

# LAMDA: TWO-PHASE MULTI-FIDELITY HPO VIA LEARNING PROMISING REGIONS FROM DATA

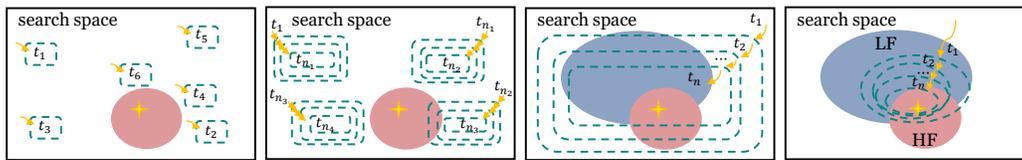
## ABSTRACT

Multi-fidelity hyperparameter optimization (HPO) combines data from both high-fidelity (HF) and low-fidelity (LF) problems during the optimization process, aiding in effective sampling and preliminary screening. To enhance its performance, approaches that incorporate expert knowledge or transfer ability into the HPO algorithm have demonstrated their superiority, while such domain knowledge or abundant data from multiple similar tasks may not always be accessible. Observing that high-quality solutions in HPO exhibit some overlap between high- and low-fidelity problems, we propose a two-phase framework `Lamda` to streamline the multi-fidelity HPO. Specifically, in the first phase, it searches in the LF landscape to identify the promising regions of LF problem. In the second phase, we leverage such promising regions to construct reliable priors to navigate the HPO. We showcase how the `Lamda` framework can be integrated with various HPO algorithms to boost their performance, and further conduct theoretical analysis towards the integrated Bayesian optimization and bandit-based Hyperband. We demonstrate the effectiveness of our framework across 56 HPO tasks.

## 1. INTRODUCTION

The performance of machine learning models is highly dependent on their hyperparameters (Bischl et al., 2023), while hyperparameter optimization (HPO) has become a popular research area in both academia and industry (Li et al., 2022a). In practice, the cost of an HPO task can be prohibitively high when dealing with large models or datasets. For instance, the training time of a specified model on large datasets can take several hours or even days (Krizhevsky et al., 2012). Various HPO methods have been developed, ranging from the well-established random search (RS) (Bergstra & Bengio, 2012) to more data-efficient Bayesian optimization (BO) (Kandasamy et al., 2018; Bergstra et al., 2011; McLeod et al., 2017). Many of these methods find solutions from a uniform global perspective as shown in Figure 1(a). To avoid directionless search with potentially low returns, variants based on localized search strategies such as the trust region Bayesian optimization (TuRBO) (Eriksson et al., 2019) have been proposed with more focused search regions illustrated in Figure 1(b). Nevertheless, all of these methods do not scale satisfactorily with the increasingly complex and costly HPO tasks. In this context, especially with deep models and large-scale datasets, fidelity management becomes more important given limited budget.

Based on the hypothesis that low-fidelity (LF) evaluation reveals a reasonable approximation of the high-fidelity (HF) performance while consuming less budget, multi-fidelity HPO methods employ various techniques to actively manipulate the evaluation fidelity, such as using subsets of dataset, reducing feature space, and decreasing the number of training epochs (Klein et al., 2017; Falkner et al., 2018). The multi-fidelity Bayesian optimization (MFBO) (Swersky et al., 2013) and bandit-based methods (Li et al., 2017) are two representative multi-fidelity HPO methods. For MFBO,



(a) Global search methods (b) Local search methods: TuRBO (c) Bandit based methods (d) Methods using promising area

Figure 1: A conceptual demonstration of how different HPO methods explore the search space. The red and blue areas represents regions of high-quality solutions for an HF and LF problem, respectively, while the yellow stars denote the optimal solution for the HPO task. The dashed lines in panel (a) show the locations of sampled points at each iteration and represent the search space of the sampling function depicted in panels (b) through (d).  $t$  represents the iteration number.

existing work primarily constructs an integrated surrogate model accommodating multi-fidelity evaluations for better acquiring candidate configurations (Poloczek et al., 2017; Kandasamy et al., 2019; Mikkola et al., 2022; Li et al., 2020b). Bandit-based methods, on the other hand, utilize data from LF problems to filter potentially good configurations for HF problems (Falkner et al., 2018; Li et al., 2022b; Awad et al., 2021). As both MFBO and bandit-based methods follow within the Bayesian framework, previous work unintentionally downplayed the role of priors. The random sampling strategy in bandit algorithms and uniform acquisition horizon of MFBO were consistently adopted for all HPO tasks, leading to uninformative priors leaving limited space to prevent performance degradation or further improve efficiency. Their trajectory routine is presented in Figure 1(c).

With the growth of data analysis techniques and the accumulation of more and more in-depth experience in HPO tasks, increasing effort has been put into the heuristics for better guidance of a single HPO task, functionally equivalent to proactively replacing the uninformative priors. The strategic search of recent HPO research has been proposed, either relying on the domain expert knowledge towards the incumbent HPO task, such as Priorband and BO with crafted prior (Souza et al., 2021; Mallik et al., 2023), or requiring the transfer similarity from multiple HPO tasks (Watanabe et al., 2023). As shown in Figure 1(d), these methods guide the search towards prior-determined promising areas to reduce budget consumption. Unfortunately, acquiring the correct expert knowledge for a specific HPO task is not often easy-to-play, and the transfer quality heavily relies on the hypothesis of task similarity and abundant meta sources. Although some work has demonstrated the optimization robustness regarding potentially misleading priors (Hvarfner et al., 2022), the additional cost for crafting the priors and unpredictable budget consumption discouraging practitioners from exhaustively determining a good prior by leveraging knowledge or transfer for their own HPO tasks.

In this paper, we endeavor to design priors for HPO algorithms with competitive heuristics and consistent budget management, without external cost or budget such as the expert cognitive load and other HPO task evaluation. This is achieved by further exploiting the relation between LF and HF landscapes of HPO tasks. We observe that in many HPO scenarios, promising regions containing good LF and HF solutions overlap to some extent (Sections E and F.3). This motivates us to construct a reasonable reliable prior from the LF evaluations. Our idea is orthogonal to that in the multi-fidelity HPO literature for two reasons. Strategically, we aim to identify the promising regions of LF landscape irrespective of the HF performance. Functionally, the identified LF regions will be used as prior for underlining HPO methods, including the multi-fidelity HPO ones. An additional advantage of our design is that budget for specifying prior can be explicitly integrated to the overall budget in multi-fidelity HPO. Table 1 shows the comparison of HPO methods in terms of budget management and search strategies. A preliminary HPO example is presented in Figure 2 considering two HPO

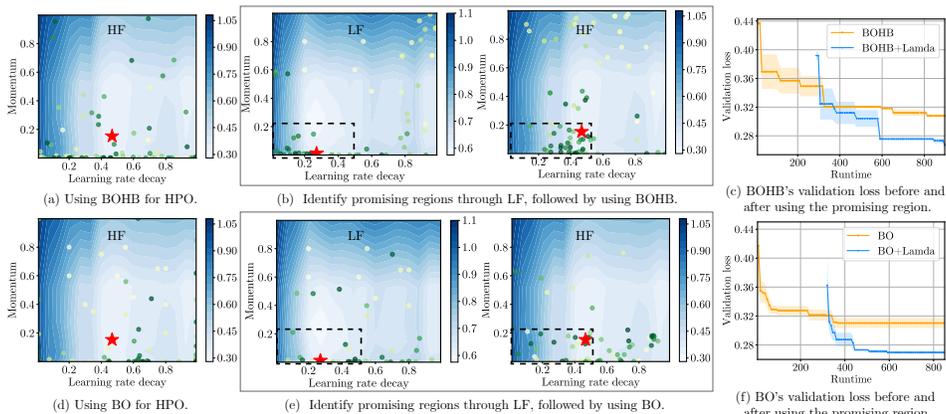


Figure 2: Using BOHB and BO for hyperparameter optimization of a WideResNet on CIFAR-100, before and after employing promising regions (denoted as with or without Lamda). The pentagrams mark the optimal solutions. In (a) and (b), the points represent sampled solutions by BOHB, while in (d) and (e), they represent sampled solutions by BO. The color gradient from yellow to green indicates the progression of sampling over time. The points represent samples from BOHB in (a) and (b), and from BO in (d) and (e), with colors transitioning from yellow to green indicating the progression of sampling over time.

Table 1: Existing methods in addressing the challenges of HPO. Methods that successfully address a challenge are marked with a checkmark (✓), while those that do not are marked with a cross (✗).

Challenges	RS	BO	TuRBO	MFBO	Bandit-based	BO with prior	Transfer search space	Priorband	Ours
Strategic Search	✗	✗	✓	✗	✗	Use expert prior	Use similar tasks	Use expert prior	Use LF
Fidelity management	✗	✗	✗	✓	✓	✗	✗	✓	✓
Consistent budget	✓	✓	✓	✓	✓	✗	✗	✗	✓

methods, BOHB (Falkner et al., 2018) and BO (Bergstra et al., 2011) with noninformative priors and LF-guided priors. In this case, the promising regions for both high- and low-fidelity problems were primarily concentrated within regions bounded by momentum values between  $[0, 0.5]$  and learning rate decay values between  $[0.2, 0.8]$ , as shown in Figure 2(b) and (e). Algorithms with LF-guided prior can **quickly** explore the real promising regions (right panel of Figure 2(b) and (e)). Moreover, in Figure 2(c) and (f), while identifying the LF-guided prior consumes a certain amount of budget, the overall efficiency is significantly improved, which highlights our motivations.

Overall, we propose a two-phase multi-fidelity HPO framework, named **Lamda (Learning promising regions from data)**, which is algorithm-agnostic and serves as a booster for existing HPO algorithms within the Bayesian routine. The contributions are threefold:

- Building on the overlapping promising regions between LF and HF landscapes, we develop a framework that first introduces LF evaluations to identify the promising regions of LF problems, constructing a reasonably reliable prior for underlining HPO algorithms, and then leverages this prior to enhance the HPO algorithms. In addition, an overlapping coefficient is introduced to quantitatively measure the extent of overlapping.
- We integrate the learned prior with various existing HPO algorithms, ranging from prior- and bandit-based methods, as well as multi-fidelity BO, to augment their performance. The rationale of this augmentation was demonstrated by showcasing theoretical analysis towards the prior-based Bayesian optimization and bandit-based Hyperband.
- We validate the competitiveness of our methods across diverse hyperparameter optimization tasks, including fully connected networks, transformers, ResNet, neural architecture search benchmarks, joint architecture and hyperparameter search cases, as well as fine-tuning pretrained image classification models.

## 2. MULTI-FIDELITY HPO BY LEARNING PROMISING REGIONS FROM DATA

The HPO problem is formulated as *minimizing* an expensive-to-evaluate objective function  $f : \mathcal{X} \rightarrow \mathbb{R}$ , where the goal is to find

$$\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}). \quad (1)$$

The configuration  $\mathbf{x}$  is selected from a search space  $\mathcal{X}$  that may include any combination of continuous, discrete, and categorical variables. In the context of HPO,  $f(\mathbf{x})$  represents the training or validation performance of a machine learning model given the hyperparameters defined by  $\mathbf{x}$ . In multi-fidelity HPO,  $f_z(\cdot)$  with  $z \in \{\ell, \ell + 1, \dots, h\}$  is introduced to denote a computation of  $f(\cdot)$  at the fidelity level  $z$ , e.g., the validation loss of a model trained for  $z$  epochs. Define  $f_h$  and  $f_\ell$  with  $\ell < h$  as the HF and LF objectives, respectively.

---

### Algorithm 1: Pseudocode for Lamda

---

In this paper, we propose a multi-fidelity HPO framework by exploiting promising regions from data (dubbed Lamda). It comprises a two-phase search strategy (as depicted in Algorithm 1): **► the first-phase search** initially identifies promising regions in the LF landscape (Lamda-1 in Algorithm 2); **► the second-**

**Input:** Total budget  $\Lambda$ , maximum first-phase budget  $B$ , configuration parameters  $l$ .

$(\varphi_{\text{pro}}(\mathbf{x}), S, \Lambda_l) \leftarrow \text{Lamda-1}(B, l);$

$\Lambda \leftarrow \Lambda - \Lambda_l;$

$x^* \leftarrow \text{Lamda-2}(\varphi_{\text{pro}}(\mathbf{x}), \Lambda_r);$

**return**  $x^*$

---

**phase search** leverages the learned information to guide the search in the HF landscape (Lamda-2 showcased in Appendix G). We will introduce the search strategies in different phases by addressing two key questions.

### 2.1 HOW TO IDENTIFY THE PROMISING REGIONS IN THE LOW-FIDELITY LANDSCAPE?

Our basic idea of the *first-phase search* is to divide the LF landscape into two parts: one consists of the promising regions while the other represents the inferior ones. This can be implemented as a binary classification problem. To train such classifier, we leverage the configurations visited so far during the HPO process in the LF landscape, denoted as  $\mathcal{S} = \{(\mathbf{x}^i, f_\ell(\mathbf{x}^i))\}_{i=1}^t$  where  $f_\ell(\cdot)$

is the LF objective function and  $t$  is the current number of function evaluations. In particular, this paper adopts the classic tree-structured Parzen estimator (TPE) method (Bergstra et al., 2011; Gramacki, 2018) as the classifier, given the scalability and supports for both mixed continuous and discrete spaces. It uses the quantile of  $\{f_\ell(\mathbf{x}) | \mathbf{x} \in \mathcal{S}\}$  to determine the classification boundary. Specifically, we divide  $\mathcal{S}$  into:  $\mathcal{S}_{\text{pro}} = \{\mathbf{x} | f_\ell(\mathbf{x}) \leq y^*, \mathbf{x} \in \mathcal{S}\}$  containing promising solutions, and  $\mathcal{S}_{\text{inf}} = \{\mathbf{x} | f_\ell(\mathbf{x}) > y^*, \mathbf{x} \in \mathcal{S}\}$  containing inferior solutions, where  $y^*$  is determined from  $\alpha = \Pr(f_\ell(\mathbf{x}) < y^*)$  quantile of  $\{f_\ell(\mathbf{x}) | \forall \mathbf{x} \in \mathcal{S}\}$ . Then we denote

$$\varphi_{\text{pro}}(\mathbf{x}) = p(\mathbf{x} | \mathcal{S}_{\text{pro}}), \quad \varphi_{\text{inf}}(\mathbf{x}) = p(\mathbf{x} | \mathcal{S}_{\text{inf}}), \quad (2)$$

where  $\varphi_{\text{pro}}(\mathbf{x})$  is the probability density function (PDF) of the promising solutions, and  $\varphi_{\text{inf}}(\mathbf{x})$  is the PDF of the inferior solutions. We will adopt the kernel density estimation for  $\varphi_{\text{pro}}(\mathbf{x})$  and  $\varphi_{\text{inf}}(\mathbf{x})$ , given its non-parametric nature and applicability to complicated distributions (Chen, 2017).

Instead of searching for the optimal configurations in the LF landscape, the purpose of the *first-phase search* is to identify the promising regions. In practice, the targeted regions are relatively scattered at the beginning and will gradually become focused around the regions that potentially cover the optima (see an illustrative example in Figure 3). Based on this observation, we hypothesize that the *first-phase search* can be terminated when the distribution of promising regions becomes stable. To keep track of the progression of such distribution, we propose to use the overlapping coefficient (OVL) (Anderson et al., 2012) as a metric to quantify the similarity between two distributions.

---

**Algorithm 2:** Pseudocode for Lambda-1
 

---

**Input:** Maximum first-phase budget  $B$ , threshold  $y^*$ ,  $\Delta$ ,  $\gamma$ , fidelity level  $l$ , budget function  $\lambda_z$ .

```

1 Initialize  $S \leftarrow \emptyset$ ,  $\Lambda_l \leftarrow 0$ ,  $t \leftarrow 0$ ,
  isStable  $\leftarrow$  False;
2 while  $\Lambda_l < B$  and not isStable do
3    $\mathbf{x}^t \leftarrow \arg \max_{\mathbf{x} \in \mathcal{X}} \text{AF}(\mathbf{x}, S)$ ;
4    $y^t \leftarrow f_\ell(\mathbf{x}^t) + \epsilon$ ;
5    $S \leftarrow S \cup \{(\mathbf{x}^t, y^t)\}$ ;
6   Update  $\mathcal{S}_{\text{pro}}$ ,  $\mathcal{S}_{\text{inf}}$ ,  $\varphi_{\text{pro}}^t(\mathbf{x})$ ,  $\varphi_{\text{inf}}^t(\mathbf{x})$ ;
7   if  $1 - \rho(\varphi_{\text{pro}}^t(\mathbf{x}), \varphi_{\text{pro}}^{t+\Delta}(\mathbf{x})) \leq \gamma$ 
8     then
9       isStable  $\leftarrow$  True;
10     $\Lambda_l \leftarrow \Lambda_l + \lambda_z(\mathbf{x}^t, l)$ ,  $t \leftarrow t + 1$ ;
11  $\varphi_{\text{pro}}(\mathbf{x}) \leftarrow \varphi_{\text{pro}}^t(\mathbf{x})$ ;
return  $(\varphi_{\text{pro}}(\mathbf{x}), S, \Lambda_l)$ 

```

---

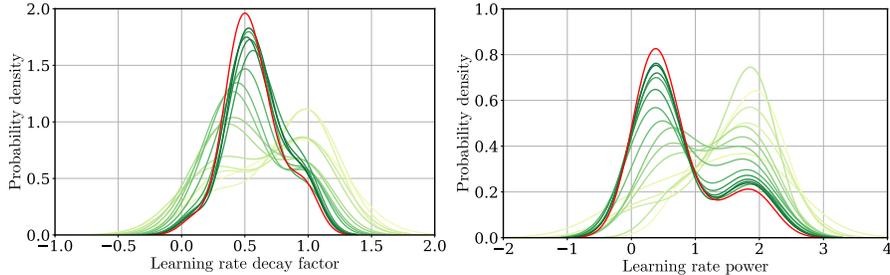


Figure 3: This figure shows the progression of PDFs for promising regions during hyperparameter optimization on a transformer model for the LM1B dataset, focusing on the LF problem. We display PDFs for two out of four hyperparameters, with colors changing from yellow to green to indicate iteration progress. The red line represents the true PDF of the promising solutions in the LF problem.

**Definition 1** (Overlapping coefficient). Let  $\varphi_1(\mathbf{x})$  and  $\varphi_2(\mathbf{x})$  be two PDFs on the search space  $\mathcal{X}$ . The overlapping coefficient  $\rho$  of the two functions is defined as:

$$\rho(\varphi_1(\mathbf{x}), \varphi_2(\mathbf{x})) = \int_{\mathbf{x} \in \mathcal{X}} \min\{\varphi_1(\mathbf{x}), \varphi_2(\mathbf{x})\} d\mathbf{x}. \quad (3)$$

Note that  $\rho(\varphi_1(\mathbf{x}), \varphi_2(\mathbf{x}))$  ranges from 0 to 1, where  $\rho = 1$  if and only if the two distributions are fully overlapped, and  $\rho = 0$  if there is no intersection at all. The *first-phase search* is terminated either if the allocated computational budget is exhausted or the OVL of estimated distributions between  $\Delta \in \mathbb{N}$  iterations is close enough:

$$1 - \rho(\varphi_{\text{pro}}^t(\mathbf{x}), \varphi_{\text{pro}}^{t+\Delta}(\mathbf{x})) \leq \gamma, \quad (4)$$

where  $\gamma$  denotes the threshold. The calculation of  $\rho$  involves a multidimensional integral, which can be numerically intractable. In practice, we employ the Monte Carlo method to estimate  $\rho$  as

$$\begin{aligned} \rho(\varphi_1(\mathbf{x}), \varphi_2(\mathbf{x})) &= \int_{\mathbf{x} \in \mathcal{X}} \min\{\varphi_1(\mathbf{x}), \varphi_2(\mathbf{x})\} d\mathbf{x} = \int_{\mathbf{x} \in \mathcal{X}} \min\left\{1, \frac{\varphi_2(\mathbf{x})}{\varphi_1(\mathbf{x})}\right\} \varphi_1(\mathbf{x}) d\mathbf{x} \\ &= \mathbb{E}\left[\min\left\{1, \frac{\varphi_2(\mathbf{x})}{\varphi_1(\mathbf{x})}\right\}\right] \approx \frac{1}{N} \sum_{i=1}^N \min\left\{1, \frac{\hat{\varphi}_2(\mathbf{x})}{\hat{\varphi}_1(\mathbf{x})}\right\}, \end{aligned} \quad (5)$$

where  $N$  is the number of samples used in the Markov Chain Monte Carlo sampling, and  $\hat{\varphi}(\cdot)$  is an approximation of  $\varphi(\cdot)$  such as using the kernel density estimation.

## 2.2 HOW TO LEVERAGE LF PROMISING REGIONS IN THE HIGH-FIDELITY LANDSCAPE?

With the identified promising regions in the LF landscape, we hypothesize that such information can be used to define the promising regions in the HF landscape. Instead of searching among the entire search space, the *second-phase search* is more focused within the regions defined below:

$$\tilde{\varphi}_{\text{pro}}(\mathbf{x}) = (1 - w) \cdot \varphi(\mathbf{x}) + w \cdot \varphi_{\text{pro}}(\mathbf{x}), \quad (6)$$

where  $\varphi(\mathbf{x})$  is the probability distribution used to guide the HPO process in the HF landscape,  $\varphi_{\text{pro}}(\mathbf{x})$  is the probability distribution of the promising regions identified from the *first-phase search* in the LF landscape, and  $w \in [0, 1]$  with is a hyperparameter that controls the trade-off between the importance of  $\varphi(\mathbf{x})$  learned *on-the-fly* and  $\varphi_{\text{pro}}(\mathbf{x})$  learned in the *first-phase search*. Figure 4 provides a conceptual visualization of leveraging equation (6) during the second-phase search. The redefined promising regions will be closer to the true optimal solution if the promising regions learned in the first phase are closer to the optimum than those before the redefinition, as proven in Proposition 1. Note that since  $\varphi(\mathbf{x})$  is progressively updated during the HPO process with new configurations evaluated and added to the dataset, it is expected that  $\tilde{\varphi}_{\text{pro}}(\mathbf{x})$  will experience a similar trend as  $\varphi_{\text{pro}}(\mathbf{x})$  in the *first-phase search*.

## 2.3 INTEGRATION AND COMPARISON WITH CURRENT MULTI-FIDELITY HPO METHODS

Instead of a standalone algorithm, Lamda plays as a booster that can be integrated with any existing multi-fidelity HPO methods with minor adaptation and thus augmenting the performance of the baseline optimizer. In a nutshell, we only need to replace the sampling strategies of the baseline optimizer with  $\tilde{\varphi}_{\text{pro}}(\mathbf{x})$ . We justify this by comparing with current multi-fidelity HPO methods.

- **Prior-Based Methods:** By using Lamda,  $\tilde{\varphi}_{\text{pro}}(\mathbf{x})$  serves as *a priori* knowledge that represents a reasonable estimation of promising regions in the *second-phase search*. Unlike prior-based methods, which depend on prescribed knowledge or experience from domain experts,  $\tilde{\varphi}_{\text{pro}}(\mathbf{x})$  is adaptively learned from data during the *first-phase search* through the HPO process in the LF landscape. As a result, our approach is resilient to ‘pathological priors’—whether misleading, lacking in informative values, or potentially adversarial—which are not uncommon when tackling new, unseen real-life black-box applications. Additionally, we expect a scenario between our data-driven priors with those elicited from experts can offer consolidated performance enhancement.
- **Bandit-Based Methods:** By using Lamda,  $\tilde{\varphi}_{\text{pro}}(\mathbf{x})$  serves as an effective alternative to the sampling distributions used in bandit-based methods. This restricts the HPO process to focus exclusively on the learned promising regions. In contrast, bandit-based methods begin with LF assessments to identify candidates for HF evaluation and then gradually shift the search focus towards these identified areas. This process, which alternates between exploration and exploitation throughout the entire search space, often leads to inefficient use of computational resources by exploring less promising regions.

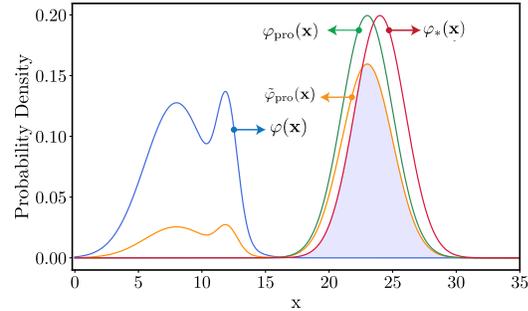


Figure 4: Conceptual visualization of leveraging promising regions:  $\varphi(\mathbf{x})$ ,  $\varphi_{\text{pro}}(\mathbf{x})$ ,  $\tilde{\varphi}_{\text{pro}}(\mathbf{x})$ , and  $\varphi_*(\mathbf{x})$  represent the original sampling distribution, the density function of the promising regions, the modified density function incorporating the promising regions with a weight of  $w = 0.8$ , and the density function of the real optimum.

- **MFBO Methods:** Similar to bandit-based methods, MFBO methods use an acquisition function learned from data collected across multiple fidelities to explore the entire search space. This can lead to unnecessary exploration of less promising regions. By using  $\text{Lamda}$ ,  $\tilde{\varphi}_{\text{pro}}(\mathbf{x})$  restricts the search space within the learned promising regions. Note that this strategy can be applied to the other BO variants.

For proof-of-concept purposes, we choose `PriorBand`, `BOHB`, `MUMBO` as representatives of the prior-, bandit-based and MFBO methods, respectively. By augmenting with `Lamda`, we have `Lamda+PriorBand`, `Lamda+BOHB`, and `Lamda+MUMBO`. In addition, it is also interesting to see whether `Lamda` can be useful for vanilla BO and even random search. To this end, we derive two other variants `Lamda+BO` and `Lamda+RS`. We provide multiple pseudocode to demonstrate how `Lamda-2` can be adapted to various algorithms, including `PriorBand`, `BOHB`, `MUMBO`, `vanilla BO`, and `random search`. The details for each integration are in [Appendix G](#). Furthermore, we provide theoretical analysis under both the prior-based BO framework and bandit-based Hyperband framework, indicating the rational of this augmentation. For `Lamda+BO`, it incorporates prior knowledge in the acquisition function:

$$\mathbf{x}^{n+1} = \arg \max \tilde{\varphi}_{\text{pro}}(\mathbf{x}) \text{AF}(\mathbf{x}, \mathcal{D}), \quad (7)$$

where  $\text{AF}$  is the acquisition function in vanilla BO such as the expected improvement (EI) considered in this paper.  $\mathcal{D} = \{\langle \mathbf{x}^i, f_h(\mathbf{x}^i) \rangle\}_{i=1}^n$  where  $f_h(\cdot)$  is the HF objective function. The Gaussian process regression is employed as the surrogate model of  $f_h(\cdot)$ . For a solution  $\tilde{\mathbf{x}}$ , the predicted mean and variance of the value distribution of  $f_h(\tilde{\mathbf{x}})$  are  $\mu_f(\tilde{\mathbf{x}})$  and  $\sigma_f^2(\tilde{\mathbf{x}})$ .

In `Lamda+BO`, we apply the widely used EI as the acquisition function in [equation \(7\)](#) given as

$$\text{EI}(\tilde{\mathbf{x}}|\mathcal{D}) = \sigma_f(\tilde{\mathbf{x}})(z\Phi_f(z) + \phi_f(z)), \quad (8)$$

where  $z = \frac{f_{\mathcal{D}}^* - \mu_f(\tilde{\mathbf{x}})}{\sigma_f(\tilde{\mathbf{x}})}$ ,  $f_{\mathcal{D}}^* = \min_{\langle \mathbf{x}, f_h(\mathbf{x}) \rangle \in \mathcal{D}} f_h(\mathbf{x})$ ,  $\Phi_f$  and  $\phi_f$  denote the cumulative distribution function and probability density function, respectively.

**Theorem 1.** Given  $\mathcal{D}_n$ ,  $\tilde{\varphi}_{\text{pro}}(\mathbf{x})$ , and applying the EI into [equation \(7\)](#), assume the GP models are non-degenerated. Let  $\mathcal{D}$  be the collected observations with  $\langle \mathbf{x}^1, f_h(\mathbf{x}^1) \rangle$  fixed while  $\{\langle \mathbf{x}^i, f_h(\mathbf{x}^i) \rangle\}_{i=2}^n$  are sequentially chosen by

$$\mathbf{x}^{n+1} = \arg \max \tilde{\varphi}_{\text{pro}}(\mathbf{x}) \text{EI}(\tilde{\mathbf{x}}|\mathcal{D}). \quad (9)$$

Then, as  $n \rightarrow \infty$ , almost surely: the acquisition function converges to zero; and the evaluated best objective  $f_{\mathcal{D}}^* \rightarrow f_h^*$ , where  $f_h^*$  represents the global optimum of  $f_h(\cdot)$ .

The proof is given in [Appendix 11A.2](#). Unlike the theory in (Hvarfner et al., 2022), the convergence property of [equation \(9\)](#) does not need a decaying factor to force exponentially decreasing of the priors of promising LF regions. Based on the dynamic update of  $\tilde{\varphi}_{\text{pro}}(\mathbf{x})$ , it is sufficient to capture the promising regions in the HF landscape according to the overlap.

For bandit-based method, we study the `Lamda+PriorBand` in the theoretical routine of Hyperband. The algorithm involves two loops. In the outer loop, at the  $k$ -th round, the algorithm allocates  $B_{k,s} = 2^k + \text{poly}(k)$  budgets and  $n_{k,s} = 2^s$  configurations randomly sampled from  $\tilde{\varphi}_{\text{pro}}(x)$ , for  $s = 0, 1, \dots, s_{\max}$ , subject to  $s_{\max} + \log_2(s_{\max}) < k$ , where  $\text{poly}(k)$  is some polynomial function w.r.t  $k$ . In the inner loop, the successive halving algorithm is leveraged to find the best arm among the  $n_{k,s}$  arms with  $B_{k,s}$  budget. In the context of multi-arm bandits, each configuration  $\mathbf{x}^i$  corresponds to an independent arm to pull, whose reward with the  $j$ -th pull is denoted by  $l_{i,j} = f_j(\mathbf{x}^i)$ . We assume there exists  $\lim_{k \rightarrow \infty} l_{i,k} = \nu_i$  for all  $\mathbf{x}^i \in \mathcal{X}$ , and denote  $\nu_* = \inf_{\mathbf{x} \in \mathcal{X}} \nu_i$ . Denote also that the distribution of  $v$  as  $F$  satisfying  $P(v - \nu_* \leq \epsilon) = F(\nu_* + \epsilon)$  for any  $\epsilon$ . The inverse function is defined by  $F^{-1}(y) = \inf\{x : F(x) \leq y\}$ . In addition, there exists a monotonically decreasing function  $\gamma(t) : N \rightarrow R$  satisfying  $\sup_i |l_{i,t} - \nu_i| \leq \gamma(t)$ .

**Theorem 2.** For fixed  $\delta \in (0, 1)$ . Let  $\hat{\nu}_B$  be the empirically best-performing arm output from successive halving of round  $k_B = \log_2(B)$  of the outer loop, and let  $s_B < k_B$ . Then, there is:

$$\hat{\nu}_B - \nu_* \leq 3 \left( F^{-1} \left( \frac{\log(4k_B^3/\delta)}{2^{s_B}} - \nu_* \right) + \gamma \left( \frac{2^{k_B-1}}{k_B} \right) \right), \quad (10)$$

where  $s_B$  satisfies  $2^{k_B} + \text{poly}(k_B) > 4s_B \mathbf{H}(F, \gamma, 2^{s_B}, 2k_B^3/\gamma)$  with  $\mathbf{H}(F, \gamma, n, \delta) = 2n \int_{p_n}^1 \gamma^{-1} \left( \frac{F^{-1}(t) - \nu_*}{4} \right) dt + \frac{10}{3} \log(2/\delta) \gamma^{-1} \left( \frac{F^{-1}(p_n) - \nu_*}{4} \right)$  and  $p_n = \frac{\log_2(2/\delta)}{n}$ .

Since all configurations for successive halving tasks are sampled randomly from a probability distribution described by  $\tilde{\varphi}_{\text{pro}}$ , the theoretical results, specifically Corollary 3 in (Li et al., 2017), still hold in this case. Different from Hyperband that relied on non-adaptive grid search exhausting  $c \log_2(B)$  overall budgets with some constant  $c$ , we sample configurations and allocate budget through both grid search and adaptive design based on  $\tilde{\varphi}_{\text{pro}}(\mathbf{x})$ . Theoretically, this requires the same order of budgets as Hyperband. It will be a quite interesting question to ask how the fact of overlapping can help avoiding the grid search of Hyperband, which will be our future work.

### 3. EXPERIMENT SETUP

#### 3.1 BENCHMARK SUITES

Our experiments consider 56 benchmarks that cover various search spaces including mixed types and log-scaled hyperparameters. Further, they involve a wide range of downstream tasks including image classification, language modeling, tabular data processing, medical applications, and translation.

They are selected from four sources. ① Tabular benchmarks include ► four cases from FCNet (Pfisterer et al., 2022), each with 6 hyperparameters; ► one from NAS-Bench-301 with 34 hyperparameters (Pfisterer et al., 2022); ► three from NAS-Bench-201, each with 6 hyperparameters (Eggenberger et al., 2021); and ► twenty benchmarks from rpart on decision tree, glmnet on elastic net, ranger on random forest, and XGBoost (Eggenberger et al., 2021). ② Surrogate benchmarks include ► four problems from PD1 benchmarks with 4 hyperparameters (Mallik et al., 2023; Wang et al., 2021); ► three problems from JAHSBench (Mallik et al., 2023; Bansal et al., 2022) with 14 mixed-type hyperparameters for tuning both the neural networks architecture and training hyperparameters. ③ Training two deep neural networks include LeNet on CIFAR-10, and ResNet-18 on CIFAR-10 and CIFAR-100 with 5 hyperparameters. ④ Two synthetic Hartmann functions (Mallik et al., 2023) with three and six variables respectively. ⑤ 20 tasks for fine-tuning pretrained image classification models (Pineda-Arango et al., 2024).

For the tabular and surrogate benchmarks, we use the number of epochs as the parameter to set the fidelity level. As for the training of deep neural networks, the size of the dataset is used as the parameter to control the fidelity level, as shown in Table 5. Our experiments have considered different scenarios where the promising regions of LF and HF landscape have varying levels of overlaps (Table 5 in Appendix C.4 provides statistics of the overlapping rates). Further detail about all benchmarks are provided in Appendix C.

#### 3.2 PEER ALGORITHMS

We choose nine peer algorithms as the baselines to validate the effectiveness of proposed approach. They are ► PriorBand (Mallik et al., 2023) and PFNS4BO (Müller et al., 2023) as prior-based methods; ► HyperBand (Li et al., 2017), BOHB (Falkner et al., 2018), and Hyper-Tune (Li et al., 2022b) as bandit-based methods; ► MUMBO (Li et al., 2021a) and DPL (Kadra et al., 2023) as MFBO methods; and ► random search (RS), BO and Turbo (Eriksson et al., 2019) as other popular HPO methods. During our experiments, we used the default values for these algorithms. For our algorithm, the parameter settings are as follows:  $\gamma = 0.1$ ,  $\Delta = 5$ ,  $\alpha = 15$  and  $w = 0.5$ . Additionally, we allocate a computational budget of  $B = 5D$  high-fidelity resources in the first-phase search, where  $D$  denotes the problem’s dimensionality.

## 4. EXPERIMENTAL RESULTS

### 4.1 EFFECTIVE OF USING PROMISING REGIONS

**Results on tabular, surrogate and synthetic benchmarks:** This experiment demonstrates how our algorithm framework improves five commonly used optimizers: PriorBand, BOHB, MUMBO, BO and RS. For BOHB, MUMBO, BO and RS, 33 tasks from the above tabular, surrogate and synthetic benchmarks are used. For PriorBand, the evaluation included eight tasks: four from the original paper (PD1-LM1B, PD1-WMT, MFH3, and MFH6) and four additional tasks in FCNet. In addition, good prior are used at PriorBand for PD1-LM1B, PD1-WMT, MFH3, and MFH6. Table 2 presents the numbers of win/lose/tie obtained by using the Wilcoxon signed-rank test and Figure 5 shows the average rank over the HPO tasks. According to Figure 5, it can be seen that using the strategy of promising regions can obtain better results on all the five algorithms. The results of the Wilcoxon signed-rank test in Table 2 also validate the efficiency of using promising regions. Using promising regions can achieve significantly better results than the baseline in more than half of the tasks. For the remaining problems, they obtain results equal to those of the baseline. Regarding PriorBand, for the four tasks with prior information, Lamda+PriorBand achieves results comparable to those of PriorBand. However, for the four tasks without prior information,

Lamda+PriorBand achieves better results. For BOHB, using the strategy of promising regions improves in rank as the consumed resources grow. The main reason for such slow-starting may be caused by the resources needed for finding promising regions. However, it quickly takes the top rank after using about 15 HF resources. MUMBO, BO and RS present similar performances in the condition of using promising regions.

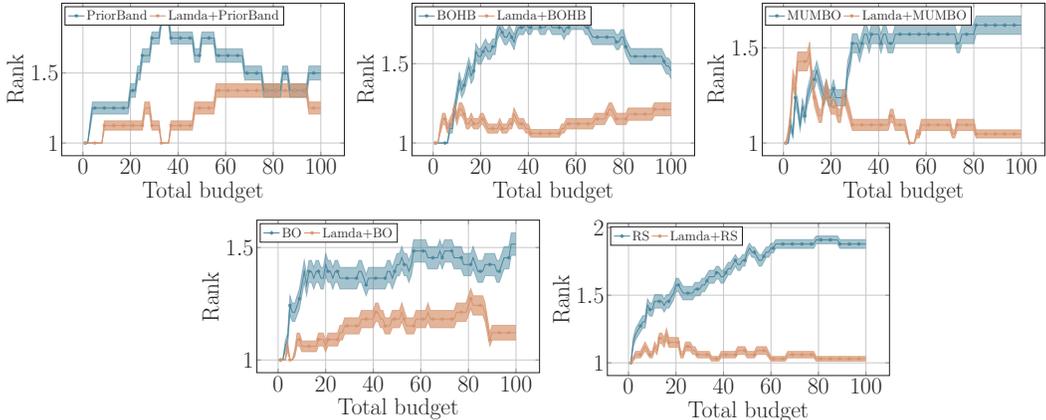


Figure 5: Comparing average relative ranks of PriorBand, BOHB, MUMBO, BO and RS under the proposed framework across 33 HPO tasks.

Table 2: Performance comparison (Win/Lose/Tie) of Lamda+PriorBand, Lamda+BOHB, Lamda+MUMBO, Lamda+BO and Lamda+RS against their baselines over 100 HF evaluations.

Lamda+BOHB	Lamda+MUMBO	Lamda+BO	Lamda+RS	Lamda+PriorBand
24/0/9	19/0/14	17/1/15	29/1/3	3/0/5

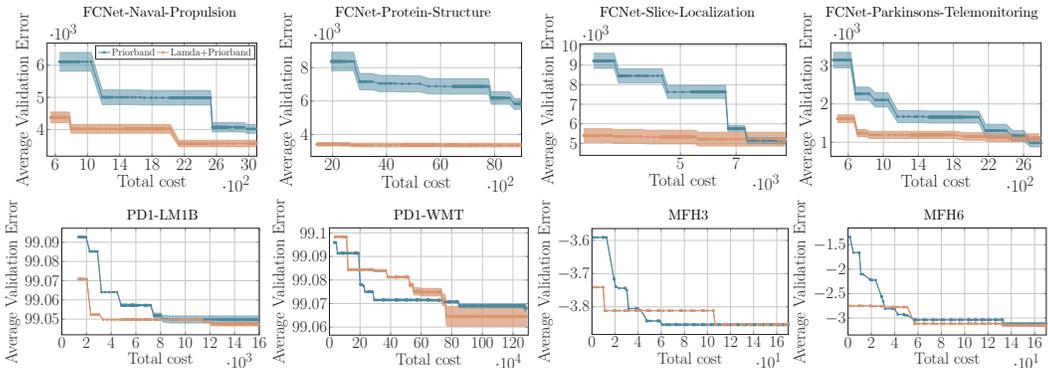


Figure 6: Validation error observed in tuning 8 HPO tasks, using PriorBand as the baseline.

Figure 6 and Figure 18 to Figure 25 shows the performance curves for each benchmark under the framework of Priorband, BOHB, MUMBO, BO and RS. Within the framework of Priorband, Lamda+Priorband converge faster on all the 8 tasks as shown in Figure 6. For other four algorithms, using the promising regions accelerates the discovery of effective solutions compared to the baseline on most of the tasks. In particular, we would like to highlight the results on JAHS-CIFAR-10, JAHS-Colorectal-Histology, and JAHS-Fashion-MNIST, whose overlaps are low. Lamda consistently enhance the baseline algorithms.

To better understand the results, we sample 10,000 hyperparameter configurations for each benchmark and evaluate their performance at high and low fidelities. The configurations are mapped into 2D space. We visualize the good solutions (Figure 14 to Figure 17) and landscape (Figure 33 to Figure 36) of the benchmarks at low and high fidelities. For FCNet, PDI, and NAS-Bench-201, the gap between Lamda+BOHB and the naive BOHB is caused by the great overlapping between the high and low fidelities as shown in Figure 14. Additionally, Figure 12 and Figure 13 shows the found good solutions by the LF. It can be seen that these solutions are close to the good solutions

432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

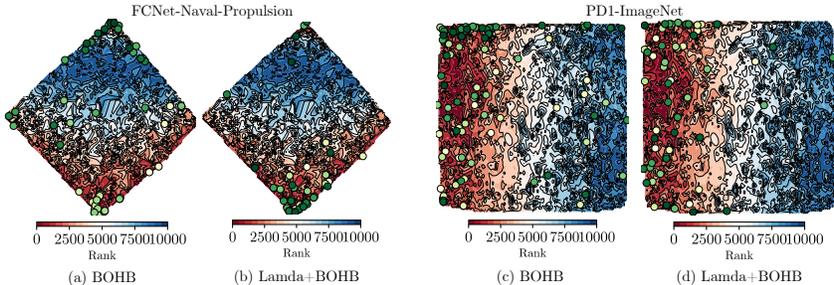


Figure 7: 2D visualization of the sampling points during the optimization process of Lamda+BOHB and BOHB on FCNet-Naval-Propulsion and PD1-ImageNet. The contour represents the dimensionality-reduced landscape of the HF problem, while the points indicate the HF samples collected during optimization. The color gradient from green to yellow indicates the order of sampling, with yellow representing later sampling stages.

of HF. Additionally, we provide a 2D visualization of the sampling points during the optimization process of Lamda+BOHB, BOHB, Lamda+MUMBO, and MUMBO on FCNet-Naval-Propulsion and PD1-ImageNet, as shown in Figure 7 and Figure 37. It can be observed that Lamda-based methods tend to focus sampling in regions with better fitness values, whereas BOHB and MUMBO allocate relatively more resources to exploring areas with moderate fitness values.

**Results on raw problems:** In this part, we evaluate BOHB and Lamda+BOHB on three raw HPO tasks. Figure 8 shows the performance curves on three vision problems. We observe that the performance of Lamda+BOHB is worse than BOHB at the initial iteration. However, it quickly outperforms BOHB after some resources. The main reason is that the promising regions at the high and low fidelity have great overlapping as shown Figure 9.

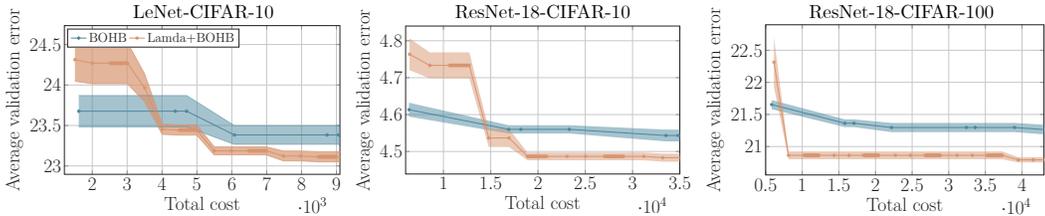


Figure 8: Validation error observed in tuning raw problems.

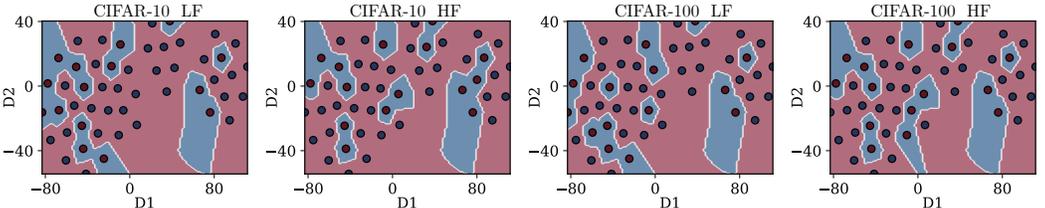


Figure 9: Visualizing the top 30% of solutions (represented by blue regions) in both high and low fidelity while optimizing the hyperparameters of ResNet-18 on CIFAR-10 and CIFAR-100 datasets.

4.2 PEER COMPARISON

**Results on tabular and surrogate benchmarks:** This experiment demonstrates how the performance of the proposed algorithm framework compared to the other methods over the number of evaluations. Figure 10 shows the average rank over the 33 tasks. According to Figure 10, we observe that Lamda+BO and Lamda+MUMBO consistently quickly take the top and keep the rank until the end. It also can be seen that Hyperband and RS are the two worst algorithms, which may be due to the random sampling strategy.

We have conducted 20 HPO tasks for fine-tuning pretrained image classification models. Experimental results (see Appendix F.4) have also demonstrated that Lamda consistently enhance the baseline algorithms.

4.3 PARAMETER ANALYSIS

We investigate the impact of parameters in Lamda within the BOHB framework, including different budgets ( $B$ ) and thresholds ( $\gamma$ ) for stopping the first phase, the interval for calculating the overlapping coefficient ( $\Delta$ ), the quantile ( $\alpha$ ) used in the TPE, and the initial weight ( $w_0$ ). Four tasks at the XGBoost benchmark, each involving a 10-dimensional hyperparameter optimization problem, are used for the experiments. The budgets are varied as  $B = D, 4D, 8D, 10D$ , and  $20D$ , and  $\gamma$  values are tested at 0.5, 0.2, and 0.1. As shown in Figure 28, the choice of  $B$  influences algorithm performance, where too high a budget ( $B > 8D$ ) leads to excessive resource consumption in the first phase. The impact of  $\gamma$  on the algorithm is relatively minor as shown in Figure 29, likely due to the constraints imposed by the maximum budget  $B$ . Additionally, we investigate the influence of the parameter  $\Delta$  on the algorithm’s performance, as illustrated in Figure 30. The settings for  $\Delta$  are 3, 5, 15, 35, and 50. The results indicate that both excessively large and excessively small values of  $\Delta$  slightly affect the algorithm’s performance. Specifically, values of  $\Delta$  that are too low may cause the algorithm to erroneously determine that promising regions have stabilized, whereas excessively high values of  $\Delta$  can lead to considerable delays in the algorithm’s capacity to verify the stability of these regions. We also analyze the influence of parameter  $\alpha$  using the values 5, 15, and 25. The results in Figure 31 indicate no significant impact on the algorithm’s performance.

Regarding the initial weight  $w_0$ , we examine settings of 1, 0.8, 0.5, 0.3, and 0.1. The results, illustrated in Figure 32, suggest that while the optimal  $w_0$  slightly vary across tasks, its overall impact on the algorithm’s performance is minimal. In addition, values above 0.1 consistently outperform the original BOHB. Further analysis of rankings show that settings with  $w > 0.1$  achieve better results compared to  $w = 0.1$ . Notably, a setting of  $w = 1$  shows superior performance during the early stages of the optimization. These results also indicate the efficiency of using the promising regions.

## 5. LIMITATIONS AND FUTURE WORK

In this work, we only consider two fidelities for a proof-of-concept purpose. Our next step is to extend the current Lamda framework for tackling multiple fidelity levels. In addition, there is a gap on theoretical underpinnings about how the involved parameters, such as the quantile threshold in LF problems and the overlapping coefficient between LF and HF landscape, impact the convergence rate or regret of the HBO methods augmented with Lamda. Last but not the least, this paper is mainly designed for multi-fidelity hyperparameter optimization. It will be interesting to explore its applications to a broader range of black-box optimization problems where multi-fidelity experiments and data are prevalent (e.g., computational fluid dynamics optimization in engineering design (Barrett et al., 2006; Liu et al., 2017), and new material (Goldfeld et al., 2005; Khatamsaz et al., 2021), or drug design (Fare et al., 2022; Greenman et al., 2021) in scientific discovery).

## 6. CONCLUSIONS

This paper highlights a common limitation in existing HPO algorithms, which often searches across the whole search space. While some methods leverage prior knowledge to constrain the search space, accessibility to the knowledge is not always guaranteed. To address this challenge, we have developed an algorithmic framework that enables algorithms to autonomously identify promising regions from the LF to accelerate the HPO process, based on the the potential overlap of promising regions between high- and low-fidelity HPO landscape. This framework is integrated with a variety of existing HPO techniques, including prior- and bandit-based methods, as well as multi-fidelity BO, to enhance their efficacy. We support the rationale behind this augmentation through theoretical analysis focused on prior-based Bayesian optimization and bandit-based Hyperband. Our empirical evaluations across diverse hyperparameter optimization tasks—such as fully connected networks, transformers, ResNets, and neural architecture search benchmarks, including joint architecture and hyperparameter searches—demonstrate the competitiveness of our methods.

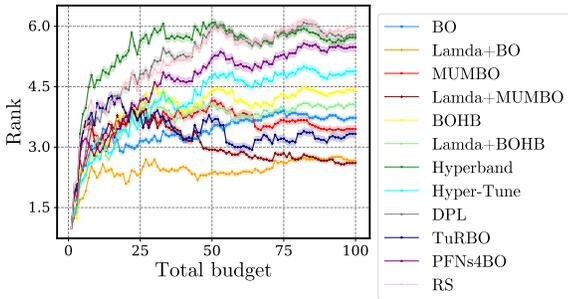


Figure 10: Comparing average relative ranks of peer algorithms across 33 HPO tasks.

## REFERENCES

- 540  
541  
542 Gordon Anderson, Oliver B. Linton, and Yoon-Jae Whang. Nonparametric estimation and inference  
543 about the overlap of two distributions. *J. Financ. Econom.*, 171:1–23, 2012.
- 544 Noor H. Awad, Neeratyoy Mallik, and Frank Hutter. DEHB: evolutionary hyperband for scalable,  
545 robust and efficient hyperparameter optimization. In *IJCAI’21: Proc. of the Thirtieth Interna-*  
546 *tional Joint Conference on Artificial Intelligence*, pp. 2147–2153. ijcai.org, 2021.
- 547 Archit Bansal, Danny Stoll, Maciej Janowski, Arber Zela, and Frank Hutter. JAHS-Bench-201: A  
548 foundation for research on joint architecture and hyperparameter search. In *NeurIPS’22: Ad-*  
549 *vances in Neural Information Processing Systems 35*, 2022.
- 550 Thomas R. Barrett, Neil W. Bressloff, and Andy J. Keane. Airfoil shape design and optimization  
551 using multifidelity analysis and embedded inverse design. *AIAA Journal*, 44:2051–2060, 2006.
- 552 Julien Bect, François Bachoc, and David Ginsbourger. A supermartingale approach to gaussian  
553 process based sequential design of experiments. *Bernoulli*, 25(4A):2883–2919, 2019.
- 554 James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *J. Mach.*  
555 *Learn. Res.*, 13:281–305, 2012.
- 556 James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-parameter  
557 optimization. In *NeurIPS’11: Advances in Neural Information Processing Systems 24*, pp. 2546–  
558 2554, 2011.
- 559 Bernd Bischl, Martin Binder, Michel Lang, Tobias Pielok, Jakob Richter, Stefan Coors, Janek  
560 Thomas, Theresa Ullmann, Marc Becker, Anne-Laure Boulesteix, Difan Deng, and Marius Lin-  
561 dauer. Hyperparameter optimization: Foundations, algorithms, best practices, and open chal-  
562 lenges. *WIREs Data. Mining. Knowl. Discov.*, 13(2), 2023.
- 563 Yen-Chi Chen. A tutorial on kernel density estimation and recent advances. *Biostat. Epidemiol.*, 1:  
564 161 – 187, 2017.
- 565 Katharina Eggenberger, Philipp Müller, Neeratyoy Mallik, Matthias Feurer, René Sass, Aaron  
566 Klein, Noor H. Awad, Marius Lindauer, and Frank Hutter. Hpobench: A collection of repro-  
567 ducible multi-fidelity benchmark problems for HPO. In *NeurIPS’21: Proc. of the Neural Infor-*  
568 *mation Processing Systems Track on Datasets and Benchmarks*, 2021.
- 569 David Eriksson, Michael Pearce, Jacob R. Gardner, Ryan Turner, and Matthias Poloczek. Scalable  
570 global optimization via local bayesian optimization. In *NeurIPS’19: Advances in Neural Infor-*  
571 *mation Processing Systems 32: Annual Conference on Neural Information Processing Systems*  
572 *2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 5497–5508, 2019.
- 573 Stefan Falkner, Aaron Klein, and Frank Hutter. BOHB: robust and efficient hyperparameter opti-  
574 mization at scale. In *ICML’18: Proc. of the 35th International Conference on Machine Learning*,  
575 *, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine*  
576 *Learning Research*, pp. 1436–1445. PMLR, 2018.
- 577 Clyde Fare, Peter Fenner, Matthew Benatan, Alessandro Varsi, and Edward O. Pyzer-Knapp. A  
578 multi-fidelity machine learning approach to high throughput materials screening. *Npj Comput.*  
579 *Mater.*, 8:1–9, 2022.
- 580 Yiska Goldfeld, Koen Vervenne, Johann Arbocz, and Fred van Keulen. Multi-fidelity optimization  
581 of laminated conical shells for buckling. *Struct. Multidiscip. Optim.*, 30:128–141, 2005.
- 582 Artur Gramacki. *Nonparametric kernel density estimation and its computational aspects*, volume 37.  
583 Springer, 2018.
- 584 Kevin Greenman, William H. Green, and Rafael Gómez-Bombarelli. Multi-fidelity prediction of  
585 molecular optical peaks with deep learning. *Chem. Sci.*, 13:1152 – 1162, 2021.
- 586 Carl Hvarfner, Danny Stoll, Artur L. F. Souza, Marius Lindauer, Frank Hutter, and Luigi Nardi.  
587  $\pi$ BO: Augmenting acquisition functions with user beliefs for Bayesian optimization. In *ICLR’22:*  
588 *Proc. of the 10th International Conference on Learning Representations*. OpenReview.net, 2022.

- 594 Arlind Kadra, Maciej Janowski, Martin Wistuba, and Josif Grabocka. Deep power laws for hyperpa-  
595 rameter optimization. *CoRR*, abs/2302.00441, 2023. URL [https://doi.org/10.48550/  
596 arXiv.2302.00441](https://doi.org/10.48550/arXiv.2302.00441).
- 597
- 598 Kirthevasan Kandasamy, Gautam Dasarathy, Jeff G. Schneider, and Barnabás Póczos. Multi-fidelity  
599 Bayesian optimisation with continuous approximations. In *ICML'17: Proc. of the 34th Inter-  
600 national Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Re-  
601 search*, pp. 1799–1808. PMLR, 2017.
- 602 Kirthevasan Kandasamy, Willie Neiswanger, Jeff Schneider, Barnabás Póczos, and Eric P. Xing.  
603 Neural architecture search with Bayesian optimisation and optimal transport. In *NeurIPS'18:  
604 Advances in Neural Information Processing Systems 31*, pp. 2020–2029, 2018.
- 605
- 606 Kirthevasan Kandasamy, Gautam Dasarathy, Junier B. Oliva, Jeff G. Schneider, and Barnabás  
607 Póczos. Multi-fidelity Gaussian process bandit optimisation. *J. Artif. Intell. Res.*, 66:151–196,  
608 2019.
- 609
- 610 Danial Khatamsaz, Abhilash Molkeri, Richard Couperthwaite, Jaylen James, Raymundo Arróyave,  
611 Douglas L. Allaire, and Ankit Srivastava. Efficiently exploiting process-structure-property rela-  
612 tionships in material design by multi-information source fusion. *Acta Materialia*, 2021.
- 613 Aaron Klein, Stefan Falkner, Simon Bartels, Philipp Hennig, and Frank Hutter. Fast Bayesian  
614 optimization of machine learning hyperparameters on large datasets. In *AISTATS'17: Proc. of the  
615 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings  
616 of Machine Learning Research*, pp. 528–536. PMLR, 2017.
- 617
- 618 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep con-  
619 volutional neural networks. In *Advances in Neural Information Processing Systems 25: 26th  
620 Annual Conference on Neural Information Processing Systems 2012. Proceedings of a meeting  
621 held December 3-6, 2012, Lake Tahoe, Nevada, United States*, pp. 1106–1114, 2012.
- 622 Cheng Li, Sunil Gupta, Santu Rana, Vu Nguyen, Antonio Robles-Kelly, and Svetha Venkatesh.  
623 Incorporating expert prior knowledge into experimental design via posterior sampling. *CoRR*,  
624 abs/2002.11256, 2020a. URL <https://arxiv.org/abs/2002.11256>.
- 625
- 626 Lisha Li, Kevin G. Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. Hyper-  
627 band: A novel bandit-based approach to hyperparameter optimization. *J. Mach. Learn. Res.*, 18:  
628 185:1–185:52, 2017.
- 629
- 630 Shibo Li, Wei Xing, Robert M. Kirby, and Shandian Zhe. Multi-fidelity Bayesian optimization via  
631 deep neural networks. In *NeurIPS'20: Proc. of the Annual Conference on Neural Information  
632 Processing Systems 2020*, 2020b.
- 633
- 634 Shibo Li, Robert M. Kirby, and Shandian Zhe. Batch multi-fidelity Bayesian optimization with deep  
635 auto-regressive networks. In *NeurIPS'21: Proc. of the Annual Conference on Neural Information  
636 Processing Systems 2021*, pp. 25463–25475, 2021a.
- 636
- 637 Yang Li, Yu Shen, Jiawei Jiang, Jinyang Gao, Ce Zhang, and Bin Cui. MFES-HB: efficient hyper-  
638 band with multi-fidelity quality measurements. In *AAAI'21: Proc. of the AAAI Conference on  
639 Artificial Intelligence*, pp. 8491–8500. AAAI Press, 2021b.
- 640
- 641 Yang Li, Yu Shen, Huaijun Jiang, Tianyi Bai, Wentao Zhang, Ce Zhang, and Bin Cui. Transfer learn-  
642 ing based search space design for hyperparameter tuning. In *KDD '22: The 28th ACM SIGKDD  
643 Conference on Knowledge Discovery and Data Mining, 2022*, pp. 967–977. ACM, 2022a.
- 644
- 645 Yang Li, Yu Shen, Huaijun Jiang, Wentao Zhang, Jixiang Li, Ji Liu, Ce Zhang, and Bin Cui. Hyper-  
646 Tune: Towards efficient hyper-parameter tuning at scale. *Proc. VLDB Endow.*, 15(6):1256–1265,  
647 2022b.
- 648
- 649 Bo Liu, Slawomir Koziel, and Nazar T. Ali. SADEA-II: A generalized method for efficient global  
650 optimization of antenna design. *J. Comput. Des. Eng.*, 4(2):86–97, 2017.

- 648 Neeratyoy Mallik, Edward Bergman, Carl Hvarfner, Danny Stoll, Maciej Janowski, Marius Lin-  
649 dauer, Luigi Nardi, and Frank Hutter. Priorband: Practical hyperparameter optimization in the  
650 age of deep learning. *CoRR*, abs/2306.12370, 2023. doi: 10.48550/ARXIV.2306.12370. URL  
651 <https://doi.org/10.48550/arXiv.2306.12370>.
- 652 Leland McInnes and John Healy. UMAP: uniform manifold approximation and projection for di-  
653 mension reduction. *CoRR*, abs/1802.03426, 2018. URL [http://arxiv.org/abs/1802.](http://arxiv.org/abs/1802.03426)  
654 [03426](http://arxiv.org/abs/1802.03426).
- 655 Mark McLeod, Michael A Osborne, and Stephen J Roberts. Practical bayesian optimization for  
656 variable cost objectives. *arXiv preprint arXiv:1703.04335*, 2017.
- 657 Petrus Mikkola, Julien Martinelli, Louis Filstroff, and Samuel Kaski. Multi-fidelity Bayesian op-  
658 timization with unreliable information sources. *CoRR*, abs/2210.13937, 2022. doi: 10.48550/  
659 arXiv.2210.13937. URL <https://doi.org/10.48550/arXiv.2210.13937>.
- 660 Henry B. Moss, David S. Leslie, and Paul Rayson. MUMBO: multi-task max-value Bayesian opti-  
661 mization. In *ECML/PKDD’20: Proc. of the 2020 European Conference on Machine Learning*  
662 *and Knowledge Discovery in Databases*, volume 12459 of *Lecture Notes in Computer Science*,  
663 pp. 447–462. Springer, 2020.
- 664 Samuel Müller, Matthias Feurer, Noah Hollmann, and Frank Hutter. Pfns4bo: In-context learning  
665 for bayesian optimization. In *International Conference on Machine Learning, ICML 2023, 23-29*  
666 *July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*,  
667 pp. 25444–25470. PMLR, 2023.
- 670 Valerio Perrone and Huibin Shen. Learning search spaces for bayesian optimization: Another view  
671 of hyperparameter transfer learning. In *NeurIPS’19: Advances in Neural Information Processing*  
672 *Systems 32*, pp. 12751–12761, 2019.
- 673 Florian Pfisterer, Lennart Schneider, Julia Moosbauer, Martin Binder, and Bernd Bischl. YAHPO  
674 gym - an efficient multi-objective multi-fidelity benchmark for hyperparameter optimization. In  
675 *AutoML’22: Proc. of 2022 International Conference on Automated Machine Learning*, volume  
676 188 of *Proceedings of Machine Learning Research*, pp. 3/1–39. PMLR, 2022.
- 677 Sebastian Pineda-Arango, Fabio Ferreira, Arlind Kadra, Frank Hutter, and Josif Grabocka. Quick-  
678 tune: Quickly learning which pretrained model to finetune and how. In *The Twelfth International*  
679 *Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenRe-  
680 view.net, 2024.
- 681 Matthias Poloczek, Jialei Wang, and Peter I. Frazier. Multi-information source optimization. In  
682 *NeurIPS’17: Advances in Neural Information Processing Systems 30*, pp. 4288–4298, 2017.
- 683 Anil Ramachandran, Sunil Gupta, Santu Rana, Cheng Li, and Svetha Venkatesh. Incorporating  
684 expert prior in Bayesian optimisation via space warping. *Knowl. Based Syst.*, 195:105663, 2020.
- 685 Artur L. F. Souza, Luigi Nardi, Leonardo B. Oliveira, Kunle Olukotun, Marius Lindauer, and Frank  
686 Hutter. Bayesian optimization with a prior for the optimum. In *PKDD’21: Machine Learning*  
687 *and Knowledge Discovery in Databases. Research Track- European Conference*, volume 12977  
688 of *Lecture Notes in Computer Science*, pp. 265–296. Springer, 2021.
- 689 Kevin Swersky, Jasper Snoek, and Ryan Prescott Adams. Multi-task Bayesian optimization. In  
690 *NeurIPS’13: Advances in Neural Information Processing Systems 26*, pp. 2004–2012, 2013.
- 691 Shion Takeno, Hitoshi Fukuoka, Yuhki Tsukada, Toshiyuki Koyama, Motoki Shiga, Ichiro Takeuchi,  
692 and Masayuki Karasuyama. Multi-fidelity Bayesian optimization with max-value entropy search  
693 and its parallelization. In *ICML’20: Proc. of the 37th International Conference on Machine*  
694 *Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 9334–9345. PMLR,  
695 2020.
- 696 Z. Wang, George E. Dahl, Kevin Swersky, Chansoo Lee, Zelda E. Mariet, Zachary Nado, Justin  
697 Gilmer, Jasper Snoek, and Zoubin Ghahramani. Pre-trained Gaussian processes for Bayesian  
698 optimization. *arXiv preprint arXiv:2109.08215*, 2021.

702 Shuhei Watanabe, Noor H. Awad, Masaki Onishi, and Frank Hutter. Speeding up multi-objective  
703 hyperparameter optimization by task similarity-based meta-learning for the tree-structured parzen  
704 estimator. In *Proceedings of the Thirty-Second International Joint Conference on Artificial In-*  
705 *telligence, IJCAI 2023, 19th-25th August 2023, Macao, SAR, China*, pp. 4380–4388. ijcai.org,  
706 2023.

707 Martin Wistuba, Nicolas Schilling, and Lars Schmidt-Thieme. Hyperparameter search space prun-  
708 ing - A new component for sequential model-based hyperparameter optimization. In *ECML*  
709 *PKDD'15: Machine Learning and Knowledge Discovery in Databases - European Conference*,  
710 volume 9285 of *Lecture Notes in Computer Science*, pp. 104–119. Springer, 2015.

711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A. THEORETICAL ANALYSIS

**Proposition 1.** Assume the OVL between  $\varphi(\mathbf{x})$  and the PDF of the true promising solutions  $\varphi_*(\mathbf{x})$  is less than the overlapping between the PDF of promising region  $\varphi_{\text{pro}}(\mathbf{x})$  and  $\varphi_*(\mathbf{x})$ . Then, the modified sampling function  $\tilde{\varphi}_{\text{pro}}(\mathbf{x}) = w_1 \cdot \varphi(\mathbf{x}) + w_2 \cdot \varphi_{\text{pro}}(\mathbf{x})$ , where  $w_1 + w_2 = 1$ , will have a greater or equal overlapping with  $\varphi_*(\mathbf{x})$  compared to the overlapping of  $\varphi(\mathbf{x})$  with  $\varphi_*(\mathbf{x})$ .

Proposition 1 suggests that incorporating the promising regions into the sampling distribution enhances its alignment with the distribution of the real optimal solutions.

*Proof.* Using the definition of  $\rho(\varphi_1, \varphi_2)$  in equation (3), the OVL between  $\varphi(\mathbf{x})$ ,  $\varphi_{\text{pro}}(\mathbf{x})$ ,  $\tilde{\varphi}_{\text{pro}}(\mathbf{x})$  and  $\varphi_*(\mathbf{x})$  are computed as follows:

$$\begin{aligned}\rho(\varphi(\mathbf{x}), \varphi_*(\mathbf{x})) &= 1 - \frac{1}{2} \int_{\mathbf{x} \in \mathcal{X}} |\varphi(\mathbf{x}) - \varphi_*(\mathbf{x})| d\mathbf{x}, \\ \rho(\varphi_{\text{pro}}(\mathbf{x}), \varphi_*(\mathbf{x})) &= 1 - \frac{1}{2} \int_{\mathbf{x} \in \mathcal{X}} |\varphi_{\text{pro}}(\mathbf{x}) - \varphi_*(\mathbf{x})| d\mathbf{x}, \\ \rho(\tilde{\varphi}_{\text{pro}}(\mathbf{x}), \varphi_*(\mathbf{x})) &= 1 - \frac{1}{2} \int_{\mathbf{x} \in \mathcal{X}} |\tilde{\varphi}_{\text{pro}}(\mathbf{x}) - \varphi_*(\mathbf{x})| d\mathbf{x}.\end{aligned}\tag{11}$$

where

$$\rho(\varphi_1(\mathbf{x}), \varphi_2(\mathbf{x})) = \int_{\mathbf{x} \in \mathcal{X}} \min\{\varphi_1(\mathbf{x}), \varphi_2(\mathbf{x})\} d\mathbf{x} = 1 - \frac{1}{2} \int_{\mathbf{x} \in \mathcal{X}} |\varphi_1(\mathbf{x}) - \varphi_2(\mathbf{x})| d\mathbf{x}.\tag{12}$$

Given that thus  $\varphi(\mathbf{x})$  has a smaller overlap with  $\varphi_*(\mathbf{x})$  than  $\varphi_{\text{pro}}(\mathbf{x})$ , it follows that:

$$\begin{aligned}\rho(\varphi(\mathbf{x}), \varphi_*(\mathbf{x})) - \rho(\varphi_{\text{pro}}(\mathbf{x}), \varphi_*(\mathbf{x})) &= 1 - \frac{1}{2} \int_{\mathbf{x} \in \mathcal{X}} |\varphi(\mathbf{x}) - \varphi_*(\mathbf{x})| d\mathbf{x} - \left(1 - \frac{1}{2} \int_{\mathbf{x} \in \mathcal{X}} |\varphi_{\text{pro}}(\mathbf{x}) - \varphi_*(\mathbf{x})| d\mathbf{x}\right) \\ &= \frac{1}{2} \int_{\mathbf{x} \in \mathcal{X}} |\varphi_{\text{pro}}(\mathbf{x}) - \varphi_*(\mathbf{x})| d\mathbf{x} - \frac{1}{2} \int_{\mathbf{x} \in \mathcal{X}} |\varphi(\mathbf{x}) - \varphi_*(\mathbf{x})| d\mathbf{x} \\ &= \frac{1}{2} \int_{\mathbf{x} \in \mathcal{X}} |\varphi_{\text{pro}}(\mathbf{x}) - \varphi_*(\mathbf{x})| - |\varphi(\mathbf{x}) - \varphi_*(\mathbf{x})| d\mathbf{x} \leq 0\end{aligned}\tag{13}$$

Given the above, if  $\rho(\varphi(\mathbf{x}), \varphi_*(\mathbf{x})) < \rho(\tilde{\varphi}_{\text{pro}}(\mathbf{x}), \varphi_*(\mathbf{x}))$ , it implies that  $\varphi(\mathbf{x})$  has a smaller overlap with  $\varphi_*(\mathbf{x})$  than  $\tilde{\varphi}_{\text{pro}}(\mathbf{x})$ . This can be further analyzed as:

$$\begin{aligned}\rho(\varphi(\mathbf{x}), \varphi_*(\mathbf{x})) - \rho(\tilde{\varphi}_{\text{pro}}(\mathbf{x}), \varphi_*(\mathbf{x})) &= 1 - \frac{1}{2} \int_{\mathbf{x} \in \mathcal{X}} |\varphi(\mathbf{x}) - \varphi_*(\mathbf{x})| d\mathbf{x} - \left(1 - \frac{1}{2} \int_{\mathbf{x} \in \mathcal{X}} |w_1 \cdot \varphi(\mathbf{x}) + w_2 \cdot \varphi_{\text{pro}}(\mathbf{x}) - \varphi_*(\mathbf{x})| d\mathbf{x}\right) \\ &= \frac{1}{2} \int_{\mathbf{x} \in \mathcal{X}} |w_1 \cdot \varphi(\mathbf{x}) + w_2 \cdot \varphi_{\text{pro}}(\mathbf{x}) - \varphi_*(\mathbf{x})| d\mathbf{x} - \frac{1}{2} \int_{\mathbf{x} \in \mathcal{X}} |\varphi(\mathbf{x}) - \varphi_*(\mathbf{x})| d\mathbf{x} \\ &\leq \frac{1}{2} \int_{\mathbf{x} \in \mathcal{X}} w_1 \cdot |\varphi(\mathbf{x}) - \varphi_*(\mathbf{x})| + w_2 \cdot |\varphi_{\text{pro}}(\mathbf{x}) - \varphi_*(\mathbf{x})| - |\varphi(\mathbf{x}) - \varphi_*(\mathbf{x})| d\mathbf{x} \\ &= \frac{1}{2} \int_{\mathbf{x} \in \mathcal{X}} w_2 \cdot |\varphi_{\text{pro}}(\mathbf{x}) - \varphi_*(\mathbf{x})| - w_2 \cdot |\varphi(\mathbf{x}) - \varphi_*(\mathbf{x})| d\mathbf{x} \\ &= \frac{w_2}{2} \int_{\mathbf{x} \in \mathcal{X}} |\varphi_{\text{pro}}(\mathbf{x}) - \varphi_*(\mathbf{x})| - |\varphi(\mathbf{x}) - \varphi_*(\mathbf{x})| d\mathbf{x} \leq 0\end{aligned}\tag{14}$$

where the inequality is obtained with the following formula:

$$\begin{aligned}|w_1 \cdot \varphi(\mathbf{x}) + w_2 \cdot \varphi_{\text{pro}}(\mathbf{x}) - \varphi_*(\mathbf{x})| &= |w_1 \cdot \varphi(\mathbf{x}) - w_1 \cdot \varphi_*(\mathbf{x}) + w_2 \cdot \varphi_{\text{pro}}(\mathbf{x}) - w_2 \cdot \varphi_*(\mathbf{x})| \\ &\leq w_1 \cdot |\varphi(\mathbf{x}) - \varphi_*(\mathbf{x})| + w_2 \cdot |\varphi_{\text{pro}}(\mathbf{x}) - \varphi_*(\mathbf{x})|\end{aligned}\tag{15}$$

This implies that  $\tilde{\varphi}_{\text{pro}}(\mathbf{x})$  has a greater overlap with  $\varphi_*(\mathbf{x})$  than  $\varphi(\mathbf{x})$  does.  $\square$

## A.2 PROOFS OF THEOREM 1

In the context of sequential design, let  $\mathcal{F}_N$  denote the  $\sigma$ -algebra generated by the random variables  $\mathbf{x}^1, Z^1, \dots, \mathbf{x}^N, Z^N$  where  $Z^i$  is the observation of  $f_i(\mathbf{x}^i)$ . Additionally, let  $\mathcal{F}_{N, \tilde{\mathbf{x}}}$  be the  $\sigma$ -algebra generated by  $\mathbf{x}^1, Z^1, \dots, \mathbf{x}^N, Z^N, \tilde{\mathbf{x}}, \tilde{Z}$  with  $\tilde{Z}$  the observation of  $f_h(\tilde{\mathbf{x}})$ . Then, the EI-based sequential design of Lamda+BO takes the following form:

$$\mathbf{x}_{N+1} = \arg \max_{\tilde{\mathbf{x}} \in \Omega} \mathbb{E}_N \left[ M_N^0 - M_{N, \tilde{\mathbf{x}}}^\rho \right], \quad (16)$$

in which  $y_h^*$  is the threshold of promising solutions in HF problems lower bounded by the  $\alpha$  quantile of  $\{f_\ell(\mathbf{x}) | \forall \mathbf{x} \in \mathcal{S}\}$  due to the overlapping between HF and LF landscape, and

$$M_{N, \tilde{\mathbf{x}}}^\rho = \min_{\mathbf{x} \in \mathcal{X}, \mathbb{P}(f_h(\mathbf{x}) < y_h^* | \mathcal{F}_{N, \tilde{\mathbf{x}}}) = 1, \sigma_f(\mathbf{x} | \mathcal{F}_{N, \tilde{\mathbf{x}}}) = 0} \tilde{f}(\mathbf{x}), \quad (17)$$

$$M_N^0 = \min_{\mathbf{x} \in \mathcal{X}, \mathbb{P}(f_h(\mathbf{x}) < y_h^* | \mathcal{F}) = 1, \sigma_f(\mathbf{x} | \mathcal{F}_N) = 0} \tilde{f}(\mathbf{x}). \quad (18)$$

When GPs are non-degenerate, i.e.,  $\sigma_f = 0$  only if  $\mathbf{x} \in \mathcal{D}$  equation (16) becomes equivalent to equation (9). Specifically,  $M_N^0$  will be equal to  $f_{\mathcal{D}}^*$  in equation (8), while  $M_{N, \tilde{\mathbf{x}}}^\rho$  will be the predicted promising HF objective value under threshold  $y_h^*$  implicitly defined by equation (6). Next, we present the criteria for *asymptotic convergence* of Lamda+BO. The proof of the first statement, i.e., the convergence of the acquisition function, comprises three steps.

**Step 1. Lamda+BO serves as a stepwise uncertainty reduction (SUR) sequential design.** For  $N \geq 2$ , a minimization version of equation (16) can be given as

$$\mathbf{x}_{N+1} = \arg \min_{\tilde{\mathbf{x}} \in \mathcal{X}} \mathbb{E}_N [H_{N, \tilde{\mathbf{x}}}], \quad (19)$$

in which

$$H_{N, \tilde{\mathbf{x}}} = M_{N, \tilde{\mathbf{x}}}^\rho - M_N^0 = \mathbb{E}_{N, \tilde{\mathbf{x}}} \left[ M_{N, \tilde{\mathbf{x}}}^\rho - \min_{\mathbf{x} \in \mathcal{X}, \mathbb{P}(f_h(\mathbf{x}) < y_h^*)} \tilde{f}(\mathbf{x}) \right].$$

The above equation holds since: *i*)  $M_N^0$  is independent from  $\tilde{\mathbf{x}}$ , and *ii*)  $\mathbb{E}_{N, \tilde{\mathbf{x}}} [M_{N, \tilde{\mathbf{x}}}^\rho] = M_{N, \tilde{\mathbf{x}}}^\rho$  for minimum operation. Therefore this strategy can be transformed into an equivalent SUR sequential design strategy for  $H_{N, \tilde{\mathbf{x}}}$ . Likewise, we define

$$H_N = \mathbb{E}_N \left[ M_N^\rho - \min_{\mathbf{x} \in \mathcal{X}, f_h(\mathbf{x}) < y_h^*} \tilde{f}(\mathbf{x}) \right]. \quad (20)$$

**Step 2. ( $H_N$ ) is a supermartingale.** For well-structured GP models and well-defined smooth functions  $\rho^i$ , we have: *i*)  $\sigma_f(\mathbf{x} | \mathcal{F}_{N+1}) \leq \sigma_f(\mathbf{x} | \mathcal{F}_N)$  (based on the definition of GP predicted variance), and *ii*)  $\mathbb{P}(f_h(\mathbf{x}) < y_h^*(\mathbf{x}) | \mathcal{F}_N) = 1$  is sufficient for  $\mathbb{P}(f_h(\mathbf{x}) < y_h^*(\mathbf{x}) | \mathcal{F}_{N+1}) = 1$  based on the non-increasing property of density estimation on an evaluated solution  $\mathbf{x}^N$ . Therefore, the following inequality holds:

$$H_N - \mathbb{E}_N [H_{N+1}] = \mathbb{E}_N [M_N^\rho - M_{N+1}^\rho] \geq 0, \quad (21)$$

which implies that  $(H_N)_{N \in \mathbb{N}}$  is a supermartingale. Consequently, there is  $H_N - \mathbb{E}_N [H_{N+1}] \rightarrow 0$  as  $N \rightarrow \infty$ , and also

$$\sup_{\tilde{\mathbf{x}} \in \mathcal{X}} [H_N - \mathbb{E}_N [H_{N, \tilde{\mathbf{x}}}] ] \rightarrow 0. \quad (22)$$

**Step 3. The acquisition function of Lamda+BO converges to zero almost surely.** Due to the lower bound, according to equation (21), as evaluations of HF objective increase,  $\tilde{\varphi}_{\text{pro}}$  tends to converge to  $\varphi$ . Note that for  $N \rightarrow \infty$ ,  $M_N^\rho = M_N^0$  as this convergence appears. Additionally, we also have

$$\sup_{\tilde{\mathbf{x}} \in \mathcal{X}} \mathbb{E}_N [M_N^\rho - M_{N, \tilde{\mathbf{x}}}^\rho] \geq \sup_{\tilde{\mathbf{x}} \in \mathcal{X}} \mathbb{E}_N [M_N^0 - M_{N, \tilde{\mathbf{x}}}^\rho] \geq \sup_{\tilde{\mathbf{x}} \in \mathcal{X}} \tilde{\varphi}_{\text{pro}}(\mathbf{x}) \text{EI}(\tilde{\mathbf{x}} | \mathcal{D}). \quad (23)$$

Therefore, with the same proof as that of Proposition 2.9 (Bect et al., 2019), for  $N \rightarrow \infty$ , equation 22 and equation 23 yield  $\tilde{\varphi}_{\text{pro}}(\mathbf{x}) \text{EI}(\tilde{\mathbf{x}} | \mathcal{D}) \rightarrow 0$ . This completes the proof for the first statement.

864 The second statement stands according to the global search ability of EI and corresponding  
 865 dense evaluated solutions in  $\mathcal{X}$ . We complete the proof by providing the following facts: *i*)  
 866  $\tilde{\varphi}_{\text{pro}}(\mathbf{x})\sigma_f(\tilde{\mathbf{x}}) \rightarrow 0$  holds from the first statement; *ii*) the lower bound  $y_h^*$  will be tight when  
 867  $N \rightarrow \infty$ ; and *iii*)  $\sigma_f(z|\mathcal{F}_N) \rightarrow 0$  for all sequences accordingly. Based on these facts, the se-  
 868 quence is almost surely dense in  $\mathcal{X}$ . As a result,  $f_{\mathcal{D}}^*$  from any sequence converges to  $f_{\mathcal{X}}^*$  almost  
 869 surely when  $N \rightarrow \infty$ .

## 871 B. RELATED WORKS

### 873 B.1 MULTI-FIDELITY BAYESIAN OPTIMIZATION

875 In the realm of MFBO, previous research primarily leverages LF to construct an accurate MF model  
 876 for guiding the sampling process (Swersky et al., 2013; Poloczek et al., 2017; Kandasamy et al.,  
 877 2019; Mikkola et al., 2022; Li et al., 2020b). One challenge in these methods is to model perfor-  
 878 mance using data from various fidelities. Solutions include employing Gaussian process regression  
 879 (GPR) models tailored to each fidelity level (Kandasamy et al., 2019), multi-task GPR models for  
 880 discrete fidelity levels (Swersky et al., 2013; Poloczek et al., 2017; Mikkola et al., 2022; Li et al.,  
 881 2020b), and GPR models for a continuous fidelity space (Klein et al., 2017; Kandasamy et al., 2017).  
 882 While Gaussian processes are commonly used for modelling the surrogate function, Li et al. (Li  
 883 et al., 2020b; 2021a) have implemented deep neural networks to represent the relationships across  
 884 different fidelities. In terms of sampling, entropy search methods (Swersky et al., 2013; Poloczek  
 885 et al., 2017; Takeno et al., 2020) are utilized, which take into account both the information gain and  
 886 the associated costs of each fidelity level. These techniques enable the sampling of new solutions at  
 887 either low or high fidelity levels, gradually leading to improved hyperparameters in the HF space.

### 888 B.2 BANDIT BASED METHOD

890 Bandit-based methods employ LF for identifying promising solutions for HF evaluations (Falkner  
 891 et al., 2018; Li et al., 2022b; Awad et al., 2021). In pioneering work in the field, Li et al. (2017)  
 892 introduced Hyperband, a method that improves upon the Successive Halving (SH) algorithm by in-  
 893 tegrating various early stopping strategies across multiple SH brackets. Each bracket starts with a  
 894 different number of solutions at varied fidelity levels. However, Hyperband’s approach of randomly  
 895 sampling solutions doesn’t utilize previous sample information (Falkner et al., 2018). To enhance  
 896 this, Falkner et al. (2018) developed BOHB, which combines BO with Hyperband, using the tree  
 897 Parzen estimator (TPE) (Bergstra et al., 2011) to build surrogate models for each fidelity level. While  
 898 BOHB is effective, it struggles with discrete dimensions and scaling to high-dimensional problems.  
 899 Addressing these challenges, Awad et al. (2021) enhanced BOHB with differential evolution for  
 900 better candidate sampling. Additionally, Li et al. (2021b) developed an MF ensemble model that  
 901 integrates information from all fidelity levels to more accurately estimate the highest fidelity. Nev-  
 902 ertheless, despite their utilization of LF data, these approaches still demand extensive exploration  
 903 throughout the entire search space.

### 904 B.3 USING PRIORS FOR OPTIMIZATION

906 These methods concentrate on promising regions to reduce the consumed resources by leveraging  
 907 prior information. A line of work uses prior information about locations of optimal solutions to ac-  
 908 celerate the process of HPO. For instance to reduce evaluations on bad regions, Souza et al. (2021)  
 909 injected priors about which parts of the input space will yield the best performance into BO’s stan-  
 910 dard probabilistic model to form a pseudo-posterior, which was shown to be more sample-efficient  
 911 than BO baselines. Further, Li et al. (2020a) incorporated Gaussian distributions of optimal solu-  
 912 tions into the posterior distribution of observed data and used Thompson sampling to obtain the next  
 913 solution. Ramachandran et al. (2020) used the prior distribution of optimal solutions to warp the  
 914 search space, expanding around high-probability regions of optimal solutions and shrinking around  
 915 low-probability regions. Truncated normal and gamma distributions were used to form the prior  
 916 distributions. Additionally, Hvarfner et al. (2022) incorporated prior Gaussian distributions about  
 917 locations of optimal solutions into the acquisition function, achieving competitive results across a  
 wide range of benchmarks. Mallik et al. (2023) also integrated prior knowledge of optimal hyperpa-  
 rameters to enhance the efficiency of Hyperband. Although using prior distributions of locations of

optimal solutions can accelerate optimization, accessing the prior for a specific task may not always be accessible.

#### B.4 TRANSFER SEARCH SPACE

In addition to leveraging experts’ prior, this method used information from previous tasks to reduce search space. For instance, Wistuba et al. (2015) pruned the bad regions of search space according to the results from previous tasks. Perrone & Shen (2019) and Li et al. (2022a) utilized previous tasks to design a sub-region of the entire search space for a new task. While these approaches have demonstrated efficiency in using promising regions instead of the entire space, they require the preparation of source tasks and the evaluation of task similarities to effectively select relevant tasks for learning promising regions, which can be challenging (Perrone & Shen, 2019).

### C. BENCHMARKS

#### C.1 TABULAR BENCHMARKS

**FCNet:** We utilized benchmarks for FCNet from Yahoo-Gym (Pfisterer et al., 2022), detailed in Table 3. The selected tasks include FCNet-Naval-Propulsion, FCNet-Protein-Structure, FCNet-Slice-Localization, and FCNet-Parkinsons-Telemonitoring.

**NAS-Bench-301:** Our NAS-Bench-301 benchmarks, also sourced from Yahoo-Gym (Pfisterer et al., 2022), focus on CIFAR-10. For hyperparameter details, refer to (Pfisterer et al., 2022).

**NAS-Bench-201:** This benchmark encompasses 6 hyperparameters for neural architecture search. It includes statistics from 15,625 CNN models across three datasets: CIFAR-10-valid, CIFAR-100, and ImageNet16-120 (Eggenesperger et al., 2021). We simplify their name as CIFAR-10, CIFAR-100, and ImageNet in this paper.

Table 3: Hyperparameter ranges for FCNet.

Parameter	Name	Type	Range
Fidelity	epoch	int	[1, 100]
Hyperparameter	batch size	int	[8, 64] (log-scale)
	initial learning rate	con	$[5e - 04, 0.1]$ (log-scale)
	dropout 1	con	[0.0, 0.6]
	dropout 2	con	[0.0, 0.6]
	number of units 1	int	[16, 512]
	number of units 2	int	[16, 512]

#### C.2 SURROGATE BENCHMARKS

We utilized four problems from the PD1 benchmarks and three from the JAHSBench surrogate benchmarks. Detailed information about these benchmarks is available in (Mallik et al., 2023).

#### C.3 RAW PROBLEMS

The hyperparameters for the deep neural networks, LeNet and ResNet-18, are detailed in Table 4.

Table 4: Hyperparameter ranges for LeNet and ResNet-18.

Parameter	Name	Type	Range
Fidelity	datasize	con	[0.3, 1]
Hyperparameter	batch size	int	[64, 512] (log-scale)
	initial learning rate	con	$[5e - 3, 0.1]$ (log-scale)
	momentum	con	[0.5, 0.99]
	weight decay	con	$[1e - 5, 1e - 2]$
	nesterov	cat	{True, False}

C.4 PARAMETERS OF LF PROBLEMS

Table 5 shows parameters of LF problems used in the first phase and the corresponding OVL between high- and low-fidelity problems.

Table 5: Parameters of LF problems used in the first phase and the corresponding OVL between high- and low-fidelity problems.

Tasks	LF parameter	OVL	Tasks	LF parameter	OVL
MFH3	epoch=4	0.713	MFH6	epoch=4	0.728
XGBoost-40981	datasize=0.3	0.885	FCNet-Naval-Propulsion	epoch=4	0.901
XGBoost-41146	datasize=0.3	0.933	FCNet-Protein-Structure	epoch=4	0.694
XGBoost-1489	datasize=0.3	0.805	FCNet-Slice-Localization	epoch=4	0.731
XGBoost-1067	datasize=0.3	0.883	FCNet-Parkinsons-Telemonitoring	epoch=4	0.854
rpart-40981	datasize=0.3	0.167	NAS-Bench-301-CIFAR-10	epoch=20	0.712
rpart-41146	datasize=0.3	0.488	NAS-Bench-201-CIFAR-10	epoch=20	0.760
rpart-1089	datasize=0.3	0.89	NAS-Bench-201-CIFAR-100	epoch=20	0.727
rpart-1067	datasize=0.3	0.642	NAS-Bench-201-ImageNet	epoch=20	0.696
ranger-40981	datasize=0.3	0.711	JAHS-CIFAR-10	epoch=4	0.574
ranger-41146	datasize=0.3	0.782	JAHS-Colorectal-Histology	epoch=4	0.523
ranger-1489	datasize=0.3	0.77	JAHS-Fashion-MNIST	epoch=4	0.532
ranger-1067	datasize=0.3	NAN	PD1-LM1B	epoch=30	0.934
glmnet-40981	datasize=0.3	0.294	PD1-WMT	epoch=4	0.926
glmnet-41146	datasize=0.3	0.962	PD1-CIFAR-100	epoch=45	0.797
glmnet-1489	datasize=0.3	0.918	PD1-ImageNet	epoch=20	0.727
glmnet-1067	datasize=0.3	0.648	NAN	NAN	NAN

D. OVL BETWEEN LOW AND HIGH FIDELITY

Figure 11 presents the overlapping coefficients between good solutions in high- and low-fidelity settings across various HPO tasks. The overlapping coefficients are computed using equation (3) and are visualized over epochs or datasets for different benchmarks. It can be observed that the overlap generally increases with the number of iterations (e.g., epochs or dataset size), especially for FCNet, NAS-Bench-201, and RecNet-18. However, for NAS-Bench-301 and JAHS, the overlap exhibits a trend of decreasing in the middle stages before rising again. Specifically:

- FCNet and RecNet-18 show an overlap that is already close to or greater than 0.7 even at smaller iteration parameters (e.g., fewer epochs or smaller datasets), indicating strong LF-HF consistency and stability at earlier stages.
- For PD1 and NAS-Bench-201, the overlap reaches or exceeds 0.7 at specific points, such as epoch = 10 for PD1 and epoch = 20 for NAS-Bench-201, suggesting that these tasks achieve good LF-HF agreement relatively early in the optimization process.
- NAS-Bench-301 and JAHS, on the other hand, maintain relatively low overlap in the initial and middle stages. The overlap only increases significantly as the settings approach the HF level, indicating that these tasks require longer training or higher fidelity to achieve substantial LF-HF alignment.

E. DISTRIBUTION OF SOLUTIONS IN THE HF AND LF.

In this section, we present the distribution of solutions in the HF and LF settings. Figure 12 and Figure 13 shows good solutions identified by Lamda+BOHB at the first phase for FCNet, NAS-Bench-301 and JAHS. It can be observed that there is an overlap between the good solutions in the HF and LF settings, and the good solutions found in the first phase are close to these overlapping regions.

Moreover, Figure 14, Figure 15, Figure 16 and Figure 17 show the distribution of good solutions in the HF and LF. Across all 15 problems, a clear overlap is observed between the good solutions identified in both fidelity levels.

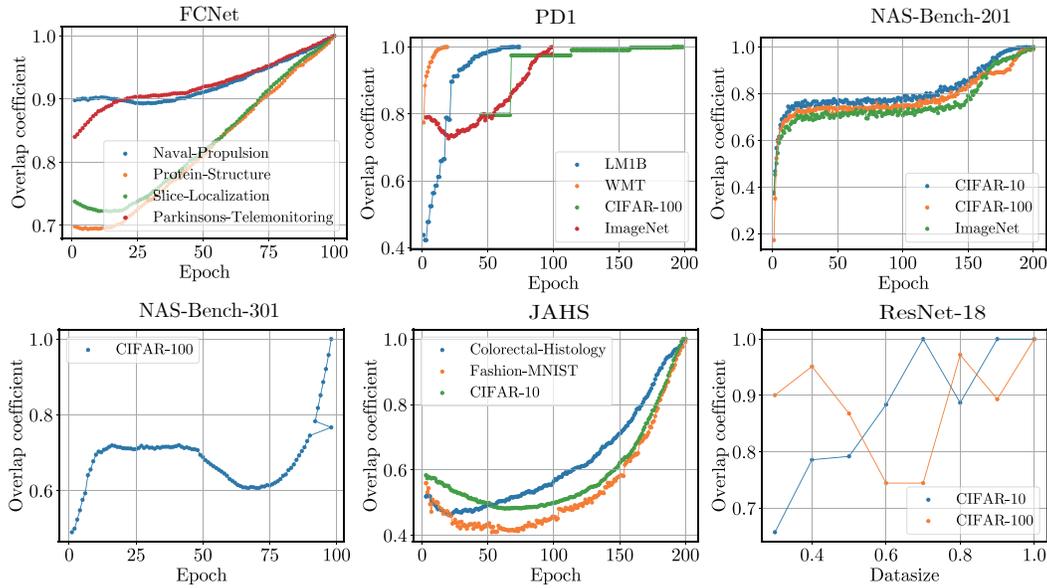


Figure 11: The overlapping coefficient of the HPO tasks.

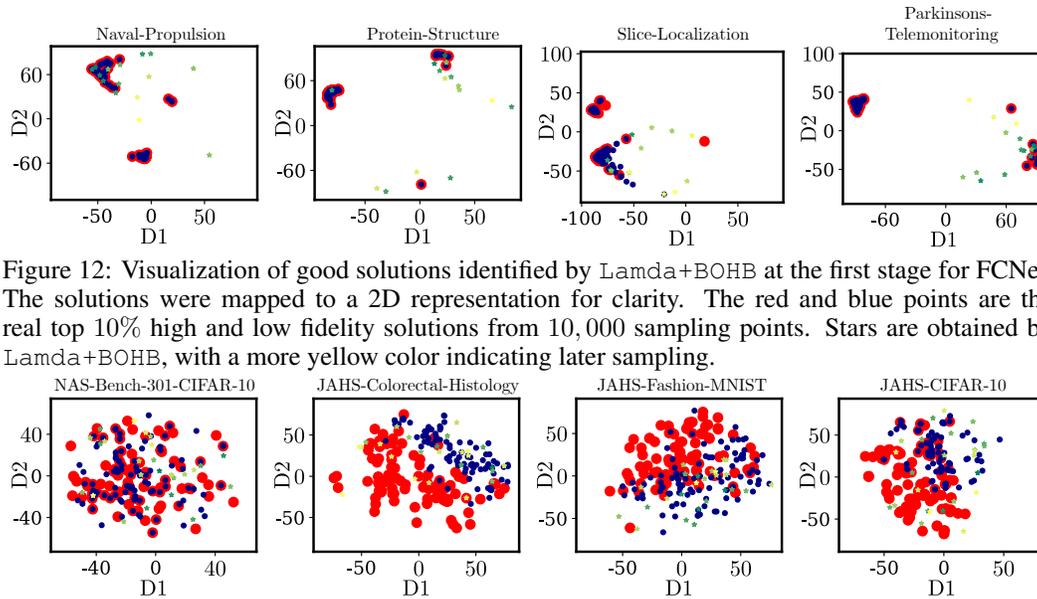


Figure 12: Visualization of good solutions identified by Lamda+BOHB at the first stage for FCNet. The solutions were mapped to a 2D representation for clarity. The red and blue points are the real top 10% high and low fidelity solutions from 10,000 sampling points. Stars are obtained by Lamda+BOHB, with a more yellow color indicating later sampling.

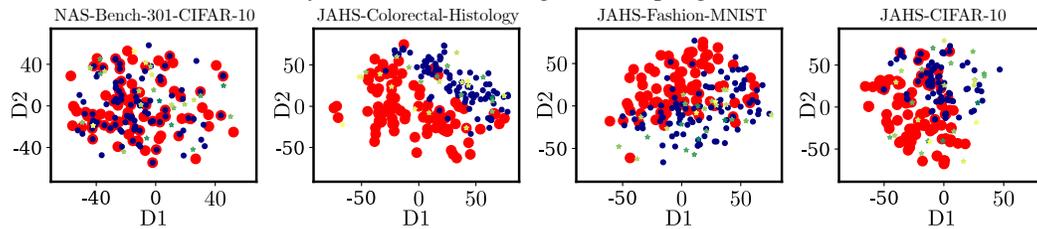


Figure 13: Visualization of good solutions identified by Lamda+BOHB at the first phase for NAS-Bench-301 and JAHS. The solutions were mapped to 2D. The red and blue points are the real top 10% high and low fidelity solutions from 10,000 sampling points. Stars are obtained by Lamda+BOHB, with a more yellow color indicating later sampling.

## F. EXPERIMENTAL RESULTS

### F.1 CONVERGENCE CURVES

This section presents the convergence curves of the proposed algorithm and its peer methods. The convergence performance of Lamda under different algorithms, including BOHB, MUMBO, BO, and RS, is illustrated in Figure 18, Figure 19, Figure 20, Figure 21, Figure 22, Figure 23, Figure 24 and Figure 25. The experimental results demonstrated that using Lamda can enhance performanc

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094

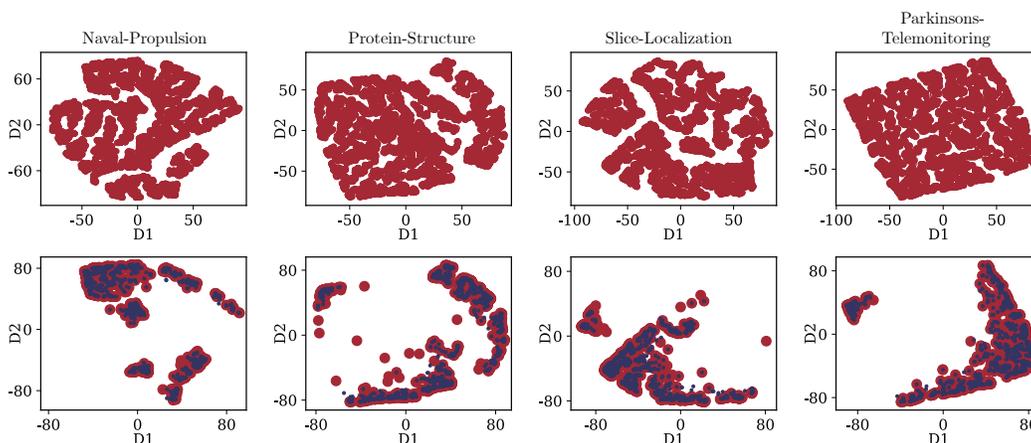


Figure 14: (Top) The visualization of 10,000 hyperparameter configurations in the 2D space for FCNet. (Bottom) The visualization of the top 10% solutions at the high and low fidelity, with high fidelity solutions marked in red. In this case, LF is obtained by setting the epoch number as four.

1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111

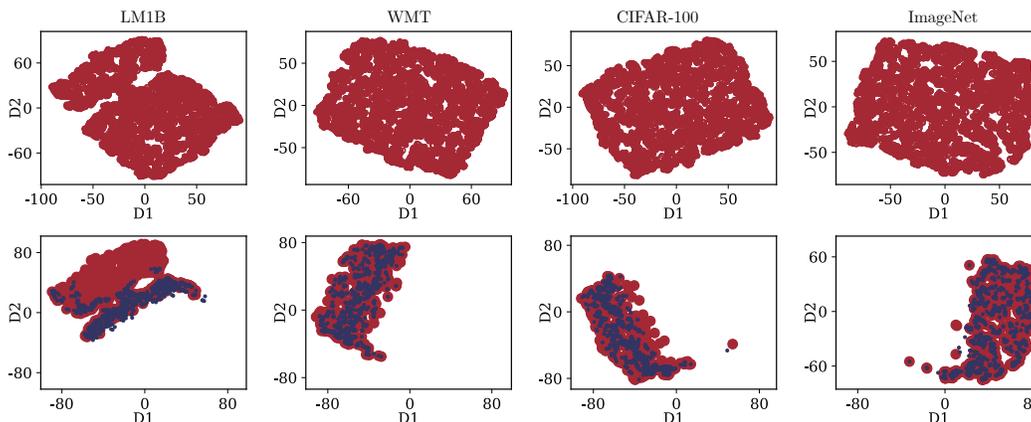


Figure 15: (Top) The visualization of the 10,000 hyperparameter configurations in the 2D space for PD1. (Bottom) The visualization of the top 10% solutions at the high and low fidelity, with high fidelity solutions marked in red. In this case, LF is obtained by setting the epoch number as 30, 4, 45, and 20 for the four tasks.

1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

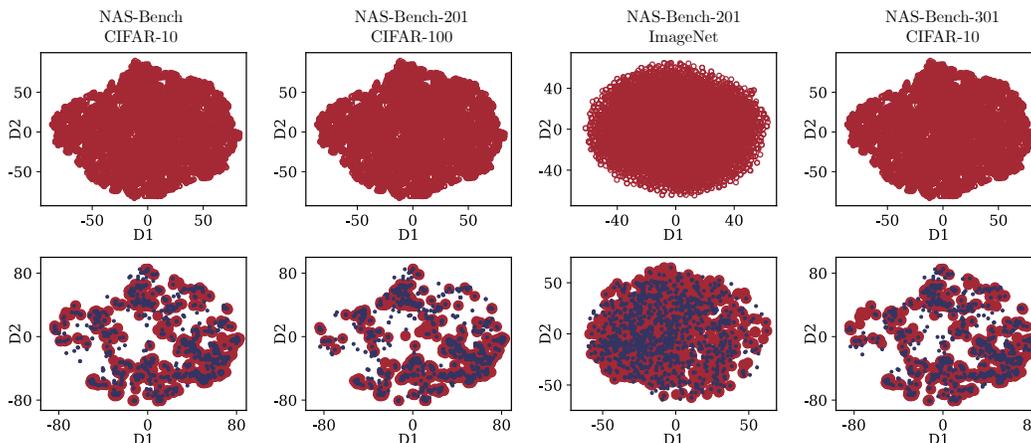


Figure 16: (Top) The visualization of 10,000 hyperparameter configurations in the 2D space for NAS-Bench-201 and NAS-Bench-301. (Bottom) The visualization of the top 10% solutions at the high and low fidelity, with high fidelity solutions marked in red. In this case, LF is obtained by setting the epoch number as 20.

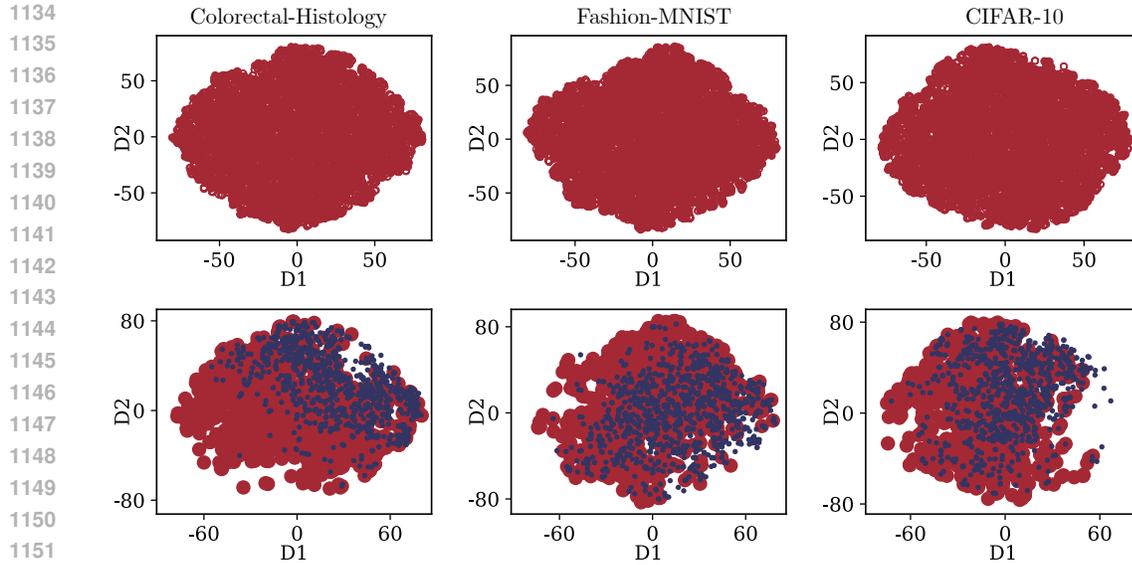


Figure 17: (Top) The visualization of 10,000 hyperparameter configurations in the 2D space for JASH. (Bottom) The visualization of the top 10% solutions at the high and low fidelity, with high fidelity solutions marked in red. In this case, LF is obtained by setting the epoch number as 20.

of the baseline algorithms. The convergence curves of peer algorithms are presented in Figure 26 and Figure 27. Lamda also performs good compared with peer algorithms.

Figure 28, Figure 29, Figure 30, Figure 31 and Figure 32 depict the results of parameter analysis including budget ( $B$ ), thresholds ( $\gamma$ ) for stopping the first phase, the interval for calculating the overlapping coefficient ( $\Delta$ ), the quantile ( $\alpha$ ) used in the TPE, and the initial weight ( $w_0$ ). These results indicate that the algorithm’s performance is robust to the hyperparameter settings.

## F.2 TABLE RESULTS OF PEER ALGORITHMS

Table 6 shows the peer algorithms’ final validation errors of the current incumbent at 100 HF evaluations horizons.

1188  
 1189  
 1190  
 1191  
 1192  
 1193  
 1194  
 1195  
 1196  
 1197  
 1198  
 1199  
 1200  
 1201  
 1202  
 1203  
 1204  
 1205  
 1206  
 1207  
 1208  
 1209  
 1210  
 1211  
 1212  
 1213  
 1214  
 1215  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241

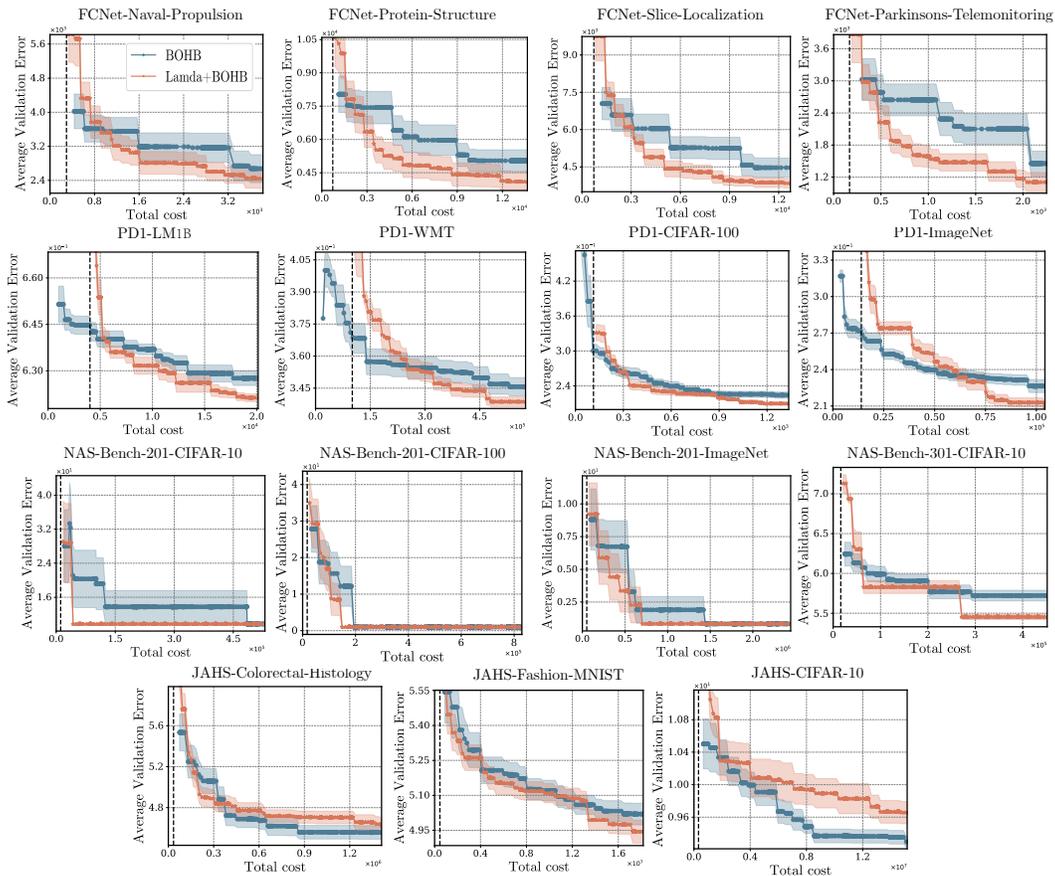


Figure 18: Validation error observed in tuning 15 HPO tasks, using BOHB as the baseline.

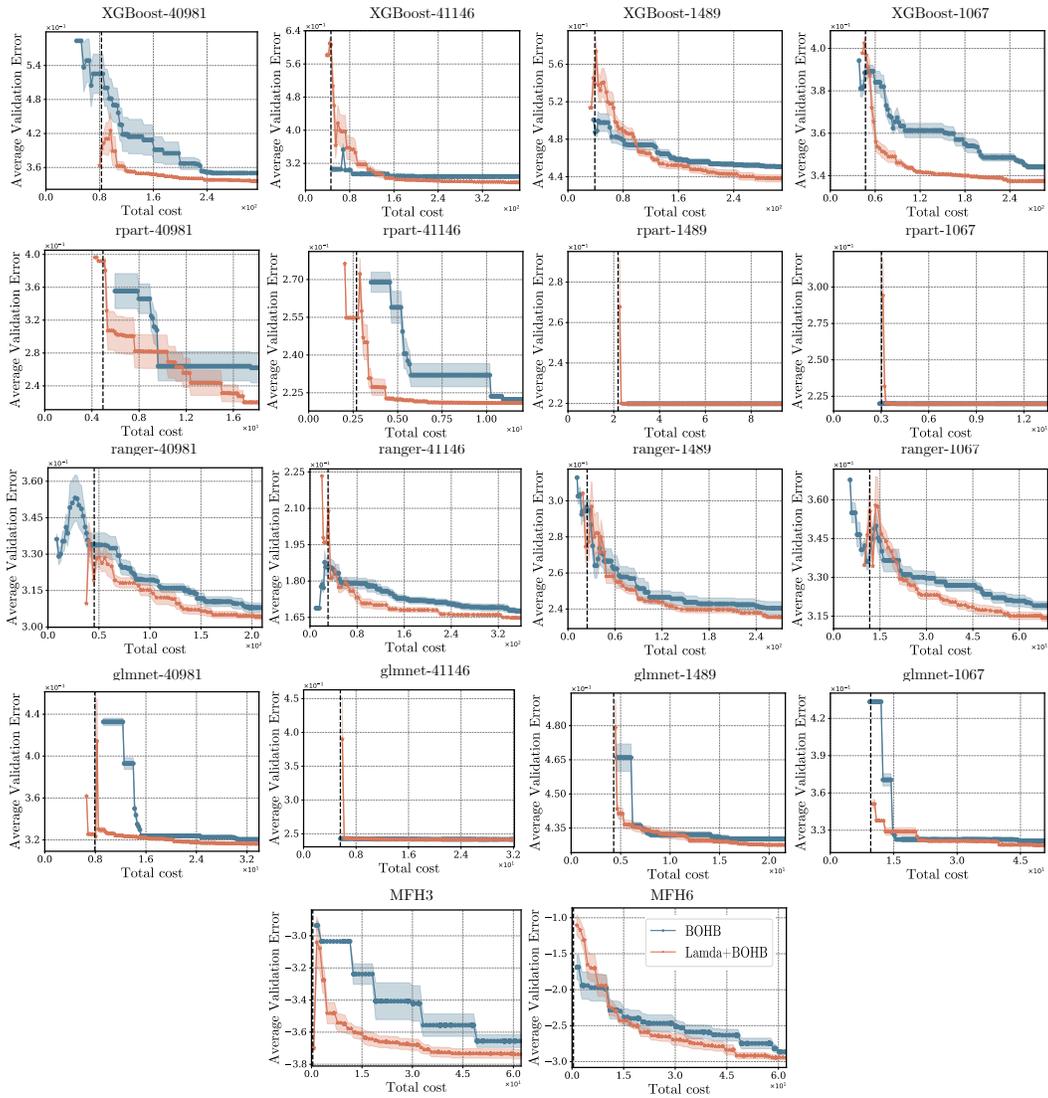


Figure 19: Validation error observed in tuning 18 HPO tasks, using BOHB as the baseline.

1296  
 1297  
 1298  
 1299  
 1300  
 1301  
 1302  
 1303  
 1304  
 1305  
 1306  
 1307  
 1308  
 1309  
 1310  
 1311  
 1312  
 1313  
 1314  
 1315  
 1316  
 1317  
 1318  
 1319  
 1320  
 1321  
 1322  
 1323  
 1324  
 1325  
 1326  
 1327  
 1328  
 1329  
 1330  
 1331  
 1332  
 1333  
 1334  
 1335  
 1336  
 1337  
 1338  
 1339  
 1340  
 1341  
 1342  
 1343  
 1344  
 1345  
 1346  
 1347  
 1348  
 1349

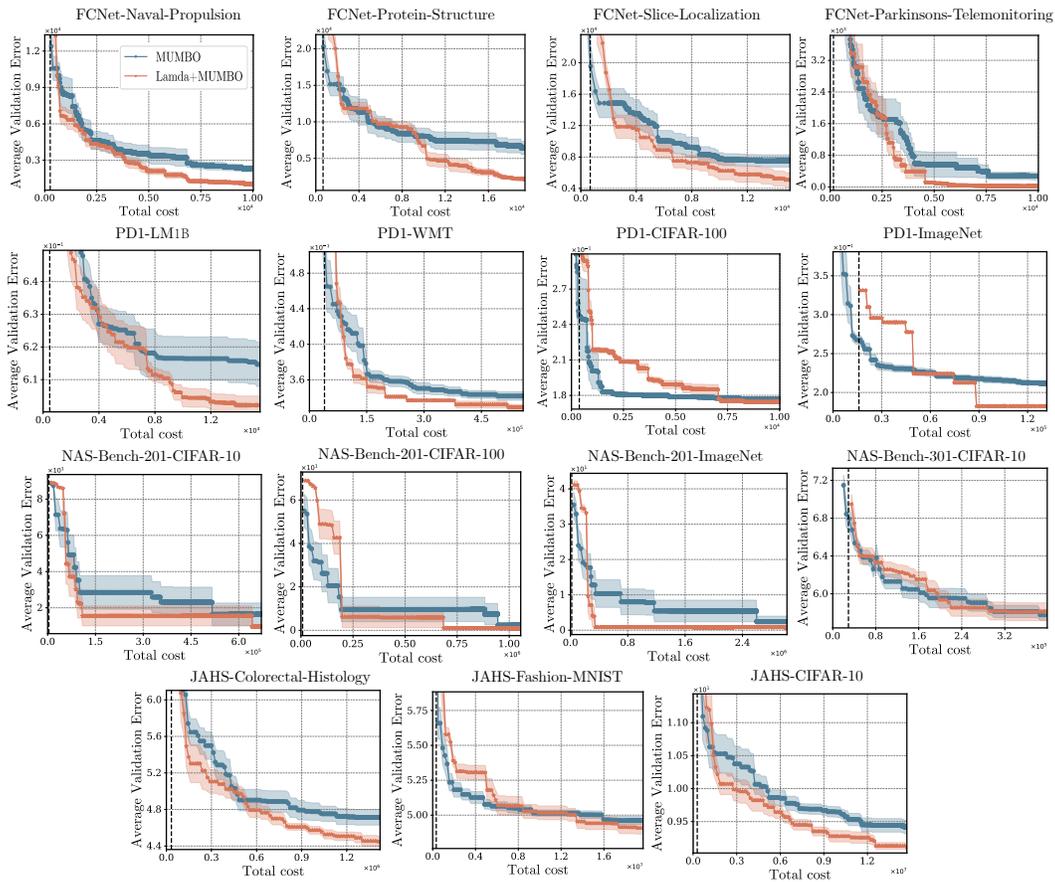


Figure 20: Validation error observed in tuning 15 HPO tasks, using MUMBO as the baseline.

1350  
 1351  
 1352  
 1353  
 1354  
 1355  
 1356  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390  
 1391  
 1392  
 1393  
 1394  
 1395  
 1396  
 1397  
 1398  
 1399  
 1400  
 1401  
 1402  
 1403

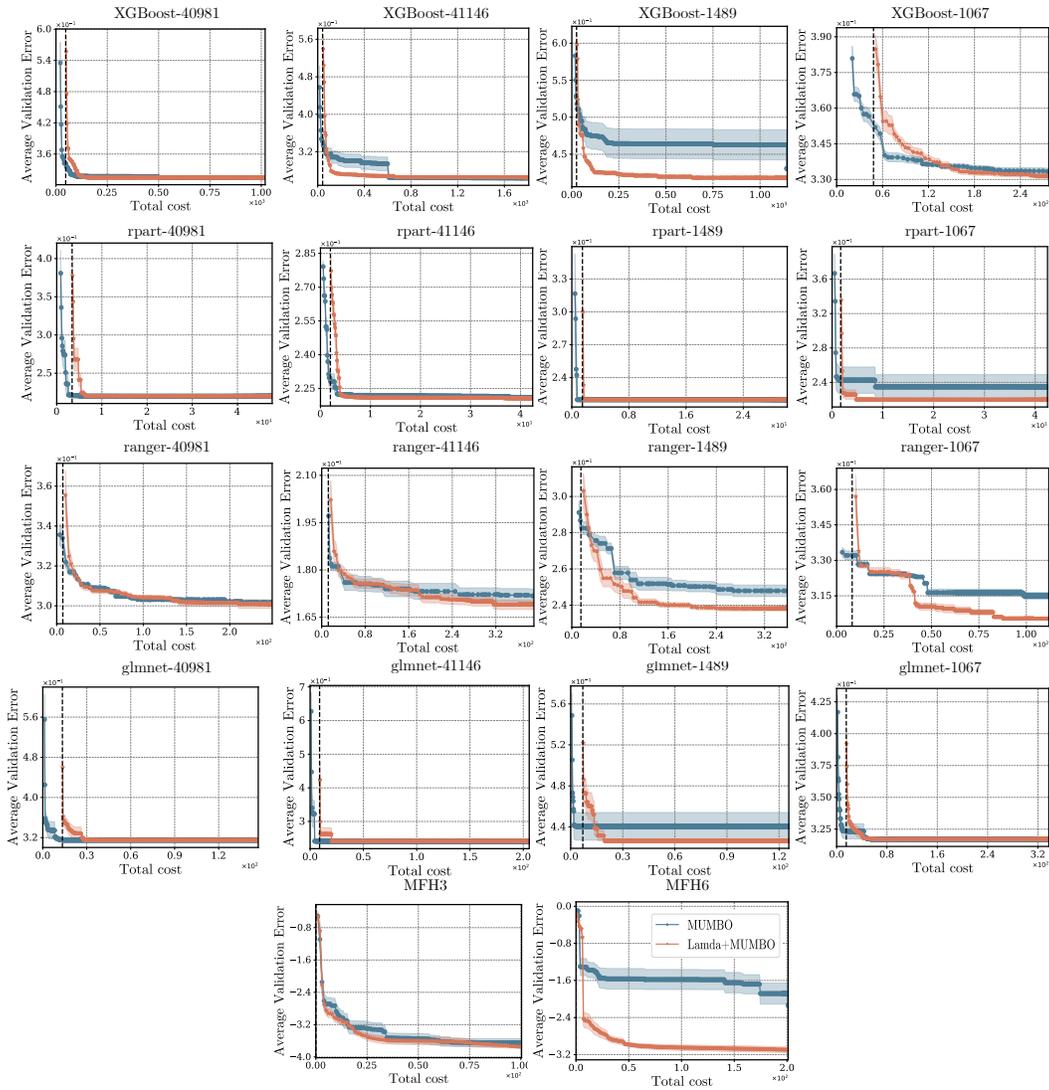


Figure 21: Validation error observed in tuning 18 HPO tasks, using MUMBO as the baseline.

1404  
 1405  
 1406  
 1407  
 1408  
 1409  
 1410  
 1411  
 1412  
 1413  
 1414  
 1415  
 1416  
 1417  
 1418  
 1419  
 1420  
 1421  
 1422  
 1423  
 1424  
 1425  
 1426  
 1427  
 1428  
 1429  
 1430  
 1431  
 1432  
 1433  
 1434  
 1435  
 1436  
 1437  
 1438  
 1439  
 1440  
 1441  
 1442  
 1443  
 1444  
 1445  
 1446  
 1447  
 1448  
 1449  
 1450  
 1451  
 1452  
 1453  
 1454  
 1455  
 1456  
 1457

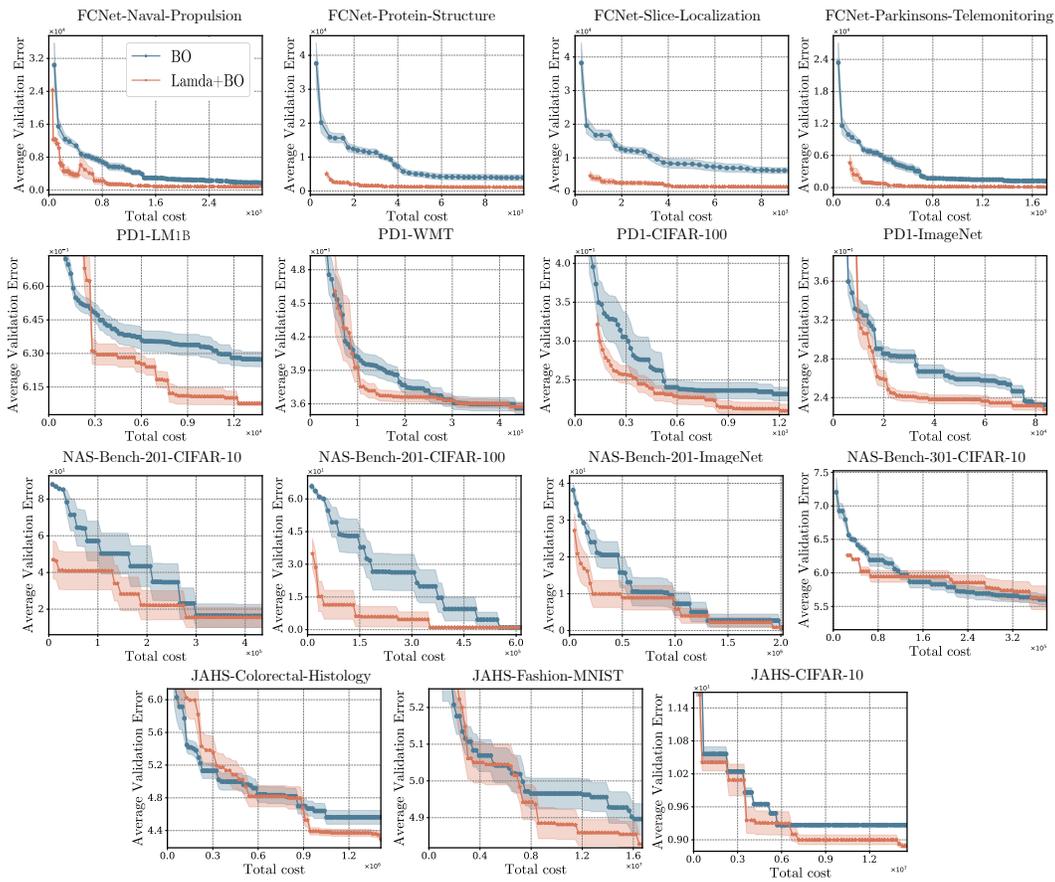


Figure 22: Validation error observed in tuning 15 HPO tasks, using BO as the baseline.

1458  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485  
 1486  
 1487  
 1488  
 1489  
 1490  
 1491  
 1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511

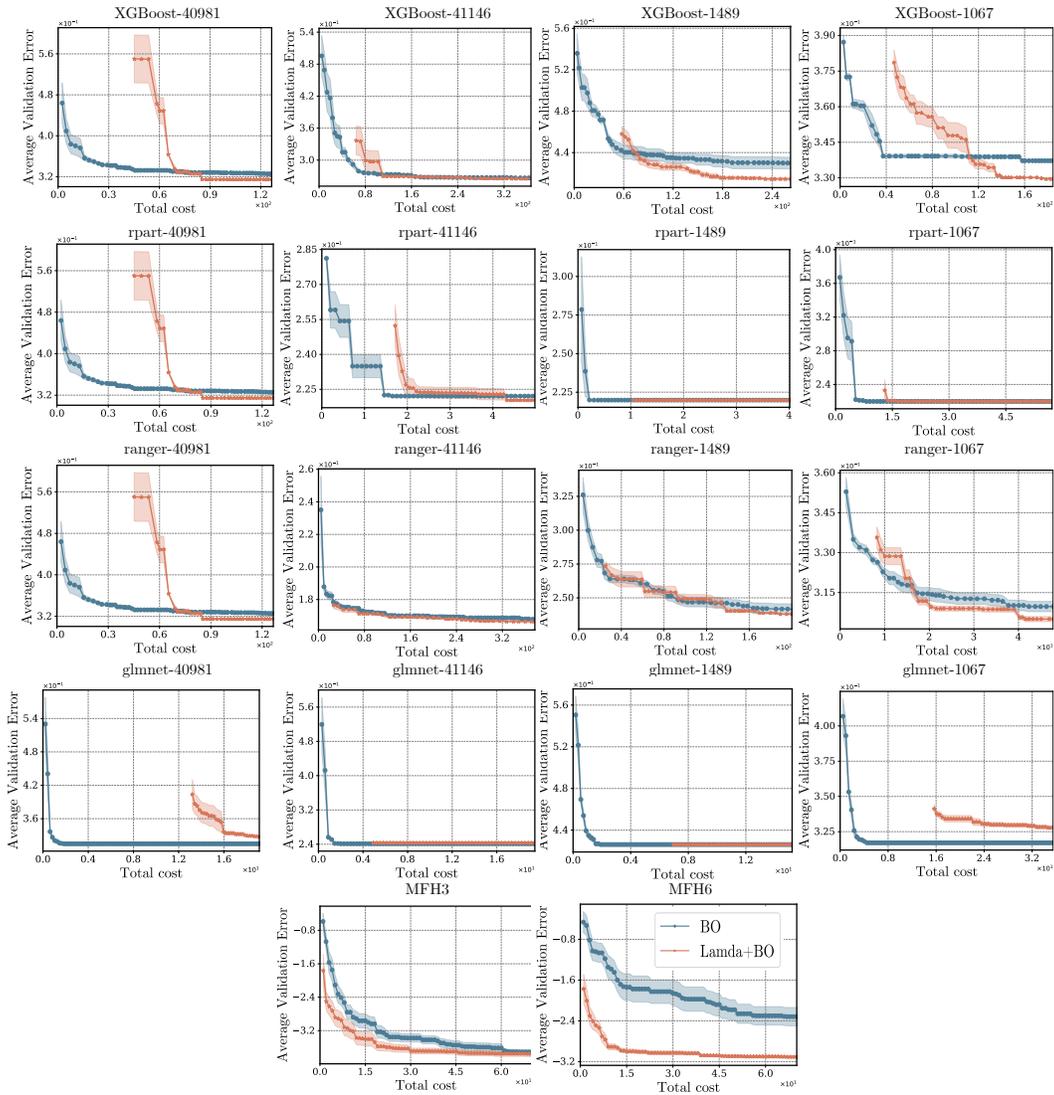


Figure 23: Validation error observed in tuning 18 HPO tasks, using BO as the baseline.

1512  
 1513  
 1514  
 1515  
 1516  
 1517  
 1518  
 1519  
 1520  
 1521  
 1522  
 1523  
 1524  
 1525  
 1526  
 1527  
 1528  
 1529  
 1530  
 1531  
 1532  
 1533  
 1534  
 1535  
 1536  
 1537  
 1538  
 1539  
 1540  
 1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552  
 1553  
 1554  
 1555  
 1556  
 1557  
 1558  
 1559  
 1560  
 1561  
 1562  
 1563  
 1564  
 1565

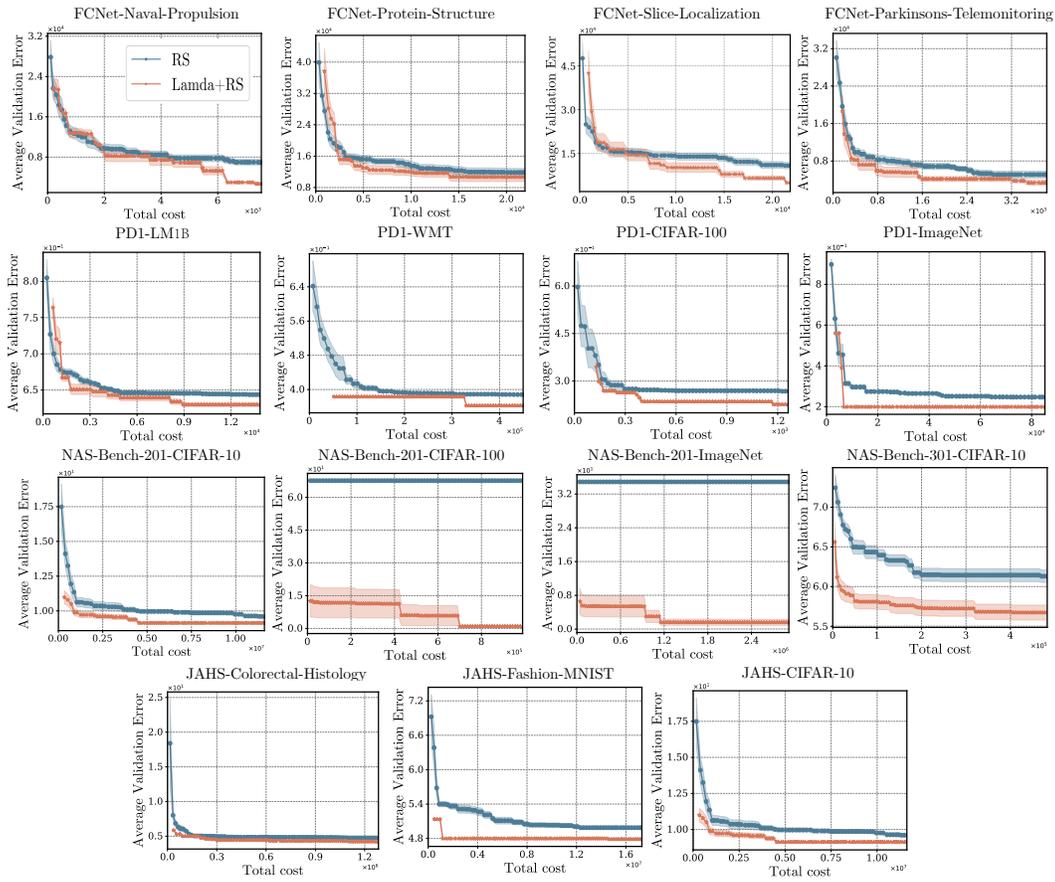


Figure 24: Validation error observed in tuning 15 HPO tasks, using RS as the baseline.

1566  
 1567  
 1568  
 1569  
 1570  
 1571  
 1572  
 1573  
 1574  
 1575  
 1576  
 1577  
 1578  
 1579  
 1580  
 1581  
 1582  
 1583  
 1584  
 1585  
 1586  
 1587  
 1588  
 1589  
 1590  
 1591  
 1592  
 1593  
 1594  
 1595  
 1596  
 1597  
 1598  
 1599  
 1600  
 1601  
 1602  
 1603  
 1604  
 1605  
 1606  
 1607  
 1608  
 1609  
 1610  
 1611  
 1612  
 1613  
 1614  
 1615  
 1616  
 1617  
 1618  
 1619

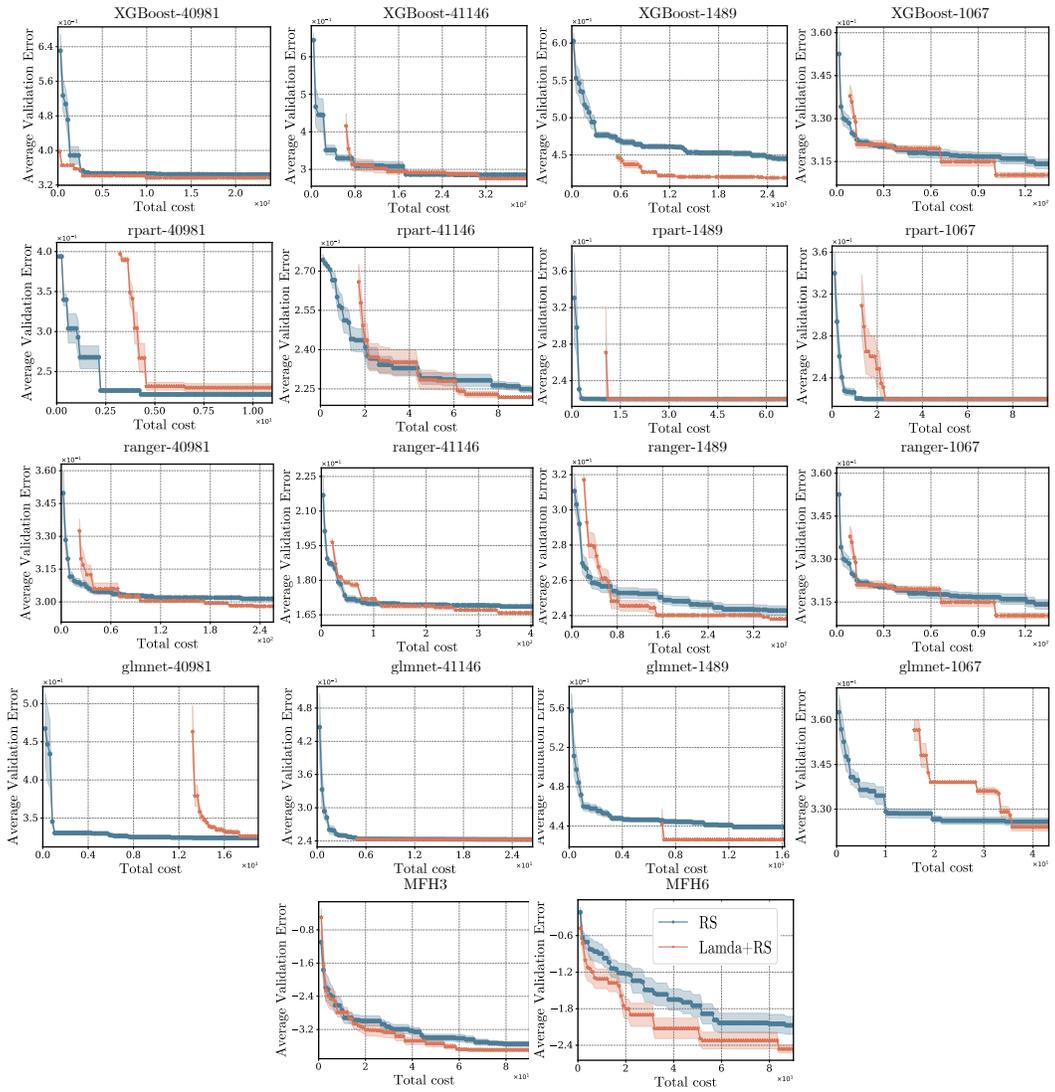


Figure 25: Validation error observed in tuning 18 HPO tasks, using RS as the baseline.

1620  
 1621  
 1622  
 1623  
 1624  
 1625  
 1626  
 1627  
 1628  
 1629  
 1630  
 1631  
 1632  
 1633  
 1634  
 1635  
 1636  
 1637  
 1638  
 1639  
 1640  
 1641  
 1642  
 1643  
 1644  
 1645  
 1646  
 1647  
 1648  
 1649  
 1650  
 1651  
 1652  
 1653  
 1654  
 1655  
 1656  
 1657  
 1658  
 1659  
 1660  
 1661  
 1662  
 1663  
 1664  
 1665  
 1666  
 1667  
 1668  
 1669  
 1670  
 1671  
 1672  
 1673

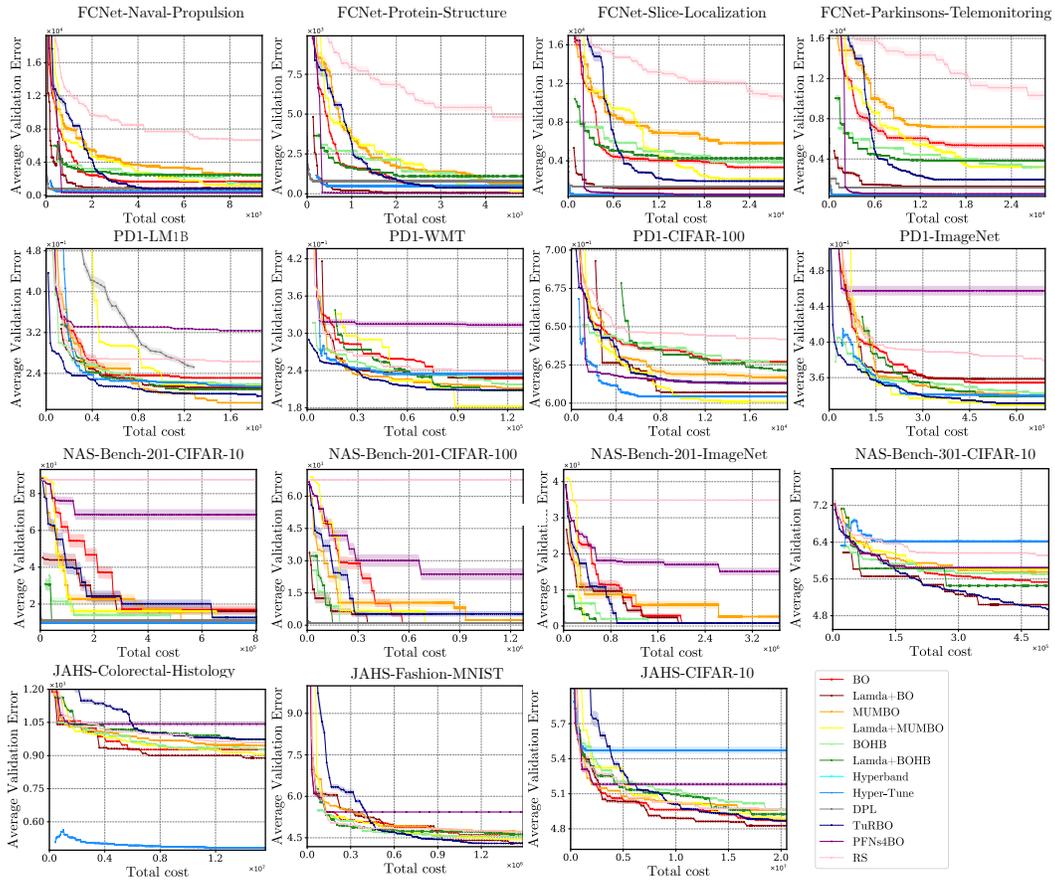


Figure 26: Validation error observed in tuning 15 HPO tasks, comparing peer algorithms.

1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727

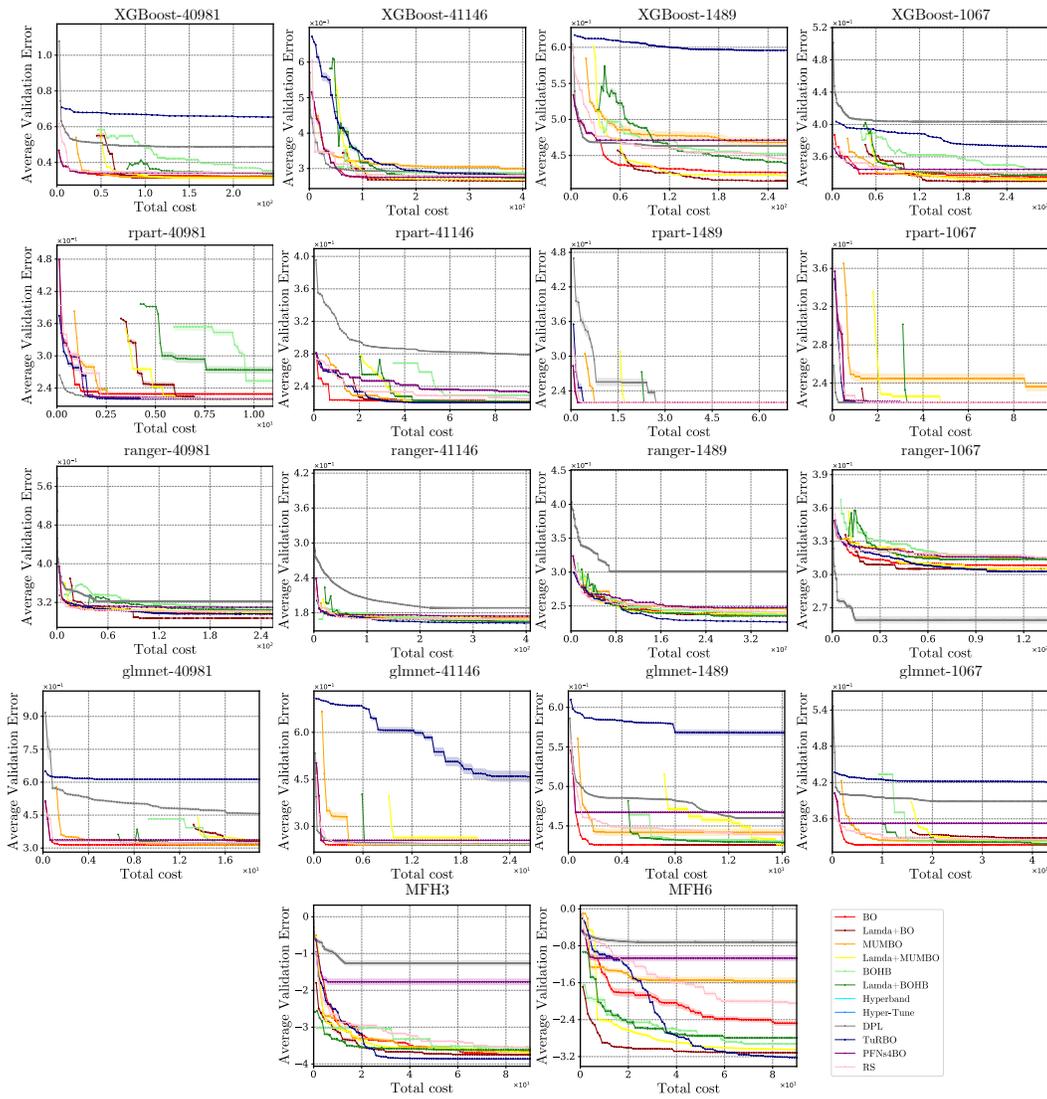


Figure 27: Validation error observed in tuning 18 HPO tasks, comparing peer algorithms.

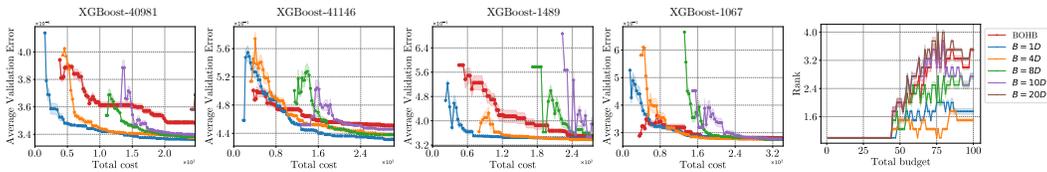


Figure 28: Validation errors of Lamda+BOHB under different parameters  $B$  used at the first phase.

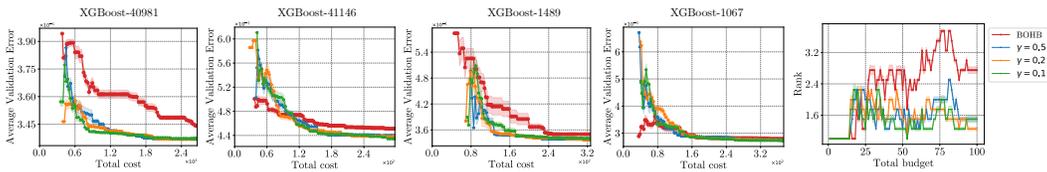


Figure 29: Validation error observed of Lamda+BOHB under different parameter  $\gamma$  at the first phase.

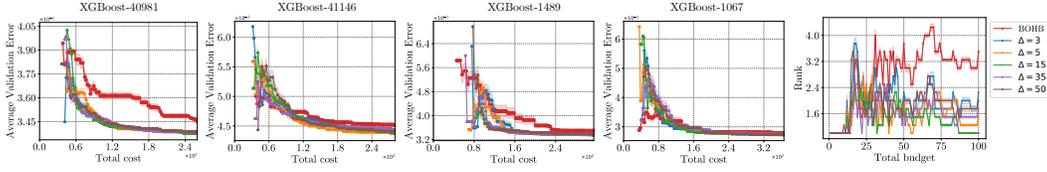


Figure 30: Validation error observed of Lamda+BOHB under different parameter  $\Delta$  at the first phase.

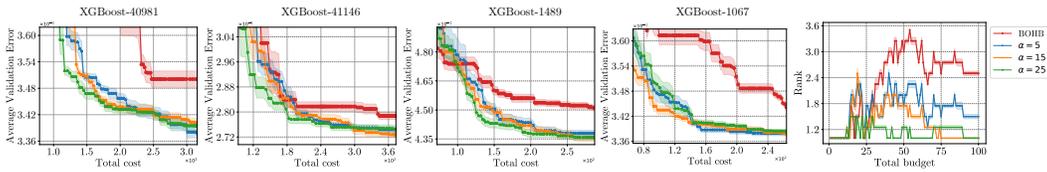


Figure 31: Validation error observed of Lamda+BOHB under different parameter  $\alpha$  at the first phase.

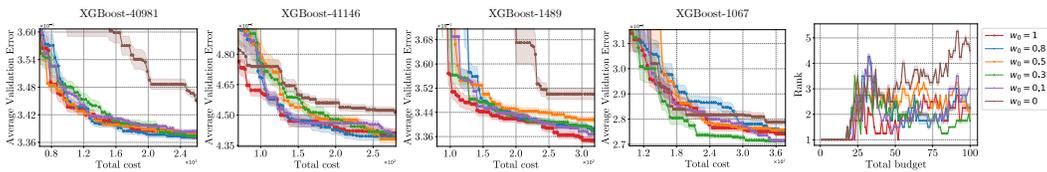


Figure 32: Validation error observed of Lamda+BOHB under different parameter  $w$  at the first phase.

1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1810  
1811  
1812  
1813  
1814

Table 6: Comparing peer algorithms’ final validation errors of the current incumbent at 100 HF evaluations horizons. Runs are averaged over 31 seeds. (The bolded part indicates: ”under the Wilcoxon rank-sum test, methods incorporating priors significantly outperform their baselines.”)

Tasks	BOHB	Lambda+BOHB	MUMBO	Lambda+MUMBO	BO	Lambda+BO	RS	Lambda+RS	PriorBand	Lambda+PriorBand
MFF3	-3.65e+00(1.339e-01)	<b>-3.716e+00(1.056e-01)</b>	-3.860e+00(3.698e-01)	-3.846e+00(2.174e-01)	-3.822e+00(2.172e-01)	-3.843e+00(5.508e-02)	<b>-3.574e+00(2.529e-01)</b>	-3.718e+00(1.323e-01)	-3.857e+00(3.553e-03)	-3.852e+00(2.401e-03)
MFF6	-2.86e+00(1.364e-01)	<b>-2.935e+00(1.427e-01)</b>	-1.872e+00(0.0)	<b>-3.178e+00(2.752e-01)</b>	-2.138e+00(0.286e-01)	<b>-3.112e+00(1.429e-02)</b>	-2.260e+00(6.577e-01)	-2.396e+00(6.227e-01)	-3.187e+00(1.506e-01)	-3.167e+00(5.198e-02)
XGBoost-40981	3.456e-01(9.817e-03)	3.337e-01(9.490e-03)	3.147e-01(2.235e-07)	3.147e-01(2.235e-07)	3.337e-01(9.493e-03)	<b>3.147e+00(1.429e-02)</b>	3.442e-01(1.143e-03)	<b>3.367e-01(1.528e-03)</b>	NA	NA
XGBoost-41146	2.883e-01(1.238e-03)	<b>2.783e-01(6.558e-03)</b>	2.634e-01(1.694e-03)	2.635e-01(1.063e-03)	2.600e-01(1.620e-03)	<b>2.646e-01(1.823e-03)</b>	2.824e-01(8.999e-03)	<b>2.768e-01(4.922e-04)</b>	NA	NA
XGBoost-1489	4.477e-01(8.205e-03)	<b>4.330e-01(2.073e-02)</b>	4.616e-01(5.269e-02)	<b>4.116e-01(1.450e-02)</b>	4.164e-01(1.462e-02)	4.116e-01(5.029e-03)	4.415e-01(1.642e-02)	<b>4.187e-01(5.146e-03)</b>	NA	NA
XGBoost-1067	3.453e-01(3.377e-03)	<b>3.377e-01(3.775e-03)</b>	3.345e-01(5.656e-03)	<b>3.301e-01(5.028e-03)</b>	3.357e-01(7.962e-04)	<b>3.298e-01(1.825e-04)</b>	3.398e-01(0.0)	<b>3.345e-01(0.0)</b>	NA	NA
rpart-40981	2.200e-01(6.758e-02)	<b>2.199e-01(8.859e-06)</b>	2.199e-01(7.525e-06)	2.199e-01(1.937e-06)	2.199e-01(7.799e-06)	<b>2.200e-01(2.866e-04)</b>	2.213e-01(8.945e-04)	<b>2.199e-01(1.214e-05)</b>	NA	NA
rpart-41146	2.222e-01(1.529e-03)	<b>2.210e-01(9.219e-04)</b>	2.210e-01(4.312e-04)	2.205e-01(3.331e-04)	2.210e-01(1.799e-03)	<b>2.204e-01(1.356e-04)</b>	2.233e-01(3.626e-03)	<b>2.216e-01(1.541e-03)</b>	NA	NA
rpart-1089	2.198e-01(0.0)	2.198e-01(0.0)	2.198e-01(0.0)	2.198e-01(0.0)	2.198e-01(0.0)	2.198e-01(0.0)	2.198e-01(0.0)	2.198e-01(0.0)	NA	NA
rpart-1067	3.087e-01(4.186e-03)	<b>3.038e-01(6.079e-03)</b>	2.993e-01(2.547e-03)	3.022e-01(3.181e-03)	3.024e-01(7.710e-03)	<b>2.835e-01(1.074e-02)</b>	3.019e-01(1.860e-03)	<b>2.967e-01(2.828e-03)</b>	NA	NA
ranger-40981	1.645e-01(6.803e-03)	1.645e-01(2.341e-03)	1.683e-01(2.645e-03)	<b>1.679e-01(4.282e-03)</b>	1.670e-01(3.793e-03)	1.673e-01(1.883e-03)	1.685e-01(1.981e-03)	<b>1.641e-01(2.814e-03)</b>	NA	NA
ranger-41146	2.328e-01(5.845e-03)	2.358e-01(2.923e-03)	2.440e-01(1.390e-03)	2.372e-01(2.984e-03)	2.336e-01(2.404e-03)	<b>2.316e-01(1.231e-03)</b>	2.464e-01(4.325e-03)	<b>2.381e-01(1.180e-03)</b>	NA	NA
ranger-1067	3.182e-01(6.926e-03)	<b>3.136e-01(6.902e-03)</b>	3.138e-01(5.727e-03)	<b>3.055e-01(5.035e-04)</b>	3.142e-01(5.990e-03)	<b>3.061e-01(5.558e-03)</b>	3.180e-01(2.325e-03)	<b>3.089e-01(3.449e-03)</b>	NA	NA
glimmer-40981	3.200e-01(2.952e-03)	<b>3.161e-01(4.786e-04)</b>	3.149e-01(0.0)	3.149e-01(0.0)	3.149e-01(0.0)	3.269e-01(6.331e-03)	3.241e-01(4.395e-04)	<b>3.247e-01(1.448e-03)</b>	NA	NA
glimmer-41146	2.415e-01(3.053e-04)	2.413e-01(3.809e-04)	2.411e-01(0.0)	4.262e-01(0.0)	2.411e-01(0.0)	2.430e-01(3.797e-04)	2.425e-01(4.839e-04)	<b>2.429e-01(5.945e-04)</b>	NA	NA
glimmer-1489	4.298e-01(2.869e-03)	<b>4.274e-01(5.247e-04)</b>	4.392e-01(0.0)	4.262e-01(0.0)	4.262e-01(0.0)	4.262e-01(0.0)	4.392e-01(6.759e-03)	<b>4.262e-01(0.0)</b>	NA	NA
glimmer-1067	3.214e-01(4.997e-04)	<b>3.178e-01(7.455e-05)</b>	3.177e-01(0.0)	3.171e-01(0.0)	3.171e-01(0.0)	3.282e-01(6.788e-03)	3.232e-01(5.589e-03)	<b>3.210e-01(4.693e-03)</b>	NA	NA
FCNet-Naval-Propulsion	2.537e+03(1.014e+03)	<b>2.411e+03(6.557e+02)</b>	2.458e+03(7.802e+02)	<b>1.002e+03(5.606e+02)</b>	2.627e+03(2.387e+02)	<b>2.121e+02(1.074e-02)</b>	6.766e+03(6.833e+02)	<b>3.570e+03(1.043e+03)</b>	4.461e+03(8.544e+02)	<b>3.956e+03(4.427e+02)</b>
FCNet-Protein-Structure	5.030e+03(2.183e+03)	<b>4.334e+03(2.098e+03)</b>	6.529e+03(4.761e+03)	<b>2.015e+03(6.318e+02)</b>	2.764e+03(3.452e+03)	<b>1.156e+03(1.424e+02)</b>	1.367e+04(6.222e+02)	<b>1.097e+04(5.827e+03)</b>	6.273e+03(4.760e+03)	<b>3.501e+03(1.508e+02)</b>
FCNet-Slice-Localization	4.418e+03(5.736e+02)	<b>4.006e+03(3.385e+02)</b>	8.197e+03(4.281e+03)	<b>3.475e+03(2.819e+03)</b>	4.825e+03(5.299e+03)	<b>1.78e+03(1.161e+02)</b>	1.388e+04(1.960e+03)	<b>5.048e+03(5.287e+03)</b>	5.275e+03(1.709e+03)	<b>4.269e+03(5.377e+03)</b>
FCNet-Parkinsons-Telemetry	1.580e+05(4.194e+02)	<b>1.139e+05(7.311e+02)</b>	5.964e+02(4.477e+02)	<b>1.197e+01(9.573e+00)</b>	4.236e+07(3.323e+02)	<b>2.623e-05(6.821e+01)</b>	5.333e+03(5.352e+03)	4.273e+03(3.180e+03)	3.997e+02(8.256e+02)	4.041e+02(1.419e+03)
NAS-Bench-301-CIFAR-10	5.760e+00(2.620e-01)	<b>5.454e+00(1.091e-01)</b>	5.761e+00(1.368e-01)	<b>5.787e+00(3.964e-01)</b>	5.492e+00(2.995e-01)	5.629e+00(5.881e-01)	6.201e+00(4.134e-01)	5.538e+00(4.347e-01)	NA	NA
NAS-Bench-201-CIFAR-100	9.712e+00(8.138e-10)	<b>9.712e+00(2.543e-10)</b>	9.792e+00(6.103e-10)	<b>9.712e+00(1.272e-10)</b>	9.712e+00(1.424e-09)	9.712e+00(9.664e-10)	9.712e+00(1.077e-09)	9.712e+00(1.077e-09)	NA	NA
NAS-Bench-201-ImageNet	8.333e+01(3.392e-10)	<b>8.333e+01(2.143e-10)</b>	8.333e+01(1.439e-10)	<b>8.100e+01(0.0)</b>	1.000e+00(0.0)	8.333e+01(2.384e-10)	8.333e+01(0.0)	<b>8.333e+01(4.230e-10)</b>	NA	NA
MPS-CIFAR-10	4.224e+00(1.828e-01)	<b>4.224e+00(3.682e-01)</b>	9.340e+00(3.337e-01)	<b>9.022e+00(1.298e-01)</b>	9.260e+00(3.337e-01)	<b>4.188e+00(2.774e-01)</b>	9.530e+00(3.553e-01)	<b>9.201e+00(4.544e-01)</b>	NA	NA
JAHNS-Genetic-Pathology	4.962e+00(6.988e-02)	<b>4.962e+00(2.062e-01)</b>	5.092e+00(1.478e-01)	<b>4.951e+00(2.517e-01)</b>	4.576e+00(2.376e-01)	<b>4.188e+00(2.774e-01)</b>	4.944e+00(3.075e-01)	<b>4.261e+00(3.044e-02)</b>	NA	NA
JAHNS-FS	6.255e+01(1.978e-02)	<b>6.209e+01(4.692e-03)</b>	6.092e+01(3.676e-02)	<b>6.045e+01(1.400e-02)</b>	6.277e+01(1.818e-02)	6.063e+01(2.018e-04)	6.459e+01(1.125e-02)	<b>6.315e+01(9.575e-03)</b>	9.905e+01(0.0)	9.905e+01(3.879e-03)
PD1-LM1B	3.438e+01(1.379e-02)	<b>3.394e+01(1.243e-02)</b>	3.357e+01(4.336e-03)	<b>3.326e+01(8.931e-03)</b>	3.578e+01(1.671e-02)	3.570e+01(2.453e-02)	3.763e+01(0.0)	<b>3.624e+01(0.0)</b>	9.907e+01(2.832e-03)	9.907e+01(2.832e-03)
PD1-LM1T	2.141e+01(3.115e-02)	<b>2.097e+01(1.186e-02)</b>	1.767e+01(2.363e-03)	<b>1.743e+01(4.128e-04)</b>	2.345e+01(4.867e-02)	<b>1.946e+01(4.658e-02)</b>	2.723e+01(5.905e-03)	<b>2.269e+01(8.842e-03)</b>	NA	NA
PD1-CIFAR-100	2.283e-01(3.206e-02)	<b>2.050e-01(1.184e-02)</b>	1.942e-01(0.0)	<b>1.824e-01(0.0)</b>	2.263e-01(1.753e-03)	<b>2.145e-01(1.641e-02)</b>	2.360e-01(0.0)	<b>1.986e-01(0.0)</b>	NA	NA
PD1-ImageNet										

### F.3 ANALYZING THE IMPACT OF PROMISING REGIONS OVERLAP IN LF AND HF LANDSCAPES ON ALGORITHM PERFORMANCE

To better illustrate the solution distributions across low- and high-fidelity problems, we visualized the distributions in a 2D space. We sampled 10,000 hyperparameter configurations and computed their fitness values under both low- and high-fidelity settings. The data was then compressed into 2D using the UMAP McInnes & Healy (2018) algorithm. To further enhance interpretability, we applied linear interpolation to generate a continuous landscape surface. Figure 33, Figure 34, Figure 35, and Figure 36 present the landscapes of FCNet, PD1, NAS-Bench-201, NAS-Bench-301, and JAHS under both high- and low-fidelity settings. The overlap between high- and low-fidelity problems is also shown in the figures.

From these visualizations, we observe that FCNet and PD1 exhibit a high overlap (ranging from 0.7 to 0.9) in their promising regions across fidelity levels, with the promising regions being relatively concentrated. In contrast, while NAS-Bench-201 and NAS-Bench-301 have moderate OVL values (ranging from 0.67 to 0.76), their solution distributions are more dispersed. JAHS demonstrates lower OVL values (ranging from 0.52 to 0.57) and also shows a more scattered solution distribution.

The results indicate that for problems where the promising regions of LF and HF landscape have varying levels of overlaps, our method (`Lamda`) achieves better performance than the original algorithms under the frameworks of BOHB, MUMBO, BO, and RS, as illustrated in the convergence curves (Figure 18, Figure 20, Figure 22, and Figure 24). In particular, we would like highlight the results on JAHS-CIFAR-10, JAHS-Colorectal-Histology, and JAHS-Fashion-MNIST, whose overlaps are low. `Lamda` consistently enhance the baseline algorithms.

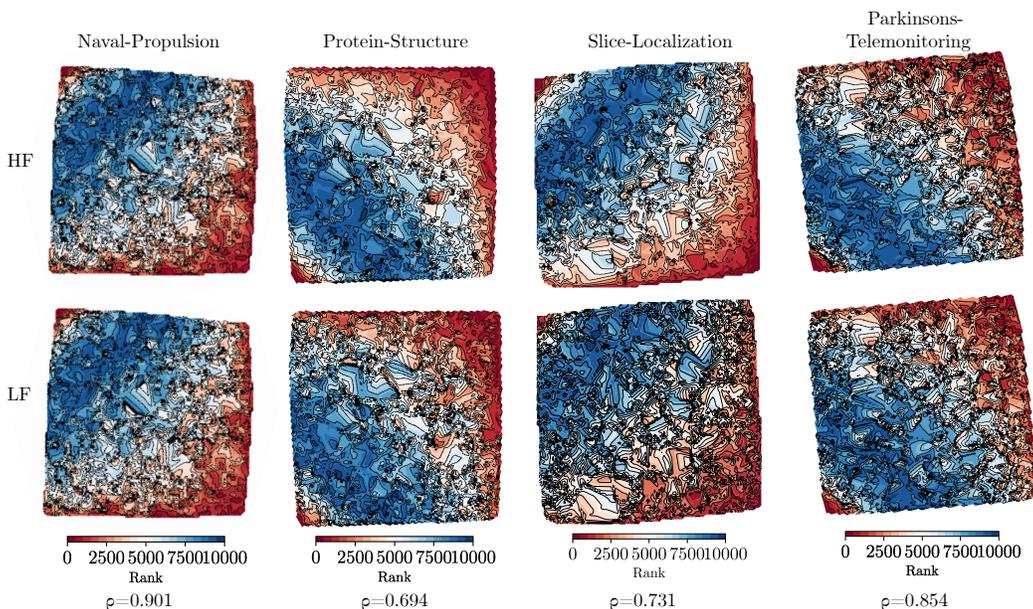


Figure 33: 2D visualization of landscapes of FCNet at high and low fidelities.

### F.4 PEER COMPARISON ON HYPERPARAMETER OPTIMIZATION FOR FINE-TUNING PRETRAINED IMAGE CLASSIFICATION MODELS

In this section, we additionally adopt the hyperparameter optimization from (Pineda-Arango et al., 2024) for fine-tuning pretrained image classification models on different datasets. A total of 20 problems are used to evaluate the performance of the algorithms, and the results are presented in Table 7. It can be observed that `textttLamda+BOHB` achieves the best performance across all problems. The overall ranking of the algorithms during the optimization process is illustrated in the Figure 38, showing that `Lamda+BOHB` consistently ranks first after consuming a portion of the resources. The convergence curves in Figure 39 further highlight its superiority in the second phase,

1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889  
1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922

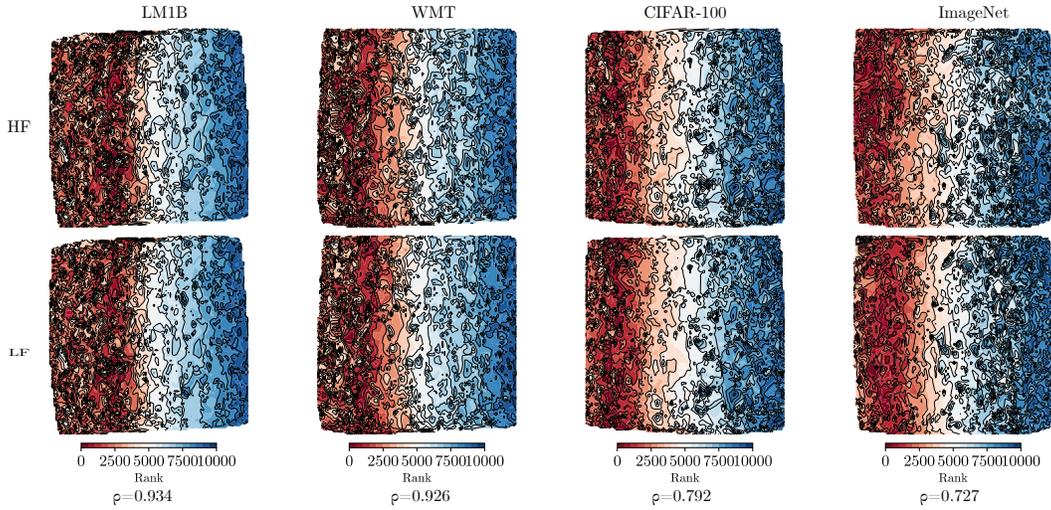


Figure 34: 2D visualization of landscapes of PD1 at high and low fidelities.

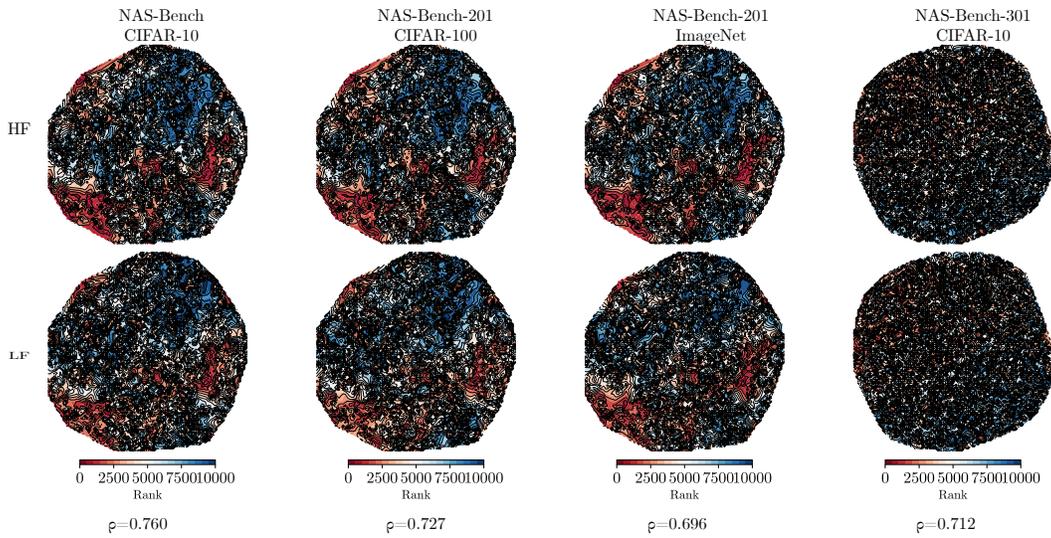


Figure 35: 2D visualization of landscapes of NAS-Banch-201 and 301 at high and low fidelities.

1923  
 1924  
 1925  
 1926  
 1927  
 1928  
 1929  
 1930  
 1931  
 1932  
 1933  
 1934  
 1935  
 1936  
 1937  
 1938  
 1939  
 1940  
 1941  
 1942  
 1943  
 1944  
 1945  
 1946  
 1947  
 1948  
 1949  
 1950  
 1951  
 1952  
 1953  
 1954  
 1955  
 1956  
 1957  
 1958  
 1959  
 1960  
 1961  
 1962  
 1963  
 1964  
 1965  
 1966  
 1967  
 1968  
 1969  
 1970  
 1971  
 1972  
 1973  
 1974  
 1975  
 1976

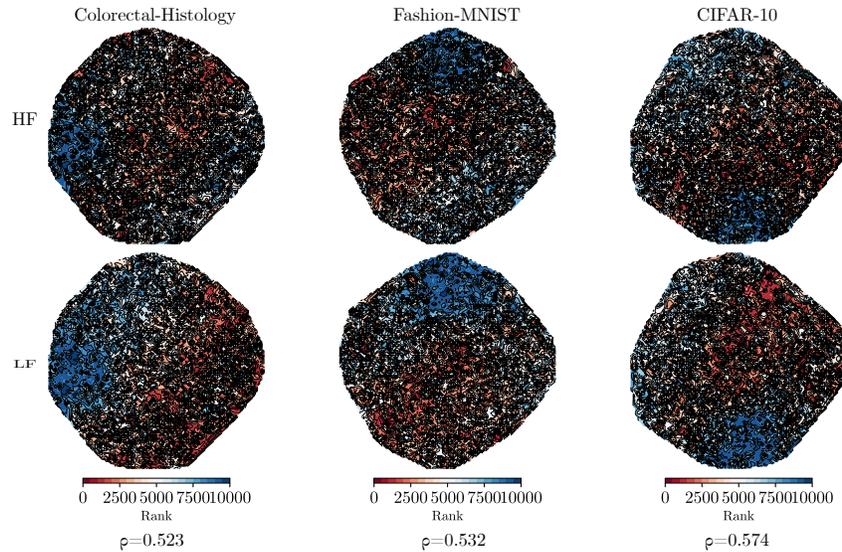


Figure 36: 2D visualization of landscapes of JAHS at high and low fidelites.

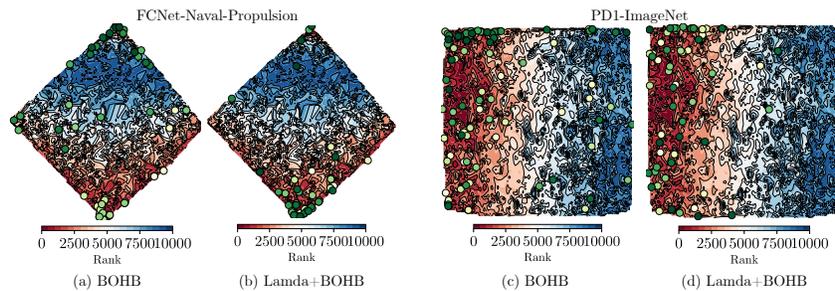
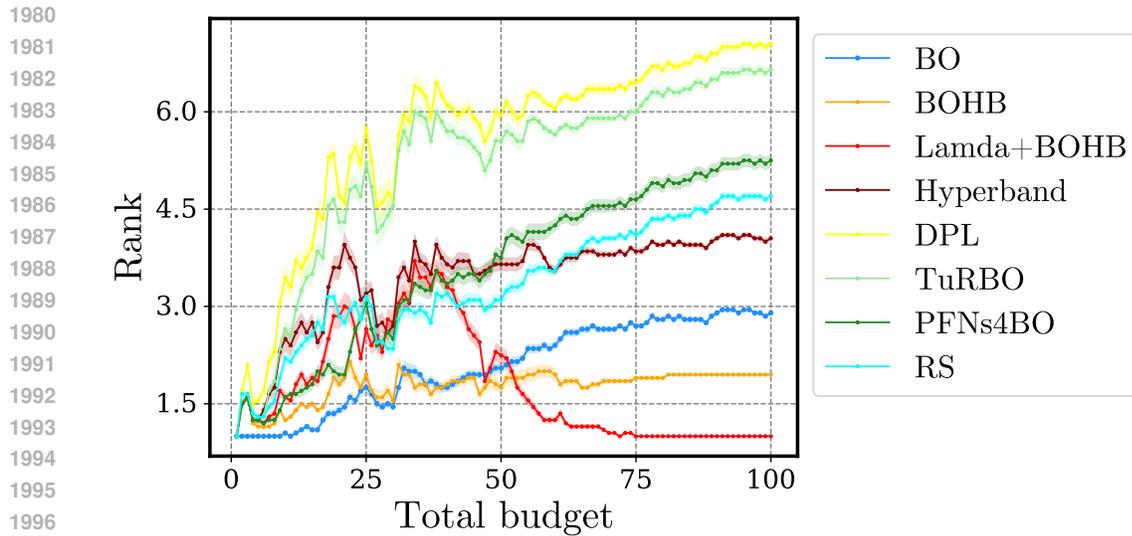


Figure 37: 2D visualization of the sampling points during the optimization process of Lamda+MUMBO and MUMBO on FCNet-Naval-Propulsion and PD1-ImageNet. The contour represents the dimensionality-reduced landscape of the HF problem, while the points indicate the HF samples collected during optimization. The color gradient from green to yellow indicates the order of sampling, with yellow representing later sampling stages.

1977 where it quickly identifies high-quality solutions and accelerates the optimization process. The ex-  
 1978 perimental results further demonstrate that Lamda consistently enhances BOHB.  
 1979



1997 Figure 38: Comparing average relative ranks of peer algorithms across 20 HPO tasks for fine-tuning  
 1998 pretrained image classification models.  
 1999

2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030

2031  
 2032  
 2033  
 2034  
 2035  
 2036  
 2037  
 2038  
 2039  
 2040  
 2041  
 2042  
 2043  
 2044  
 2045  
 2046  
 2047  
 2048  
 2049  
 2050  
 2051  
 2052  
 2053  
 2054  
 2055  
 2056  
 2057  
 2058  
 2059  
 2060  
 2061  
 2062  
 2063  
 2064  
 2065  
 2066  
 2067  
 2068  
 2069  
 2070  
 2071  
 2072  
 2073  
 2074  
 2075  
 2076  
 2077  
 2078  
 2079  
 2080  
 2081  
 2082  
 2083  
 2084

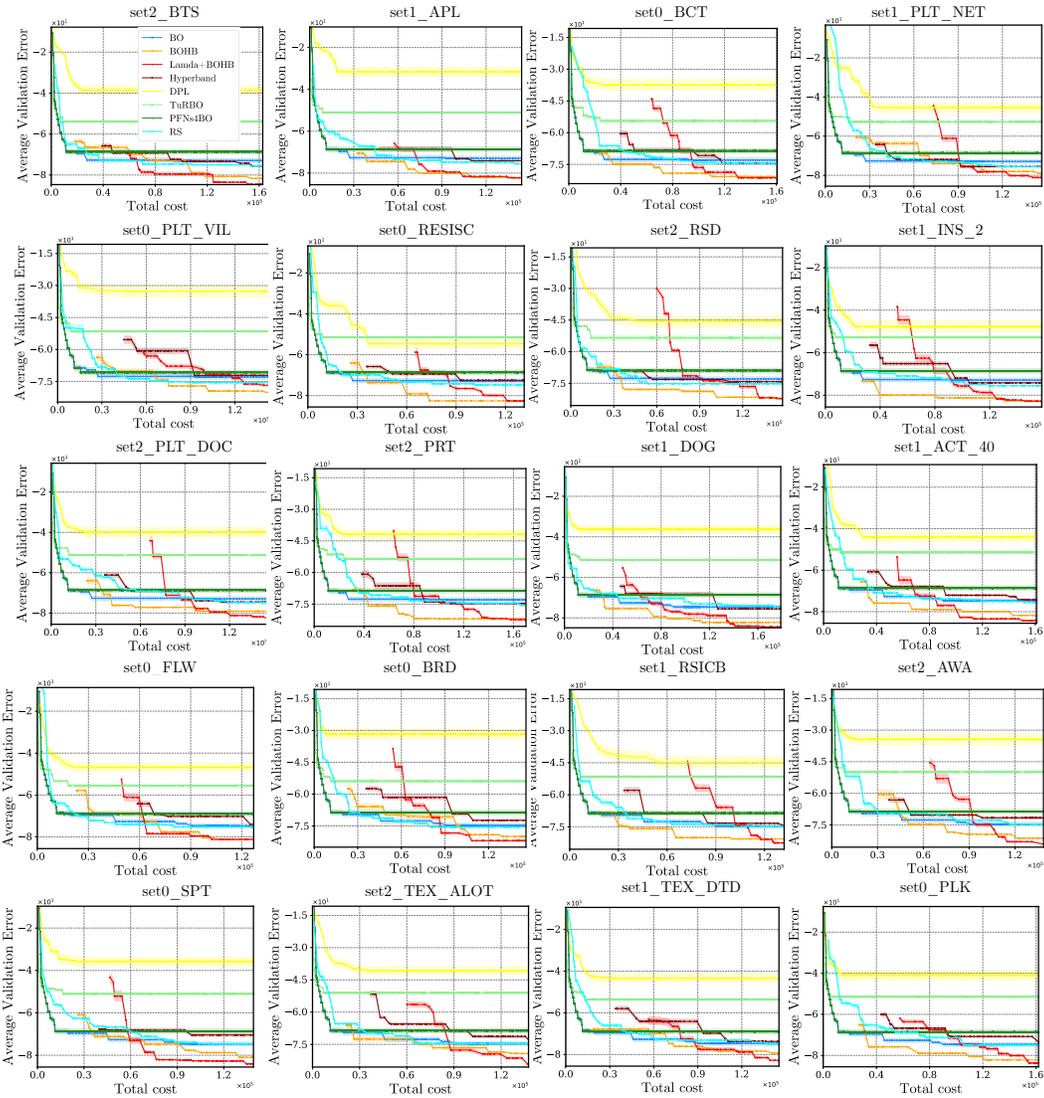


Figure 39: Validation error observed in tuning 20 HPO tasks for fine-tuning pretrained image classification models, comparing peer algorithms.

2085  
2086  
2087  
2088  
2089  
2090  
2091  
2092  
2093  
2094  
2095  
2096  
2097  
2098  
2099  
2100  
2101  
2102  
2103  
2104  
2105  
2106  
2107  
2108  
2109  
2110  
2111  
2112  
2113  
2114  
2115  
2116  
2117

Table 7: Comparing peer algorithms’ final validation errors of the current incumbent at 100 HPO tasks for fine-tuning pretrained image classification models). Runs are averaged over 31 seeds. (The bolded part indicates: ”under the Wilcoxon rank-sum test, methods incorporating priors significantly outperform their baselines.”)

Data Name	BO	BOHB	Lambda+BOHB	Hyperband	DPL	TuRBO	PFNs+4BO	RS
mtlhm_extended.set2.BTS	-7.295e+01(1.282e+01)	-8.256e+01(0.000e+00)	<b>-8.466e+01(2.709e+00)</b>	-7.556e+01(0.000e+00)	-3.857e+01(2.394e+01)	-5.392e+01(1.042e+01)	-6.862e+01(1.273e+01)	-7.556e+01(0.000e+00)
mtlhm_extended.set1.APL	-7.295e+01(1.282e+01)	-8.256e+01(0.000e+00)	<b>-8.556e+01(0.000e+00)</b>	-7.556e+01(0.000e+00)	-3.164e+01(2.416e+01)	-5.105e+01(7.422e+00)	-6.862e+01(1.273e+01)	-7.556e+01(0.000e+00)
mtlhm_extended.set0.BCT	-7.295e+01(1.282e+01)	-8.256e+01(0.000e+00)	<b>-8.409e+01(3.032e+00)</b>	-7.556e+01(0.000e+00)	-3.747e+01(2.900e+01)	-5.446e+01(1.274e+01)	-6.862e+01(1.273e+01)	-7.455e+01(3.034e+00)
mtlhm_extended.set1.PLT.NET	-7.295e+01(1.282e+01)	-8.166e+01(2.709e+00)	<b>-8.556e+01(0.000e+00)</b>	-7.556e+01(0.000e+00)	-4.541e+01(2.757e+01)	-5.270e+01(8.711e+00)	-6.862e+01(1.273e+01)	-7.556e+01(0.000e+00)
mtlhm_extended.set0.PLT.VIL	-7.295e+01(1.282e+01)	-8.256e+01(0.000e+00)	<b>-8.466e+01(2.709e+00)</b>	-7.556e+01(0.000e+00)	-3.273e+01(2.761e+01)	-5.142e+01(8.522e+00)	-7.078e+01(1.197e+01)	-7.556e+01(0.000e+00)
mtlhm_extended.set0.RESISC	-7.295e+01(1.282e+01)	-8.256e+01(0.000e+00)	<b>-8.556e+01(0.000e+00)</b>	-7.556e+01(0.000e+00)	-5.465e+01(2.450e+01)	-5.142e+01(8.522e+00)	-6.862e+01(1.273e+01)	-7.478e+01(2.349e+00)
mtlhm_extended.set2.RSD	-7.295e+01(1.282e+01)	-8.256e+01(0.000e+00)	<b>-8.556e+01(0.000e+00)</b>	-7.556e+01(0.000e+00)	-4.570e+01(2.934e+01)	-5.342e+01(1.467e+01)	-6.896e+01(1.288e+01)	-7.523e+01(9.984e-01)
mtlhm_extended.set1.INS.2	-7.295e+01(1.282e+01)	-8.256e+01(0.000e+00)	<b>-8.556e+01(0.000e+00)</b>	-7.556e+01(0.000e+00)	-4.778e+01(2.946e+01)	-5.270e+01(8.980e+00)	-6.862e+01(1.273e+01)	-7.556e+01(0.000e+00)
mtlhm_extended.set2.PLT.DOC	-7.295e+01(1.282e+01)	-8.223e+01(9.984e-01)	<b>-8.556e+01(0.000e+00)</b>	-7.556e+01(0.000e+00)	-3.996e+01(2.678e+01)	-5.131e+01(8.203e+00)	-6.862e+01(1.273e+01)	-7.478e+01(2.349e+00)
mtlhm_extended.set2.PRT	-7.295e+01(1.282e+01)	-8.256e+01(0.000e+00)	<b>-8.556e+01(0.000e+00)</b>	-7.556e+01(0.000e+00)	-4.190e+01(2.468e+01)	-5.356e+01(1.065e+01)	-6.862e+01(1.273e+01)	-7.556e+01(0.000e+00)
mtlhm_extended.set1.DOG	-7.467e+01(1.215e+01)	-8.256e+01(0.000e+00)	<b>-8.489e+01(1.331e+00)</b>	-7.556e+01(0.000e+00)	-3.633e+01(2.196e+01)	-5.142e+01(8.522e+00)	-6.862e+01(1.273e+01)	-7.462e+01(2.825e+00)
mtlhm_extended.set1.ACT.40	-7.467e+01(1.215e+01)	-8.256e+01(0.000e+00)	<b>-8.556e+01(0.000e+00)</b>	-7.556e+01(0.000e+00)	-4.404e+01(2.534e+01)	-5.142e+01(8.392e+00)	-6.862e+01(1.273e+01)	-7.556e+01(0.000e+00)
mtlhm_extended.set0.FLW	-7.467e+01(1.215e+01)	-8.256e+01(0.000e+00)	<b>-8.556e+01(0.000e+00)</b>	-7.556e+01(0.000e+00)	-4.677e+01(2.420e+01)	-5.547e+01(1.106e+01)	-6.896e+01(1.288e+01)	-7.556e+01(0.000e+00)
mtlhm_extended.set0.BRD	-7.467e+01(1.215e+01)	-8.256e+01(0.000e+00)	<b>-8.500e+01(1.694e+00)</b>	-7.556e+01(0.000e+00)	-3.177e+01(2.669e+01)	-5.392e+01(1.042e+01)	-6.862e+01(1.273e+01)	-7.556e+01(0.000e+00)
mtlhm_extended.set1.RSICB	-7.467e+01(1.215e+01)	-8.256e+01(0.000e+00)	<b>-8.556e+01(0.000e+00)</b>	-7.556e+01(0.000e+00)	-4.478e+01(2.931e+01)	-5.146e+01(8.506e+00)	-6.862e+01(1.273e+01)	-7.500e+01(1.694e+00)
mtlhm_extended.set2.AWA	-7.467e+01(1.215e+01)	-8.256e+01(0.000e+00)	<b>-8.556e+01(0.000e+00)</b>	-7.556e+01(0.000e+00)	-3.442e+01(2.638e+01)	-4.993e+01(9.759e+00)	-6.862e+01(1.273e+01)	-7.466e+01(2.709e+00)
mtlhm_extended.set0.SPT	-7.467e+01(1.215e+01)	-8.178e+01(2.349e+00)	<b>-8.556e+01(0.000e+00)</b>	-7.523e+01(9.984e-01)	-3.574e+01(2.475e+01)	-5.105e+01(7.422e+00)	-6.862e+01(1.273e+01)	-7.462e+01(2.825e+00)
mtlhm_extended.set2.TEX.ALOT	-7.467e+01(1.215e+01)	-8.256e+01(0.000e+00)	<b>-8.556e+01(0.000e+00)</b>	-7.556e+01(0.000e+00)	-4.067e+01(2.077e+01)	-5.105e+01(7.422e+00)	-6.862e+01(1.273e+01)	-7.462e+01(2.825e+00)
mtlhm_extended.set1.TEX.DTD	-7.467e+01(1.215e+01)	-8.223e+01(9.984e-01)	<b>-8.556e+01(0.000e+00)</b>	-7.556e+01(0.000e+00)	-4.332e+01(2.276e+01)	-5.356e+01(1.065e+01)	-6.896e+01(1.288e+01)	-7.500e+01(1.694e+00)
mtlhm_extended.set0.PLK	-7.467e+01(1.215e+01)	-8.256e+01(0.000e+00)	<b>-8.556e+01(0.000e+00)</b>	-7.556e+01(0.000e+00)	-4.097e+01(2.802e+01)	-5.146e+01(8.506e+00)	-6.862e+01(1.273e+01)	-7.500e+01(1.694e+00)

## G. ALGORITHM DETAILS

In this section, we outline the workflow of our Lamda framework, which operates in two phases: the first-phase low-fidelity search (dubbed as Lamda-1) and the second-phase optimization (dubbed as Lamda-2). The overall workflow is depicted in Algorithm 1.

The first phase, Lamda-1, focuses on identifying promising regions in the search space using LF evaluations. The pseudocode for Lamda-1 is provided in Algorithm 2. In particular, the sampling strategy in Lamda-1 is algorithm-agnostic and can be incorporated with most HPO algorithms.

The second phase, Lamda-2, allows the use of different algorithms to guide the search. We provide multiple pseudocode to demonstrate how Lamda-2 can be adapted to various algorithms, including PriorBand, BOHB, MUMBO, vanilla BO, and random search. The details for each integration are outlined below, with modifications from the original algorithms highlighted in red:

- For PriorBand, Lamda-2 replaces the expert prior  $p_\pi(\mathbf{x})$  from (Mallik et al., 2023) with the learned prior  $\varphi_{\text{pro}}(\mathbf{x})$  obtained in Lamda-1.
- For BOHB, Lamda-2 modifies the uniform sampling distribution in Steps 7 and 8 of Algorithm 3 into the incumbent distribution determined by  $\tilde{\varphi}_{\text{pro}}(\mathbf{x})$ .
- For MUMBO and BO, Lamda-2 combines  $\tilde{\varphi}_{\text{pro}}(\mathbf{x})$  with their acquisition functions. The pseudocode for these adaptations is shown in Algorithm 4 and Algorithm 5, respectively.
- For Random Search, Lamda-2 replaces the uniform sampling strategy by the incumbent sampling strategy defined by  $\tilde{\varphi}_{\text{pro}}(\mathbf{x})$ , as illustrated in Algorithm 6.

---

### Algorithm 3: Pseudocode for sampling in Lamda+BOHB

---

**Input:** Observations  $D$ , fraction of random runs  $\rho$ , percentile  $q$ , number of samples  $N_s$ , minimum number of points  $N_{\min}$  to build GP models, and bandwidth factor  $b_w$

**Output:** next configuration to evaluate

```

1 Initialize  $b \leftarrow \arg \max\{D_b : |D_b| \geq N_{\min} + 2\}$ ,  $\tilde{\rho} \leftarrow \text{Rand}(0, 1)$ ;
2 if  $\tilde{\rho} < \rho$  or  $b = \emptyset$  then
3   return randomly sampled configuration;
4 else
5   Compute  $l(\mathbf{x})$  and  $g(\mathbf{x})$  as Eqs. (2) and (3) in Falkner et al. (2018);
6   Draw  $N_s$  configurations according to  $\tilde{\varphi}_{\text{pro}}(\mathbf{x})$  in equation (6);
7   return configuration with highest ratio  $l(\mathbf{x})/g(\mathbf{x})$ 

```

---

2172  
2173  
2174  
2175  
2176  
2177  
2178  
2179  
2180  
2181  
2182  
2183  
2184  
2185  
2186  
2187  
2188  
2189  
2190  
2191  
2192  
2193  
2194  
2195  
2196  
2197  
2198  
2199  
2200  
2201  
2202  
2203  
2204  
2205  
2206  
2207  
2208  
2209  
2210  
2211  
2212  
2213  
2214  
2215  
2216  
2217  
2218  
2219  
2220  
2221  
2222  
2223  
2224  
2225

**Algorithm 4:** Second-phase search with BO

**Input:** Input space  $\mathcal{X}$ ,  $\varphi_{\text{pro}}(\mathbf{x})$ ,  $w$ ,  $N$  solution for the initial design of GPs, budget  $\Lambda_r$ , fidelity level  $h$ , budget function  $\lambda_z$ .

**Output:** Optimized design  $x^*$ .

/\* Initialization \*/;

- 1 Sample  $\{\mathbf{x}^i\}_{i=1}^n$  from distribution given by  $\varphi_{\text{pro}}(\mathbf{x})$ ;
- 2  $y^i \leftarrow f_h(\mathbf{x}^i) + \epsilon^i$ , where  $\epsilon^i \sim \mathcal{N}(0, \sigma^2)$ ;  $\lambda^i \leftarrow \lambda_z(\mathbf{x}^i, h)$ ;
- 3  $\Lambda_r \leftarrow \Lambda_r - \sum_{i=1}^n \lambda^i$ ;
- 4  $D \leftarrow \{(\mathbf{x}^i, y^i)\}_{i=1}^n$ ;
- 5 **while**  $\Lambda_r > 0$  **do**
- 6      $\varphi(\mathbf{x}) \leftarrow p(\mathbf{x}|D)$ ;
- 7      $\tilde{\varphi}_{\text{pro}}(\mathbf{x}) \leftarrow (1-w) \cdot \varphi(\mathbf{x}) + w \cdot \varphi_{\text{pro}}(\mathbf{x})$ ;
- 8      $\mathbf{x}^{n+1} \leftarrow \arg \max_{\mathbf{x} \in \mathcal{X}} \tilde{\varphi}_{\text{pro}}(\mathbf{x}) \text{AF}(\mathbf{x}, D)$ ;
- 9      $y^{n+1} \leftarrow f_h(\mathbf{x}^{n+1}) + \epsilon$ ;
- 10     Update  $D \leftarrow D \cup \{(\mathbf{x}^{n+1}, y^{n+1})\}$ ;
- 11      $\Lambda_r \leftarrow \Lambda_r - \lambda_z(\mathbf{x}^{n+1}, h)$ ;
- 12      $n \leftarrow n + 1$ ;
- 13 **return**  $x^* \leftarrow \arg \min_{(\mathbf{x}^i, y^i) \in D} y^i$

**Algorithm 5:** Second-phase search with MUMBO

**Input:** Input space  $\mathcal{X}$ , prior obtained in the first phase  $\varphi_{\text{pro}}(\mathbf{x})$ , prior confidence parameter  $w$ , size  $n$  of the initial design, budget for first phase  $\Lambda_r$ .

**Output:** Optimized design  $x^*$ .

- 1 Sample  $\{\mathbf{x}^i\}_{i=1}^n \sim \varphi_{\text{pro}}(\mathbf{x})$  and randomly assign fidelity levels  $\{z^i\}_{i=1}^n$  with  $z^i \sim \text{Uniform}(\{\ell, \ell + 1, \dots, h\})$ ;
- 2 Compute  $y^i \leftarrow f_z(\mathbf{x}^i, z^i) + \epsilon^i$ , where  $\epsilon^i \sim \mathcal{N}(0, \sigma^2)$ ;  $\lambda^i \leftarrow \lambda_z(\mathbf{x}^i, z^i)$ ;
- 3  $\Lambda_r \leftarrow \Lambda_r - \sum_{i=1}^n \lambda^i$ ;
- 4 Initialize  $D \leftarrow \{(\mathbf{x}^i, z^i), y^i\}_{i=1}^n$ ;
- 5 **while**  $\Lambda_r > 0$  **do**
- 6     Fit GP to the collected observations  $D$ ,  $\varphi(\mathbf{x}) \leftarrow p(\mathbf{x}|D)$ ;
- 7     Simulate  $N$  samples of  $g^*|D$ ;
- 8     Prepare  $\alpha_{n-1}^{\text{MUMBO}}(\mathbf{x}, z)$  as given by Eq. (5) in Moss et al. (2020);
- 9     Update  $\tilde{\varphi}_{\text{pro}}(\mathbf{x}) \leftarrow (1-w) \cdot \varphi(\mathbf{x}) + w \cdot \varphi_{\text{pro}}(\mathbf{x})$ ;
- 10     Find the next point and fidelity to query
- 11      $(\mathbf{x}^{n+1}, z^{n+1}) \leftarrow \arg \max_{(\mathbf{x}, z)} \tilde{\varphi}_{\text{pro}}(\mathbf{x}) \frac{\alpha_{n-1}^{\text{MUMBO}}(\mathbf{x}, z)}{\lambda_z(\mathbf{x}, z)}$
- 12     Collect the new evaluation  $y^{n+1} \leftarrow f_z(\mathbf{x}^{n+1}, z^{n+1}) + \epsilon^{n+1}$ ,  $\epsilon^{n+1} \sim \mathcal{N}(0, \sigma^2)$ ;
- 13     Append new evaluation to observation set  $D \leftarrow D \cup \{(\mathbf{x}^{n+1}, z^{n+1}), y^{n+1}\}$ ;
- 14     Update spent budget  $\Lambda_r \leftarrow \Lambda_r - \lambda_z(\mathbf{x}^{n+1}, z^{n+1})$ ;
- 15 **return**  $x^* \leftarrow \arg \min_{\{(\mathbf{x}^i, z^i), y^i\} \in D, z^i = h} y^i$

2226  
 2227  
 2228  
 2229  
 2230  
 2231  
 2232  
 2233  
 2234  
 2235  
 2236  
 2237  
 2238  
 2239  
 2240  
 2241  
 2242  
 2243  
 2244  
 2245  
 2246  
 2247  
 2248  
 2249  
 2250  
 2251  
 2252  
 2253  
 2254  
 2255  
 2256  
 2257  
 2258  
 2259  
 2260  
 2261  
 2262  
 2263  
 2264  
 2265  
 2266  
 2267  
 2268  
 2269  
 2270  
 2271  
 2272  
 2273  
 2274  
 2275  
 2276  
 2277  
 2278  
 2279

---

**Algorithm 6:** Second-phase search with Random Search

---

**Input:** Input space  $\mathcal{X}$ , prior obtained in the first phase  $\varphi_{\text{pro}}(\mathbf{x})$ , prior confidence parameter  $w$ , budget for first phase  $\Lambda_r$ , uniform distribution  $p_U$ .

**Output:** Optimized design  $x^*$ .

```

1  $\varphi(\mathbf{x}) \leftarrow p_U(\mathbf{x});$ 
2 while  $\Lambda_r > 0$  do
3    $\tilde{\varphi}_{\text{pro}}(\mathbf{x}) \leftarrow (1 - w) \cdot \varphi(\mathbf{x}) + w \cdot \varphi_{\text{pro}}(\mathbf{x});$ 
4   Sample  $\mathbf{x}^{n+1} \sim \tilde{\varphi}_{\text{pro}}(\mathbf{x});$ 
5    $y^{n+1} \leftarrow f_h(\mathbf{x}^{n+1}) + \epsilon;$ 
6   Update  $D \leftarrow D \cup \{(\mathbf{x}^{n+1}, y^{n+1})\};$ 
7    $\Lambda_r \leftarrow \Lambda_r - \lambda_z(\mathbf{x}^{n+1}, h);$ 
8 return  $x^* \leftarrow \arg \min_{(\mathbf{x}^i, y^i) \in D} y^i$ 

```

---