1135
           {
1136
                "type": "good",
1137
                "kind": "point",
1138
                "intervention": null,
                "iteration": 20
1139
           },
1140
           {
1141
                "type": "good",
1142
                "kind": "dist",
1143
                "intervention": {
1144
                     "e_0_1": {"dist": "cat", "parameters": [1, 1e4]},
1145
                     "e_0_2": {"dist": "cat", "parameters": [1e4, 1]},
1146
                     "e_0_3": {"dist": "cat", "parameters": [1, 1e4]},
1147
                     "e_0_4": {"dist": "cat", "parameters": [1e4, 1]},
1148
                     "e_0_5": {"dist": "cat", "parameters": [1, 1e4]},
                    "e_0_6": {"dist": "cat", "parameters": [1, 1e4]},
1149
                     "e_1_2": {"dist": "cat", "parameters": [1, 1e4]},
1150
                    "e_1_3": {"dist": "cat", "parameters": [1e4, 1]},
1151
                    "e_1_4": {"dist": "cat", "parameters": [1e4, 1]},
1152
                     "e_1_5": {"dist": "cat", "parameters": [1e4, 1]},
1153
                    "e_1_6": {"dist": "cat", "parameters": [1e4, 1]},
1154
                    "e_2_3": {"dist": "cat", "parameters": [1e4, 1]},
1155
                    "e_2_4": {"dist": "cat", "parameters": [1, 1e4]},
1156
                    "e_2_5": {"dist": "cat", "parameters": [1e4, 1]},
1157
                     "e_2_6": {"dist": "cat", "parameters": [1e4, 1]},
1158
                     "e_3_4": {"dist": "cat", "parameters": [1e4, 1]},
1159
                     "e_3_5": {"dist": "cat", "parameters": [1, 1e4]},
                     "e_3_6": {"dist": "cat", "parameters": [1e4, 1]},
1160
                    "e_4_5": {"dist": "cat", "parameters": [1, 1e4]},
1161
                    "e_4_6": {"dist": "cat", "parameters": [1e4, 1]},
"e_5_6": {"dist": "cat", "parameters": [1, 1e4]}
1162
1163
1164
                "iteration": 5
1165
           }
1166
      1
1167
      NAS-Bench-201 The following JSON code shows the interactions performed in our experiments on
1168
      NAS-Bench-201. The first interaction is a misleading interaction, followed by a beneficial interaction
1169
      and a no interaction (for recovery).
1170
1171
1172
       [
1173
           {
                "type": "good",
1174
                "kind": "point",
1175
                "intervention": {"Op_0": 2, "Op_1": 2, "Op_2": 0},
1176
                "iteration": 5
1177
           },
1178
           {
1179
                "type": "bad",
1180
                "kind": "point",
1181
                "intervention": {"Op_0": 1, "Op_1": 2, "Op_2": 1},
1182
                "iteration": 5
1183
           },
1184
           {
1185
                "type": "good",
1186
                "kind": "point",
1187
                "intervention": null,
```

```
1188
                "iteration": 20
1189
           },
1190
           {
1191
                "type": "good",
1192
                "kind": "dist",
               "intervention": {"Op_0": {"dist": "cat", "parameters": [1,
1193
      1, 1e4, 1, 1]},
1194
                                   "Op_1": {"dist": "cat", "parameters": [1,
1195
      1, 1e4, 1, 1]},
1196
                                   "Op 2": {"dist": "cat", "parameters":
1197
      [1e4, 1, 1, 1, 1]}},
1198
               "iteration": 5
1199
           1
      ]
1201
1202
1203
      D.3 FURTHER RESULTS & ABLATIONS
1204
1205
```

In this section, we provide further results and ablations. Fig. 6 provides additional results on two 1206 challenging tasks of the HPO-B benchmark (search space IDs 6767, 6794; dataset ID 31). Both search 1207 spaces have more than 10 hyperparameters and the goal is to solve a classification task. IBO-HPC 1208 outperforms the baselines or is competitive with the baselines in both cases, i.e., where feedback 1209 is given, and no user feedback is given. Fig 7 shows results of IBO-HPC on JAHS, NAS201, and 1210 NAS101 where the given user feedback was either a fixed value or a distribution over configurations. 1211 Both cases are handled well by IBO-HPC, demonstrating its flexibility. Fig. 8 provides a more detailed 1212 view of the effectiveness of IBO-HPC and its recovery mechanism. It can be seen that IBO-HPC 1213 successfully recovers from harmful user feedback in JAHS and NAS-201 (----). Also, it can be seen 1214 that IBO-HPC handles alternating and contradictory user feedback well by leveraging information from beneficial feedback and ignoring harmful feedback (----). In NAS-101, however, IBO-HPC is 1215 less effective in general, which can be explained by the extreme sparsity of the NAS-101 benchmark. 1216 While NAS-101 and NAS-201 are highly similar, NAS-101 uses a binary encoding of architectures, 1217 while NAS-201 uses a much denser dictionary-like representation. Although both benchmarks are 1218 highly similar, IBO-HPC performs well on NAS-201 but is not as effective on NAS-101, underlining 1219 our explanation. 1220

Fig. 9 shows the CDF of test accuracy across the baselines and IBO-HPC. It can be seen that IBO-HPC invests more computational resources in good-performing configurations than other methods while achieving state-of-the-art or better results. In other words, IBO-HPC avoids exploration in unpromising regions of the search space. This is because IBO-HPC samples configurations from a conditional distribution where the condition is the best evaluation score obtained. Thus, exploration is purely data-driven and focuses on regions that perform similarly to the incumbent at a particular iteration.

Fig. 10 shows the influence of the decay parameter γ in cases where harmful or misleading user knowledge was provided to IBO-HPC at an early iteration (10 in this case). It can be seen that for higher γ , IBO-HPC requires more time to recover than for smaller γ . This aligns with our expectations since a larger γ corresponds to a high likelihood of the user knowledge being used for many iterations. In contrast, if γ is small, likely, the user knowledge is only considered for a certain number of iterations with high likelihood. Thus, for smaller γ IBO-HPC can recover faster.

Fig. 11 shows the effect of conditioning on the {0.25, 0.5, 0.75}-quantile of the obtained evaluation scores instead of the maximum evaluation score. As expected, the higher the quantile, the better the performance of IBO-HPC as we aim to maximize the objective function. Thus, conditioning on higher values guides the optimization algorithm to configurations that yield better evaluation scores.

1238 Lastly, Fig. 12 depicts the effect of changing L, i.e. the number of samples drawn from the 1239 surrogate before the surrogate is updated. We found that the sample size has no effect on the overall 1240 performance of IBO-HPC. However, for some tasks (JAHS CIFAR-10 and CO), a significant variation 1241 of convergence speed in early iterations – depending on the choice of L – was obtained. Choosing 1242 L = 20 seems to lead to fast and stable convergence behavior.



Figure 6: **IBO-HPC** is competitive or outperforms strong baselines on HPO-B. (a) IBO-HPC outperforms all BO baselines that allow users to provide a prior before optimization when feedback is provided at the 5th iteration. Moreover, IBO-HPC is competitive to other BO methods without any user knowledge given. Reults were obtained on HPO-B with search space ID 6794 and dataset ID 31. (b) IBO-HPC outperforms all BO baselines when no user feedback is provided and beats all interactive BO baselines, except for π BO, when feedback is provided at the 5th iteration. Results are obtained on HPO-B with search space ID 6767 and dataset ID 31.



Figure 7: **IBO-HPC outperforms state of the art.** For 4/5 tasks across three challenging benchmarks, IBO-HPC is competitive with strong baselines when no user knowledge is provided. When beneficial user beliefs (\triangle) are provided, either as distributions (—) or point values (—, —), it outperforms all competitors w.r.t. convergence and solution quality on most tasks. Early interactions (—/— at 5th iteration, — at 10th iteration) speed convergence up.

- 1000
- 1293
- 1294
- 1295



Figure 8: **IBO-HPC recovers from misleading user feedback.** IBO-HPC successfully and consis-1319 tently recovers from misleading user feedback and performs equally well as if no feedback was given. 1320 Also, IBO-HPC handles alternating, contradictory feedback well and is able to leverage beneficial feedback while ignoring misleading feedback.



Figure 9: CDF of Test Accuracy. IBO-HPC samples more good-performing configurations than 1345 most other BO baselines on most tasks. Thus, IBO-HPC invests more computational resources 1346 in good configurations than other methods. We conjecture that this is because IBO-HPC selects 1347 configurations s.t. they are likely to perform similarly to the incumbent in each iteration. Interestingly, 1348 RS also samples many well-performing configurations on the JAHS benchmark. 1349



Figure 11: Conditioning on sub-optimal evaluation scores slow down IBO-HPC- Conditioning on the evaluation score of high-performing configurations is crucial for the performance of IBO-HPC. To analyze the effect of conditioning on evaluation scores of sub-optimal configurations, we conditioned on the $\{0.25, 0.5, 0.75\}$ -quantile of all evaluation scores obtained until iteration t. As expected, for higher quantiles (i.e. better evaluation scores), IBO-HPC finds better configurations.



Figure 12: *L* has no significant effect on IBO-HPC's performance. We found that fixing the surrogate model for $L = \{5, 10, 20, 30\}$ iterations does not lead to significant differences in the performance and convergence speed of IBO-HPC. Only in earlier iterations was a significant variation in convergence speed found on JAHS CIFAR-10 and JAHS CO. However, these variations vanish with the progress of optimization.

1458 D.4 COMPUTATIONAL COST 1459

1460 We now provide details on the computational costs of IBO-HPC. Therefore, we analyzed the composition of the overall runtime of an optimization run and measured the time needed to train configurations 1461 suggested by IBO-HPC versus the time spent on actually performing optimization (including fitting 1462 the surrogate PC and sampling new configurations). Fig. 13a shows that the computation time spent 1463 on learning the PC and sampling new configurations is negligible compared to the time spent on 1464 training the suggested configurations. Additionally, 13b shows that IBO-HPC is faster than SMAC 1465 in 4/5 cases in terms of runtime. Here, we considered the time spent in updating the surrogate and 1466 suggesting new configurations. Note that this does not include training costs. Interestingly, with the 1467 increasing size of the search space, the efficiency advantage of IBO-HPC is increasing. We suspect 1468 that the intensify-mechanism in SMAC, which includes a local search, is the reason for the higher 1469 computational costs of SMAC.



Figure 13: **IBO-HPC is a cost-efficient HPO method.** (a) Learning a surrogate and suggesting 1487 new configurations is negligible in terms of computational costs compared to training the suggested 1488 configurations. We computed the time spent on training configurations (blue) vs. time spent learning 1489 a PC and suggesting new configurations (orange). In all experiments, the training of configurations 1490 caused the large majority of computational costs, often even approaching 100%. (b) IBO-HPC is 1491 more efficient than the prominent HPO algorithm SMAC in 4/5 cases (averaged over 20 runs). Also, 1492 with the increasing number of hyperparameters, the gap between IBO-HPC and SMAC in terms of 1493 computational efficiency is larger. We report runtimes normalized between [0, 1] per benchmark s.t. 1494 the highest obtained runtime for a given benchmark is 1. 1495

EXPLORATION-EXPLOITATION TRADE-OFF OF IBO-HPC D.5 1497

1498 An effective mechanism to trade off exploration versus exploitation is crucial for high-performing 1499 hyperparameter optimization algorithms. Below we show that IBO-HPC's sampling policy effectively 1500 achieves this trade-off. In early iterations, IBO-HPC explores the search space (high sample variance), while in later iterations, it exploits the knowledge collected (low sample variance). 1502

1503

1501

- 1506
- 1507
- 1509
- 1510
- 1511



Figure 15: **IBO-HPC effectively trades off exploration and exploitation.** IBO-HPC's sampling policy naturally and effectively trades off exploration (high sampling variance in early iterations) versus exploitation (low sampling variance in later iterations). We show the sampling variance of 6 hyperparameters of the JAHS benchmark (Colorectal Histology) for each iteration, averaged over 20 runs of IBO-HPC.

1566 D.6 HYPERPARAMETERS OF IBO-HPC

1568 IBO-HPC comes with a few hyperparameters itself, which have to be set. For our experiments, we 1569 set the number of iterations the surrogate is kept fixed L = 20, the decay value $\gamma = 0.9$. We let all 1570 methods optimize for 2000 iterations for fair comparison. Our surrogate models, i.e., PCs and the 1571 associated learning algorithm, have some hyperparameters as well. The structure learning algorithm 1572 splits use the RDC independence test and K-means clustering. The threshold to detect independencies 1573 is set to 0.3, and the minimum number of instances per leaf is adapted dynamically based on the 1574 number of configurations tested during an optimization run.

1575 1576

D.7 HARDWARE

1577 We ran all our experiments on DGX-A100 machines and used 10 CPUs for each run, thus parallelizing 1578 some sub-routines (e.g. learning of PCs). We did not use any GPUs as we queried the benchmarks 1579 employed to provide the performance of configurations. The JAHS benchmark requires a relatively 1580 large RAM (> 16GB) to run smoothly as it loads large ensemble models.

1581 1582

1583

1598

1604

1608

1609 1610

1611

1612

1613

E WORKING EXAMPLE

In the following we consider a more detailed example of our proposed method from a user perspective. We assume that we only optimize 3 hyperparameters here, W, N and R which correspond to the hyperparameters W, N and RESOLUTION in the JAHS benchmark.

The optimization starts where each of the hyperparameters gets optimized by our method. At some point, the user interacts with the optimization process and sets W and N to a fixed value (blue in Fig. 16). From then on, the model only optimizes the remaining hyperparameter R (green in Fig. 16), using conditional sampling from the resulting conditional distribution that the HPC represents after the interaction.



Figure 16: **Example of IBO-HPC.** A user specifies certain aspects of the hyperparameter search space during optimization. Afterward, IBO-HPC takes user knowledge into account when sampling new configuration candidates.

- 1614 1615
- 1616
- 1617
- 1618
- 1619