

# DYNAMIC PRIORS IN BAYESIAN OPTIMIZATION FOR HYPERPARAMETER OPTIMIZATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

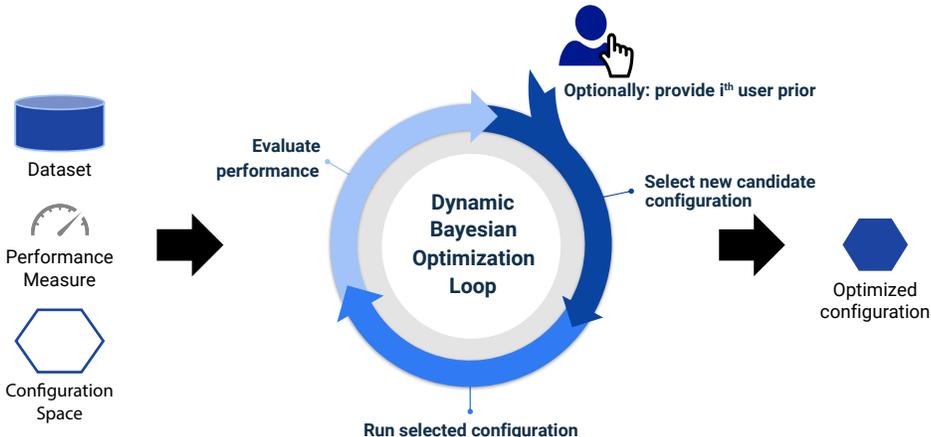
Hyperparameter optimization (HPO), for example, based on Bayesian optimization (BO), supports users in designing models well-suited for a given dataset. HPO has proven its effectiveness on several applications, ranging from classical machine learning for tabular data to deep neural networks for computer vision and transformers for natural language processing. However, HPO still sometimes lacks acceptance by machine learning experts due to its black-box nature and limited user control. Addressing this, first approaches have been proposed to initialize BO methods with expert knowledge. However, these approaches do not allow for online steering *during* the optimization process. In this paper, we introduce a novel method that enables repeated interventions to steer BO via user input, specifying expert knowledge and user preferences at runtime of the HPO process in the form of prior distributions. To this end, we generalize an existing method,  $\pi$ BO, preserving theoretical guarantees. We also introduce a misleading prior detection scheme, which allows protection against harmful user inputs. In our experimental evaluation, we demonstrate that our method can effectively incorporate multiple priors, leveraging informative priors, whereas misleading priors are reliably rejected or overcome. Thereby, we achieve competitiveness to unperturbed BO.

## 1 INTRODUCTION

Hyperparameter optimization (HPO) is concerned with automatically optimizing the hyperparameters of machine learning algorithms for a given task. A task often consists of a dataset, a configuration space, and a performance measure (Hutter et al., 2019). HPO is effective on classical machine learning (Eggenesperger et al., 2021; Bansal et al., 2022; Pfisterer et al., 2022) as well as on deep neural networks and transformers for computer-vision and natural language processing (Müller et al., 2023; Wang et al., 2024; Rakotoarison et al., 2024; Pineda Arango et al., 2024). Nevertheless, users often disregard HPO in favor of maintaining manual hyperparameter configuration control (Bouthillier & Varoquaux, 2020; Van der Blom et al., 2021; Kannengießer et al., 2025). The reluctance is amplified in the case of knowledgeable experts (Kannengießer et al., 2025). At the same time, explainability methods for HPO (Wang et al., 2019; Sass et al., 2022; Zöller et al., 2023) seek to provide users with insights into the optimization process, often motivating adaptations to HPO. However, approaches that allow practitioners to steer the optimization process, instead of restarting it, remain largely underexplored. Nevertheless, they could tackle expert users’ perceived lack of control (Kannengießer et al., 2025), thus possibly furthering HPO methods’ acceptance.

Recently, initial steps have been taken toward collaborative automated machine learning (AutoML) approaches, subsuming HPO, in view of a more human-centered AutoML paradigm (Lindauer et al., 2024). More specifically, Hvarfner et al. (2022) and Mallik et al. (2023) propose user-centric interfaces for Bayesian optimization (BO) (Jones et al., 1998) and Hyperband (Li et al., 2017), respectively. Using explicit user priors on the location of the optimal configuration, the interfaces enable users to bias the optimization process. Hvarfner et al. (2024) extend this idea by building a general framework that enables users to encode properties of the optimized function, such as the maximum achievable performance, into the optimization process. These approaches show that informative priors can significantly enhance performance, while misleading priors do not break the optimization process, and the associated performance deterioration can be recovered from in the long run. Although these developments represent progress toward human-centered BO, the proposed level of interaction remains limited, as the approaches lack the possibility of user-enabled online control.

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068



069  
070  
071  
072  
073  
074

Figure 1: Overview of the proposed dynamic Bayesian optimization (DynaBO) method. Provided a dataset, a performance measure, a configuration space, and an optional initial prior, the loop iteratively selects new hyperparameter configurations. At each step, a candidate configuration is evaluated, and it is assessed. The process continues until an optimized configuration is identified. The framework allows users to steer the optimization process by dynamically adding priors at runtime.

075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085

In this paper, we propose an approach for dynamic Bayesian optimization, dubbed DynaBO, to advance the interactivity of BO-based HPO by enabling, but not requiring, users to inject knowledge on potentially well-suited regions of the hyperparameter configuration space. This user knowledge is provided in the form of user priors and can be injected at any point during the process, as illustrated in Figure 1. When multiple priors are supplied simultaneously or in close succession, they are stacked and their effects are combined. This (i) allows users to continuously steer the optimization process, and (ii) enables them to use HPO in the typical rapid-prototyping workflows of machine learning practitioners (Studer et al., 2021) where users oversee and direct the construction of the continuously updated model. These advantages are enhanced by an informed user, e.g., through utilizing explainability methods for HPO (Hutter et al., 2014; Moosbauer et al., 2021; Segel et al., 2023), which provide actionable insights.

086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096

Concretely, DynaBO generalizes the work of Hvarfner et al. (2022) from a single to multiple priors provided by users over time. Importantly, DynaBO maintains  $\pi$ BO’s speedup for informative priors, and retains its theoretical convergence properties despite misleading priors. To tackle the nevertheless significant anytime performance deterioration due to misleading priors, we propose an additional mechanism for their detection and potential rejection. This mechanism is based on an assessment of the provided priors using acquisition functions and the already trained surrogate model, therefore requiring no significant overhead. Our experimental evaluation on a multitude of real-world HPO tasks reveals substantial speedups of DynaBO over vanilla BO and  $\pi$ BO (Hvarfner et al., 2022) for informative priors. At the same time, while  $\pi$ BO suffers from performance deterioration due to a single misleading prior, our prior rejection mechanism enables DynaBO to achieve competitiveness with vanilla BO even when exposed to misleading priors repeatedly.

097  
098  
099

## 2 HYPERPARAMETER OPTIMIZATION

100  
101  
102  
103  
104  
105  
106  
107

Hyperparameter optimization (HPO) is concerned with finding a suitable hyperparameter configuration of a learner for a given task, typically defined by a labeled dataset  $D$  and a loss function  $\ell$  that evaluates predictive quality (Bischl et al., 2023). Given an input space  $\mathcal{X}$  and an output (label) space  $\mathcal{Y}$ ,  $D = \{(x^{(k)}, y^{(k)})\}_{k=1}^N$  is a finite sample drawn from a joint probability distribution  $P(\mathcal{X}, \mathcal{Y})$ .  $D$  is split into a training set  $D_T$  and validation set  $D_V$ . We consider a loss function  $\ell : \mathcal{Y} \times P(\mathcal{Y}) \rightarrow \mathbb{R}$  that evaluates the quality of predictions in terms of the true label and a probability distribution over  $\mathcal{Y}$ . Given a hyperparameter configuration  $\lambda \in \Lambda$  from a configuration space  $\Lambda$ , and dataset  $D_T$ , the learner produces a hypothesis  $h_{\lambda, D_T} \in \mathcal{H} := \{h \mid h : \mathcal{X} \rightarrow P(\mathcal{Y})\}$ , that is, a mapping from inputs to probability distributions over  $\mathcal{Y}$ .

Since the hyperparameter configuration  $\lambda$  influences both the hypothesis space  $\mathcal{H}$  and the inductive behavior of the learner, selecting a well-suited hyperparameter configuration with respect to the dataset and the loss function is key to achieving peak performance. The objective of HPO is to find a configuration  $\lambda^*$  that leads to a hypothesis based on  $D_T$  with strong generalization to  $D_V$ :

$$\lambda^* \in \arg \min_{\lambda \in \Lambda} f(\lambda) := \arg \min_{\lambda \in \Lambda} \mathbb{E}_{(D_T, D_V) \sim D} \left[ \frac{1}{|D_V|} \sum_{(x, y) \in D_V} \ell(y, h_{\lambda, D_T}(x)) \right]. \quad (1)$$

While the HPO problem can be tackled naïvely via grid search or random search (Bergstra & Bengio, 2012), more sophisticated techniques, such as Bayesian optimization, are recommended today due to their higher efficiency and effectiveness (Turner et al., 2021; Bischl et al., 2023).

**Bayesian Optimization (BO)** (Moćkus, 1975; Jones et al., 1998; Shahriari et al., 2016) is a model-based sequential optimization technique widely-used for sample-efficient HPO (Snoek et al., 2012; Falkner et al., 2018; Cowen-Rivers et al., 2022; Makarova et al., 2022; Bischl et al., 2023). BO is particularly well-suited for optimizing black-box functions  $f$ , which are costly to evaluate and have neither a closed-form solution nor any gradient information available.

The BO process begins with an initial design that aims to cover the hyperparameter configuration space diversely. It then proceeds, alternating between two main steps: (i) fitting a probabilistic surrogate model  $\hat{f}$  to the set of observed evaluations, and (ii) selecting the next configuration

$$\lambda^{t+1} \in \arg \max_{\lambda \in \Lambda} \alpha_{\hat{f}}(\lambda) \quad (2)$$

where the acquisition function  $\alpha_{\hat{f}} : \Lambda \rightarrow \mathbb{R}$  quantifies the utility of candidate points in balancing exploration and exploitation. In practice, and in line with standard GP-based BO (Rasmussen & Williams, 2006), the surrogate is typically trained on normalized observations (e.g., scaled to the  $[0, 1]$  interval). This normalization stabilizes GP hyperparameter estimation and ensures that commonly used acquisition functions yield non-negative values by construction. A prominent example is Expected Improvement (EI) (Jones et al., 1998), which selects points expected to yield improvements over the best previously observed configuration, also called the incumbent. Formally, let  $\hat{\lambda}$  denote the current incumbent of the black-box function, then EI at a hyperparameter configuration  $\lambda$  is defined as

$$\alpha_{\hat{f}}^{EI}(\lambda) := \mathbb{E} \left[ \max \left( f(\hat{\lambda}) - \hat{f}(\lambda), 0 \right) \right],$$

where  $\hat{f}(\lambda)$  is treated as a random variable representing the probability distribution at  $\lambda$  as modeled by  $\hat{f}$ . In HPO, the black-box function  $f$  typically corresponds to the empirical generalization error, estimating the expected loss as in Eq. 1, for example, via hold-out validation or cross-validation.

### 3 RELATED WORK

Approaches related to our work can be broadly grouped into (1) data-driven priors, (2) explicit user generated priors, and (3) interactive HPO frameworks.

**Data-Driven Priors** A substantial body of research leverages prior experience to configure the HPO process. This includes transfer learning across tasks (Swersky et al., 2013; Feurer et al., 2015; van Rijn & Hutter, 2018; Feurer et al., 2022), configuration space design (Perrone et al., 2019) and surrogate model configuration (Feurer et al., 2018). While effective, such approaches do not enable direct user interaction. Instead, they require users to trust automated knowledge transfer.

**Learning from User Priors** Another line of research explicitly incorporates user-specified beliefs. Bergstra et al. (2011) introduce fixed priors over the configuration space, whereas Souza et al. (2021b) estimate posterior-driven models, though both approaches are limited by their acquisition function compatibility. Ramachandran et al. (2020) warp the configuration space to emphasize promising regions, but this requires invertible priors and struggles with misleading ones. More recently, Hvarfner et al. (2022) propose  $\pi$ BO, which augments acquisition functions with priors in a flexible and robust way, whereas Mallik et al. (2023) extend HyperBand (Li et al., 2017) to balance

random sampling with user-defined priors. However, all of these methods restrict user input to the initialization phase and offer neither continuous control nor oversight over the optimization process.

**Interactive Hyperparameter Optimization** Moving beyond priors provided at the beginning of the optimization process, some recent works integrate users more directly into the optimization process. Xu et al. (2024) allow users to reject candidate evaluations, while Adachi et al. (2024) let users choose among proposed alternatives. Seng et al. (2025) propose an alternative to the common BO framework that is based on probabilistic circuits and does not rely on an acquisition function. They sample the use of the prior in each iteration based on a binomial distribution decaying over time. Priors can be given at any time, but in contrast to our approach, they override previous information entirely, and the approach suffers from extrapolation issues. Complementary to these, Chang et al. (2025) introduce LLINBO working with LLM-generated candidate configurations to augment BO, with rejection schemes for ensuring robustness. Likewise, cooperative design optimization (Niwa et al., 2025) explores natural language interfaces where LLMs propose configurations, optionally guided by user input. Their additional user studies demonstrate that such interaction mitigates over-tuning to local optima while maintaining user agency.

In contrast, our proposed method DynaBO enables users to provide priors at any time during the HPO process, while remaining model-agnostic and acquisition function-compatible. Unlike previous methods, DynaBO not only integrates repeated user input but also detects misleading priors and can safeguard against those, thus advancing toward a more collaborative paradigm for HPO.

## 4 DYNAMIC PRIORS IN BAYESIAN OPTIMIZATION

In this section, we present how dynamic priors can be leveraged in Bayesian optimization (BO), providing the first component of our proposed method, DynaBO. To this end, we first generalize prior-weighted acquisition functions from a single prior to stacking multiple priors (Section 4.1). Afterward, we devise a method to detect misleading priors and safeguard against them (Section 4.2).

### 4.1 PRIOR-WEIGHTED ACQUISITION FUNCTION

Following Hvarfner et al. (2022), we integrate user-provided prior information on the location of the optimum by weighting the original acquisition function with a prior distribution. Thereby, we enable the optimization process to incorporate external knowledge dynamically, which may be provided by the user or taken from some other source of knowledge.

Given an acquisition function  $\alpha$ , and a user-specified prior  $\pi : \Lambda \rightarrow (0, 1]$ , the next point to be evaluated with respect to  $f$  (e.g, as defined in Eq. 1) at time  $t$  is selected as follows:

$$\alpha_{\hat{f}}^{\pi\text{BO}}(\lambda) := \alpha_{\hat{f}}(\lambda) \cdot \pi(\lambda)^{\beta/t},$$

where  $\beta \in \mathbb{R}^+$  is a scaling hyperparameter. For  $t \rightarrow \infty$ , the weight induced by  $\pi$  converges to 1 independent of the configuration  $\lambda$  and  $\beta$ , that is, the effect of the prior diminishes over time.

In contrast to the work of Hvarfner et al. (2022), we suppose a finite sequence of user-specified priors  $\{\pi^{(m)}\}_{m=1}^M$  provided at times  $\{t^{(m)}\}_{m=1}^M$ ,  $t^{(1)} < \dots < t^{(M)} \leq T$ . We define a dynamically-adapted acquisition function  $\alpha_{\text{dyna}} : \Lambda \rightarrow \mathbb{R}$  by multiplying the sum of the priors:

$$\alpha_{\hat{f}}^{\text{dyna}}(\lambda) := \alpha_{\hat{f}}(\lambda) \cdot \sum_{m=1}^M \pi^{(m)}(\lambda)^{\beta/(t-t^{(m)})}$$

By stacking priors, we can incorporate the information given at different time steps, and the priors are faded individually based on their age; that is, older priors are considered less important, as shown in Figure 2. The flexibility in incorporating multiple priors sets DynaBO apart from  $\pi\text{BO}$  (Hvarfner et al., 2022) and Priorband (Mallik et al., 2023), which consider a single initial prior, as well as Seng et al. (2025)’s approach, which considers one prior at a time.

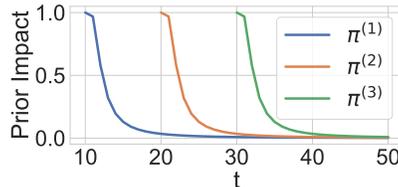


Figure 2: Acquisition function impact of priors  $\pi^1, \pi^2, \pi^3$ , provided at  $t = 10, 20$ , and  $30$ , with  $\pi(\lambda) = 0.5$ .

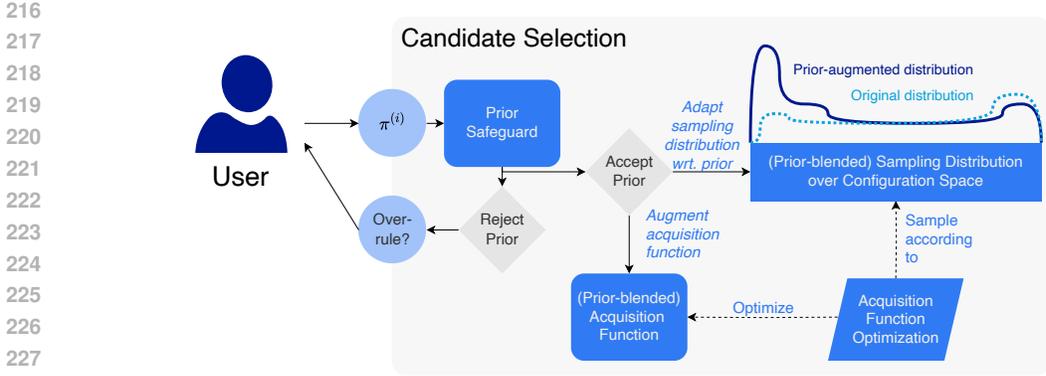


Figure 3: Illustration of the candidate selection process in DynaBO, incorporating a user-provided prior  $\pi^{(i)}$ . A safeguard mechanism evaluates the prior, determining whether to accept or reject it. If accepted, the candidate selection is biased by the prior; otherwise, the user can overrule the rejection.

**Acquisition Function Optimization** If user priors are peaked in a small area of the configuration space, the resulting acquisition function  $\alpha_{\text{dyna}}$  could be harder to optimize, see Equation (2). To ensure that the acquisition function optimizer covers both user-suggested regions and generally unexplored areas of the configuration space adequately, we adopt a modified version of Hutter et al. (2011)’s combined local and random search. To ensure that a sufficient number of candidates are sampled close to the peaked prior, we adapt the sampling of the starting points of the local search to the prior distributions. To this end, each prior  $\pi^{(m)}$  is assigned a weight  $\omega_m = e^{\phi(t-t^{(m)})}$ , which also decays with time. Then, a fraction  $\omega_m \cdot \min\{\sum_{i=1}^m w_i, 0.9\}$  of candidate configurations are sampled according to prior  $\pi^{(m)}$ . For more details, we refer to Appendix B.2.

## 4.2 REJECTING PRIORS

To safeguard against misleading priors, potentially slowing down DynaBO, we propose a mechanism to reject priors based on their feasibility to yield improvements given the surrogate model  $\hat{f}$  and an acquisition function  $\xi$ , which can but need not be the same as  $\alpha$ . For an overview of how the safeguard is embedded in DynaBO’s candidate selection routine, please refer to Figure 3. Specifically, we assess the promisingness of a prior  $\pi^{(m)}$  by comparing the potential of the suggested region against the region around the best configuration found previously, that is, the current incumbent  $\hat{\lambda}$ .

Sampling a number of candidate configurations from normal distributions centered around the incumbent and the prior (denoted by  $\mathcal{N}_{\hat{\lambda}}$  and  $\mathcal{N}_{\pi^{(m)}}$ , respectively), we compare the expected acquisition function values of both samples based on the surrogate model  $\hat{f}$ . We accept the prior if and only if the quality of the prior exceeds the quality of the incumbent area by a given threshold  $\tau$ :

$$\mathbb{E}_{\lambda \sim \mathcal{N}_{\pi^{(m)}}} [\xi_{\hat{f}}(\lambda)] - \mathbb{E}_{\lambda \sim \mathcal{N}_{\hat{\lambda}}} [\xi_{\hat{f}}(\lambda)] \geq \tau. \quad (3)$$

The higher  $\tau$  is set, the fewer priors are accepted by DynaBO, but also the more misleading priors can be removed. Setting  $\tau$  to a lower value lets DynaBO embrace user-provided priors more often, while making it more prone to misleading priors. For more details regarding our choice of  $\tau$  and the rejection of priors, refer to Appendix B.3.

In practice, the acquisition function  $\xi$  lives in the space of loss values, e.g., Lower Confidence Bound (LCB) or EI (Snoek et al., 2012). Since we want to reject priors based on their potential, we recommend utilizing LCB and setting  $\tau$  with respect to the user’s beliefs on the remaining optimization potential  $f(\hat{\lambda}) - f(\lambda^*)$ . In summary, this approach ensures that optimization resources are not wasted on poor suggestions, that is, regions which are already known to perform inferior to the incumbent region. We also expect that the prior rejection scheme would be used to warn against priors, but users would have the chance to overrule their rejection. For a sensitivity analysis on  $\tau$ , refer to Section 6.4.

## 270 5 THEORETICAL ANALYSIS

271  
272 In the theoretical analysis of DynaBO, we establish convergence and robustness properties under  
273 multiple dynamically provided user priors, and quantify the conditions under which informative  
274 priors yield accelerated convergence. Extending the convergence results of  $\pi$ BO (Hvarfner et al.,  
275 2022), we demonstrate that our approach retains the almost sure convergence behavior of BO even  
276 when given misleading priors. At the same time, our method is capable of leveraging informative  
277 priors effectively to accelerate convergence. Note that we assume a *finite* prior set, and the utilization  
278 of UCB as an acquisition function. All proofs are provided in Appendix A. We follow standard  
279 convergence results for BO (Srinivas et al., 2012) augmented by the dynamic prior influence.

280 **Almost Sure Convergence of DynaBO** We analyze the asymptotic behavior of DynaBO given  
281 assumptions on objective function regularity, finiteness and vanishing influence of priors, introduced  
282 formally in Appendix A, allowing priors to be selected dynamically and to vary in quality over time.

283 **Theorem 1** (Almost Sure Convergence of DynaBO). *Under the assumptions in Appendix A, the*  
284 *sequence of query points selected by DynaBO,  $\{\lambda_t\}_{t=1}^\infty \subset \Lambda$ , satisfies almost sure convergence*  
285 *to the global optimum; that is, irrespective of the variation in priors, the method converges to an*  
286 *optimal configuration  $\lambda^*$  with probability one:*

$$287 \mathbb{P}\left(\lim_{t \rightarrow \infty} f(\lambda_t) = f(\lambda^*)\right) = 1.$$

288  
289 **Robustness to Misleading Priors** Asymptotically, the algorithm does not suffer degradation in  
290 performance even when faced with misleading priors. Formally, it is guaranteed that the best ob-  
291 jective value found after  $t$  iterations is, in the limit, no worse than that of vanilla BO, providing a  
292 safeguard against user-supplied noise or misguidance. This is described in the following theorem:

293 **Theorem 2** (Robustness to Misleading Priors). *Let  $f_t^*$  be the optimal value of  $f$  found after  $t$  itera-*  
294 *tions using DynaBO, and let  $f_{t,BO}^*$  denote the corresponding value for standard BO. Then:*

$$295 \limsup_{t \rightarrow \infty} (f_{t,BO}^* - f_t^*) \leq 0.$$

296  
297 **Acceleration of Convergence with Informative Priors** While being robust to misleading priors,  
298 DynaBO is also designed to benefit from informative ones. When user-provided priors place high  
299 probability mass in regions near or containing the global optimum, the optimization process can  
300 concentrate evaluations more effectively within these promising areas. This focused search results  
301 in an improved convergence rate, captured by a regret bound dependent on the prior’s concentration.  
302 The following result formalizes this intuition by showing that the regret bound depends only on the  
303 information gain within the neighborhood of the optimum:

304 **Theorem 3** (Acceleration of Convergence with Informative Priors). *Suppose there exists a prior*  
305  *$\pi^{(m)}$  such that the prior mass in a neighborhood  $U_\epsilon(\lambda^*)$  of the global optimum  $\lambda^*$  satisfies*

$$306 \int_{U_\epsilon(\lambda^*)} \pi^{(m)}(\lambda) d\lambda \geq 1 - \delta,$$

307 *for some small  $\delta > 0$ . Then, the expected cumulative regret  $\mathbb{E}[R_T]$  of our method after  $T$  iterations*  
308 *satisfies*

$$309 \mathbb{E}[R_T] = \mathcal{O}\left(\sqrt{T\beta_{UCB_T}\gamma_T(U_\epsilon)}\right),$$

310 *with high probability, where  $\beta_{UCB_T}$  is a confidence parameter of the UCB algorithm.  $\gamma_T(U_\epsilon)$  is the*  
311 *maximum information gain restricted to the neighborhood  $U_\epsilon(\lambda^*)$ , and  $\gamma_T(U_\epsilon) < \gamma_T(\Lambda)$ .*

## 312 6 EMPIRICAL EVALUATION

313  
314 In this section, we present the empirical evaluation of DynaBO, analyzing its anytime performance  
315 across various **black-box** benchmark scenarios and for different qualities of priors. Section 6.1  
316 describes how priors are constructed. The experiment setup is detailed in Section 6.2 and the results  
317 are presented in Section 6.3. To investigate DynaBO’s sensitivity with respect to the prior rejection  
318 threshold  $\tau$ , we conduct a sensitivity analysis in Section 6.4.

## 6.1 PRIOR CONSTRUCTION: EXPERT, INFORMATIVE, LOCAL, ADVERSARIAL

Inspired by the evaluation protocols of Souza et al. (2021a); Hvarfner et al. (2022); Mallik et al. (2023), and Seng et al. (2025), we construct artificial, data-driven priors. To ground our investigation in an analysis of well and poorly performing areas of the configuration space, we conduct an extensive search for every benchmark scenario and cluster the found configuration-loss pairs  $(\lambda, \ell_\lambda) := (\lambda, f(\lambda))$  hierarchically via Gower’s distance (Gower, 1971) into  $n$  clusters. For each of the  $n$  clusters,  $c_1, \dots, c_n$ , we compute its centroid  $\bar{c}_i$  and the median loss  $\bar{\ell}_{c_i}$  of configurations contained. In the following, we assume the clusters to be ordered according to their median loss, that is,  $\bar{\ell}_{c_i} \leq \bar{\ell}_{c_j}$  for  $i < j$ . For more information on the construction of clusters, we refer to Appendix B.

To obtain dynamic priors, considering the current incumbent  $\hat{\lambda}$  and its loss  $\ell_{\hat{\lambda}}$ , we select the cluster  $c^+$  and configuration  $\lambda^+$  according to the following four prior policies simulating different aspects and levels of informativeness. The chosen configuration  $\lambda^+$  is then used as the center of a normal distribution over the configuration space to guide the optimization process toward its broader region.

**Expert Priors** bias toward clusters spanning significantly-better regions of the configuration space, that is,  $\ell_{\hat{\lambda}} \geq \bar{\ell}_{c_i}$ . A cluster  $c^+$  is sampled with probability  $\mathbb{P}(c_i) \propto e^{0.1 \cdot i}$ . From this cluster, we choose the best configuration  $\lambda^+ \in \arg \min_{\lambda \in c^+} \ell_\lambda$  as the center of the normal distribution of our prior.

**Advanced Priors** bias toward clusters spanning better-performing regions of the configuration space, that is,  $\ell_{\hat{\lambda}} \geq \bar{\ell}_{c_i}$ . A cluster  $c^+$  is sampled with probability  $\mathbb{P}(c_i) \propto e^{0.15 \cdot i}$ . From this cluster, we sample  $\lambda^+ \in c^+$  randomly as the center of the normal distribution of our prior.

**Local Priors** bias toward well-performing clusters close to the current incumbent. To this end, the incumbents’ Gower’s distance (Gower, 1971) to each cluster  $D^{gower}(\hat{\lambda}, \bar{c}_i)$  is utilized to select the 10 closest, later considered clusters  $C^+$ . The cluster with the lowest median loss  $c^+ \in \arg \min_{c \in C^+} \bar{\ell}_c$  is selected, and the prior center  $\lambda^+ \in c^+$  is sampled randomly. In contrast to previous evaluations of prior-guided BO (Souza et al., 2021a; Hvarfner et al., 2022; Mallik et al., 2023), we hypothesize that such local priors, exploiting knowledge about the location of current incumbents, are more similar to human behavior. Simultaneously, local priors are closer to already observed data used for fitting the surrogate model, facilitating informed acquisition.

**Adversarial Priors** bias toward sampling configurations in poorly performing regions of the configuration space. For that,  $c^+$  is randomly sampled from the five clusters with the worst median loss. The center of the prior is set to  $\lambda^+ \in \arg \max_{\lambda \in c^+} \ell_\lambda$ . This is meant to simulate the worst case in which a human user provides priors based on wrong assumptions.

## 6.2 EXPERIMENT SETUP

Although our theoretical analysis focuses on LCB, we conduct our experimental evaluation using the common EI acquisition function. As baselines, we consider the state-of-the-art approach  $\pi$ BO and vanilla Bayesian optimization (vanilla BO), as implemented in the SMAC3 library (Lindauer et al., 2022).  $\pi$ BO allows the user to provide prior information before the optimization process, whereas vanilla BO is without any user guidance. We chose to disregard BoPro (Souza et al., 2021b) as a baseline, since its substantially dominated by  $\pi$ BO. We evaluated the approaches both with Gaussian processes (GPs) and random forests (RFs). In accordance with the results of Eggenesperger et al. (2021), SMAC with RFs showed stronger results overall and can handle mixed spaces natively; the results with GPs are provided in Appendix D.1. We utilize LCB as an acquisition function for prior rejection. Details on the competitors are provided in Appendix C.2.

In our experiments, we evaluate various models, ranging from traditional to complex deep learning models, discussed in Appendix C.1. For each model and dataset, we repeat every HPO run with 30 different seeds. Given the same seed, DYNABO and  $\pi$ BO will sample identical (first) priors. For simple models, we allow for 200 trials and four priors after 50, 90, 130, and 170 trials. Considering the comparatively large run-time of complex models, we reduce their trial budget to 50 with priors after iterations 13, 23, 33, and 43. Experiments for random prior times are presented in Appendix D.3. Following the recommendations by Hvarfner et al. (2022),  $\beta$  is initialized as  $N/10$  with  $N$  denoting the number of overall trials. All experiments discussed in this paper were executed on HPC nodes

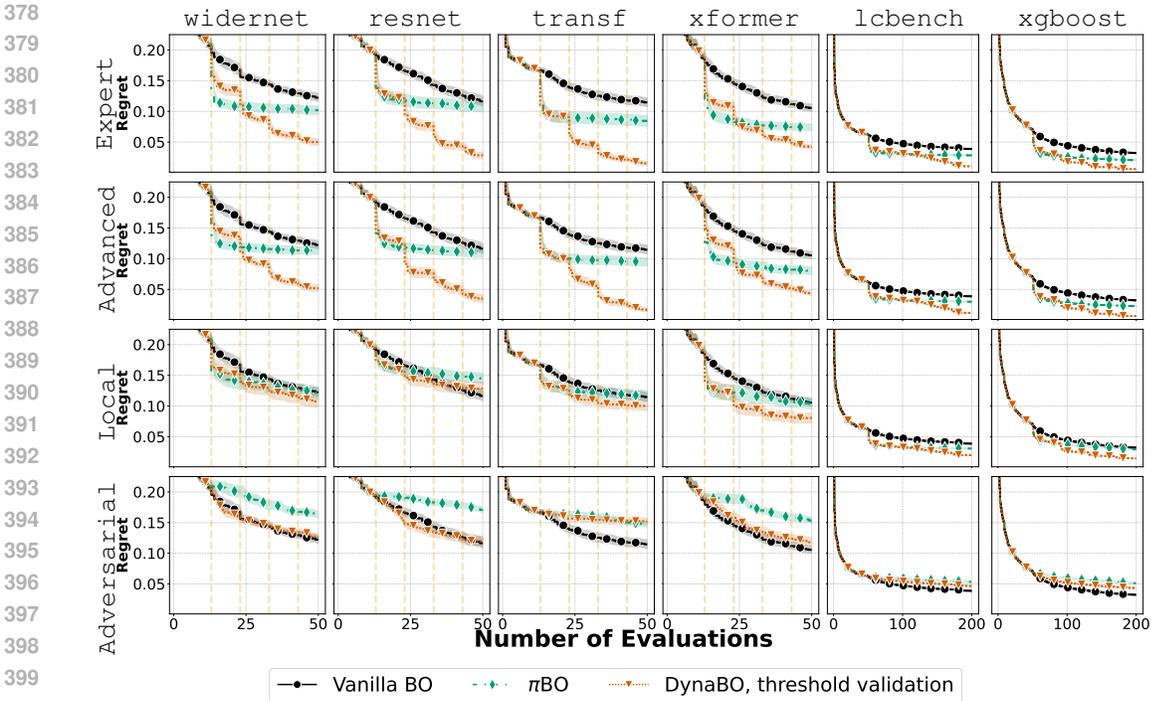


Figure 4: Mean regret for lcbench, xgboost, and PD1 using Expert, Advanced, Local, and Adversarial priors. Priors are provided at vertical lines. The shaded areas visualize the standard error. For lcbench and xgboost, the plots average all datasets. The results indicate DynaBO outperforming piBO and remaining competitive to vanilla BO for adversarial priors.

equipped with 2 Intel(R) Xeon(R) Platinum 8470 @2.0GHz processors and 488GiB RAM, of which 2 CPU cores and 6GB RAM were allocated per run.

### 6.3 RESULTS

In Figure 4, we present anytime performance plots where the mean regret of the best found incumbent, wrt. preceding exploration detailed in Appendix B, is plotted over the number of evaluations of the black box function  $f$ .

**Expert and Advanced Priors** Generally speaking, DynaBO outperforms vanilla BO, and piBO significantly in both anytime and final performance. On widernet and xformer, DynaBO is predominated by piBO until the second prior is provided, due to overly-cautious rejection of priors. Generally, expert priors provide a larger performance boost than advanced priors.

**Local Priors** For local priors, we see similar results, but with a reduced gain through the provided priors. In the case of resnet, local priors reduce rather than improve performance for both piBO and DynaBO. However, this impacts piBO more than DynaBO.

**Adversarial Priors** Generally speaking, adversarial priors result in performance degradation for piBO, as well as DynaBO. However, our rejection mechanism results in a significant performance boost except for transformer\_lm1b. The recovery from adversarial priors is analyzed in Figure 5 using an increased budget. Here DynaBO outperforms piBO if equipped with Gaussian processes, and DynaBO’s prior rejection scheme is necessary if random forests are used as surrogate model. Additional plots, considering each scenario individually, are provided in Appendix D.4.

As a general trend, DynaBO performs equal or superior to piBO. Additionally, DynaBO’s rejection mechanism boosts performance for both Local and Adversarial priors, while only marginally reducing solution quality for Expert and Advanced priors. Other results using GPs for the PD1 benchmarks (Wang et al., 2024), presented in Appendix D.1, further support the analysis above. In

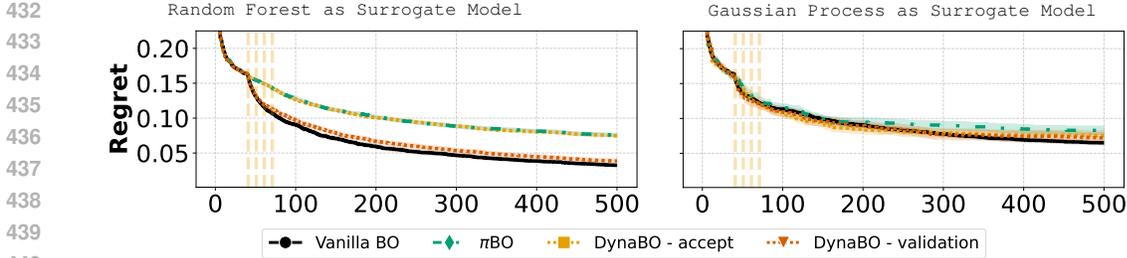


Figure 5: Anytime regret for PD1 averaged over 30 seeds, and scenarios comparing vanilla BO,  $\pi$ BO, DynaBO-accept all priors, and DynaBO with validation (DynaBO-validation).

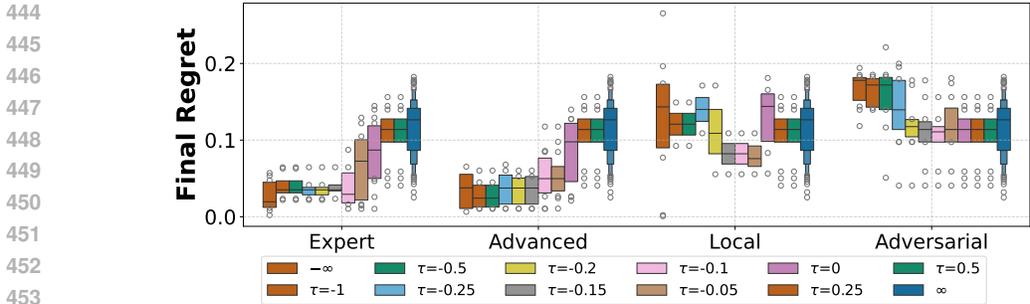


Figure 6: Sensitivity analysis of different thresholds  $\tau$ . Setting  $\tau = -\infty$  accepts all, while setting  $\tau = \infty$  rejects all priors. The boxplots contain the merged results from all scenarios.

Appendix D.7 we investigate the effect of decaying the prior faster, or slower. The results show that a linear prior decay is an adequate choice.

#### 6.4 SENSITIVITY ANALYSIS OF THE PRIOR REJECTION CRITERION

Our prior rejection scheme utilizes a threshold  $\tau$ , encoding the minimum estimated average improvement over the current incumbent needed to accept a prior. This improvement is quantified with LCB, one option for Equation (3). When  $\tau < 0$ , priors are accepted, even if they are predicted to be misleading;  $\tau > 0$  ensures that only priors of ample potential are accepted.

To study the impact of  $\tau$  on DynaBO, we conduct a sensitivity analysis on PD1’s (Wang et al., 2024) optimization scenarios, shown in Figure 6 (further results in Appendix D). As anticipated,  $\tau$  enables a tradeoff between being permissive to potentially helpful and rejecting misleading priors. While  $\tau$  could in principle be customized to reflect the user’s confidence or expertise level, we find that setting  $-0.25 \leq \tau \leq -0.05$  strikes a good balance. The preceding experiments utilize  $\tau = -0.15$ .

### 7 CONCLUSION

In this work, we presented DynaBO, a method that seamlessly incorporates user beliefs into Bayesian optimization (BO) for hyperparameter optimization. To this end, DynaBO can adapt any acquisition function by incorporating priors at any point in the process. To ensure robustness, DynaBO includes a safeguard that rejects misleading priors. Empirical results demonstrate that DynaBO outperforms both vanilla BO and  $\pi$ BO with a static prior, while also offering theoretical guarantees of robustness and effectiveness—even without the safeguard. Nonetheless, the safeguard strengthens empirical resilience against misleading user guidance.

However, while our empirical evaluation considers different kinds of user priors, they are generated synthetically. Future work should focus on expanding the types of priors supported, carrying out user studies, and developing improved rejection mechanisms, for example, by evaluating prior-based configurations. Additionally, analyzing the effect of different surrogate models on prior handling in detail and extending DynaBO to multi-fidelity optimization similar to Mallik et al. (2023) are

worth exploring. Lastly, we aim to enhance DynaBO with explainable AI (XAI) and large language models (LLMs) to support more transparent, user-friendly, low/no-code interfaces, fostering more intuitive collaboration between users and HPO approaches. Lastly, as in related papers (Hvarfner et al., 2022; Seng et al., 2025), DynaBO focuses on black-box optimization and should be extended to multi-fidelity optimization.

**Reproducibility Statement** We have taken several steps to ensure the reproducibility of our results. We follow standard methodological practices, including the use of best practices for random seeding. An additional, in-depth description of the experimental setup is provided in Appendix C. A detailed description of the conducted data generation runs is provided in Appendix B. Furthermore, we provide a reference implementation via GitHub: <https://anonymous.4open.science/r/DynaBO-EBB3/>.

**Ethical Considerations.** We believe that our work does not raise any specific ethical concerns beyond the standard considerations associated with research in machine learning.

## REFERENCES

- M. Adachi, B. Planden, D. A. Howey, M. A. Osborne, S. Orbell, N. Ares, K. Muandet, and S. L. Chau. Looping in the human: Collaborative and explainable bayesian optimization. In S. Dasgupta, S. Mandt, and Y. Li (eds.), *Proceedings of the 27th International Conference on Artificial Intelligence and Statistics (AISTATS'24)*. Proceedings of Machine Learning Research, 2024.
- R. Agrawal. Sample mean based index policies by  $o(\log n)$  regret for the multi-armed bandit problem. *Advances in Applied Probability*, 1995.
- A. Bansal, D. Stoll, M. Janowski, A. Zela, and F. Hutter. JAHS-bench-201: A foundation for research on joint architecture and hyperparameter search. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Proceedings of the 36th International Conference on Advances in Neural Information Processing Systems (NeurIPS'22)*. Curran Associates, 2022.
- J. Bergstra and Y. Bengio. Random search for hyper-parameter optimization. *Journal of Machine Learning Research*, 13:281–305, 2012.
- J. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl. Algorithms for hyper-parameter optimization. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K. Weinberger (eds.), *Proceedings of the 25th International Conference on Advances in Neural Information Processing Systems (NeurIPS'11)*, pp. 2546–2554. Curran Associates, 2011.
- B. Bischl, M. Binder, M. Lang, T. Pielok, J. Richter, S. Coors, J. Thomas, T. Ullmann, M. Becker, A.-L. Boulesteix, D. Deng, and M. Lindauer. Hyperparameter optimization: Foundations, algorithms, best practices, and open challenges. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, pp. e1484, 2023.
- O. Bojar, R. Chatterjee, C. Federmann, B. Haddow, M. Huck, C. Hokamp, P. Koehn, V. Logacheva, C. Monz, M. Negri, M. Post, C. Scarton, L. Specia, and M. Turchi. Findings of the 2015 workshop on statistical machine translation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, 2015.
- X. Bouthillier and G. Varoquaux. Survey of machine-learning experimental methods at NeurIPS2019 and ICLR2020. Research report [hal-02447823], Inria Saclay Ile de France, 2020.
- C. Chang, M. Azvar, C. Okwudire, and R. A. Kontar. Llinbo: Trustworthy llm-in-the-loop bayesian optimization. *arXiv:2505.14756 [cs.LG]*, 2025.
- T. Chen and C. Guestrin. XGBoost: A scalable tree boosting system. In B. Krishnapuram, M. Shah, A. Smola, C. Aggarwal, D. Shen, and R. Rastogi (eds.), *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'16)*, pp. 785–794. ACM Press, 2016.

- 540 A. Cowen-Rivers, W. Lyu, R. Tutunov, Z. Wang, A. Grosnit, R. Griffiths, A. Maraval, H. Jianye,  
541 J. Wang, J. Peters, and H. Ammar. HEBO: Pushing the limits of sample-efficient hyper-parameter  
542 optimisation. *Journal of Artificial Intelligence Research*, 74:1269–1349, 2022.
- 543  
544 J. Deng, W. Dong, R. Socher, L. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical  
545 Image Database. In *Proceedings of the International Conference on Computer Vision and Pattern  
546 Recognition (CVPR’09)*, pp. 248–255. ieeecs, IEEE, 2009.
- 547 K. Eggenberger, P. Müller, N. Mallik, M. Feurer, R. Sass, A. Klein, N. Awad, M. Lindauer, and  
548 F. Hutter. HPOBench: A collection of reproducible multi-fidelity benchmark problems for HPO.  
549 In J. Vanschoren and S. Yeung (eds.), *Proceedings of the Neural Information Processing Systems  
550 Track on Datasets and Benchmarks*. Curran Associates, 2021.
- 551 S. Falkner, A. Klein, and F. Hutter. BOHB: Robust and efficient Hyperparameter Optimization  
552 at scale. In J. Dy and A. Krause (eds.), *Proceedings of the 35th International Conference on  
553 Machine Learning (ICML’18)*, volume 80, pp. 1437–1446. Proceedings of Machine Learning  
554 Research, 2018.
- 555 M. Feurer, J. Springenberg, and F. Hutter. Initializing Bayesian Hyperparameter Optimization via  
556 meta-learning. In B. Bonet and S. Koenig (eds.), *Proceedings of the Twenty-ninth AAAI Confer-  
557 ence on Artificial Intelligence (AAAI’15)*, pp. 1128–1135. AAAI Press, 2015.
- 558 M. Feurer, B. Letham, and E. Bakshy. Scalable meta-learning for bayesian optimization using  
559 ranking-weighted gaussian process ensembles. In R. Garnett, F. Hutter, J. Vanschoren, P. Brazdil,  
560 R. Caruana, C. Giraud-Carrier, I. Guyon, and B. Kégl (eds.), *ICML workshop on Automated  
561 Machine Learning (AutoML workshop 2018)*, 2018.
- 562 M. Feurer, B. Letham, F. Hutter, and E. Bakshy. Practical transfer learning for bayesian optimization.  
563 *arXiv:1802.02219v4 [stat.ML]*, 2022.
- 564 J. Gower. A general coefficient of similarity and some of its properties. *Biometrics*, pp. 857–871,  
565 1971.
- 566 I. Guyon, M. Lindauer, M. van der Schaar, F. Hutter, and R. Garnett (eds.). *Proceedings of the  
567 First International Conference on Automated Machine Learning*, 2022. Proceedings of Machine  
568 Learning Research.
- 569 K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings  
570 of the International Conference on Computer Vision and Pattern Recognition (CVPR’16)*, pp.  
571 770–778. Computer Vision Foundation and IEEE Computer Society, IEEE, 2016.
- 572 F. Hutter, H. Hoos, and K. Leyton-Brown. Sequential model-based optimization for general algo-  
573 rithm configuration. In C. Coello (ed.), *Proceedings of the Fifth International Conference on  
574 Learning and Intelligent Optimization (LION’11)*, volume 6683 of *Lecture Notes in Computer  
575 Science*, pp. 507–523. Springer, 2011.
- 576 F. Hutter, H. Hoos, and K. Leyton-Brown. An efficient approach for assessing hyperparameter  
577 importance. In E. Xing and T. Jebara (eds.), *Proceedings of the 31th International Conference on  
578 Machine Learning (ICML’14)*, pp. 754–762. Omnipress, 2014.
- 579 F. Hutter, L. Kotthoff, and J. Vanschoren (eds.). *Automated Machine Learning: Methods, Systems,  
580 Challenges*. Springer, 2019. Available for free at <http://automl.org/book>.
- 581 C. Hvarfner, D. Stoll, A. Souza, L. Nardi, M. Lindauer, and F. Hutter.  $\pi$ BO: Augmenting Acquisition  
582 Functions with User Beliefs for Bayesian Optimization. In *The Tenth International Conference  
583 on Learning Representations (ICLR’22)*. ICLR, 2022. Published online: [iclr.cc](https://iclr.cc).
- 584 C. Hvarfner, F. Hutter, and L. Nardi. A general framework for user-guided bayesian optimization.  
585 In *The Twelfth International Conference on Learning Representations (ICLR’24)* ICL (2024).  
586 Published online: [iclr.cc](https://iclr.cc).
- 587 *Proceedings of the International Conference on Learning Representations (ICLR’24)*, 2024. ICLR.  
588 Published online: [iclr.cc](https://iclr.cc).
- 589  
590  
591  
592  
593

- 594 D. Jones, M. Schonlau, and W. Welch. Efficient global optimization of expensive black box func-  
595 tions. *Journal of Global Optimization*, 13:455–492, 1998.
- 596
- 597 J. Jost. *Postmodern analysis*. Springer, 2005.
- 598 N Kannengießer, N. Hasebrook, F. Morsbach, M. Zöller, J. K. H. Franke, and M. Lindauer F. Hut-  
599 ter A. Sunyaev. Practitioner motives to use different hyperparameter optimization methods. *ACM*  
600 *Transactions on Computer-Human Interaction*, 37(4), 2025. doi: 10.1145/3745771.
- 601
- 602 A. Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University  
603 of Toronto, 2009.
- 604 L. Li, K. Jamieson, G. DeSalvo, A. Rostamizadeh, and A. Talwalkar. Hyperband: Bandit-based  
605 configuration evaluation for Hyperparameter Optimization. In *The Fifth International Conference*  
606 *on Learning Representations (ICLR’17)*. ICLR, 2017. Published online: [iclr.cc](https://iclr.cc).
- 607
- 608 M. Lindauer, K. Eggensperger, M. Feurer, A. Biedenkapp, D. Deng, C. Benjamins, T. Ruhkopf,  
609 R. Sass, and F. Hutter. SMAC3: A versatile bayesian optimization package for Hyperparameter  
610 Optimization. *Journal of Machine Learning Research*, 23(54):1–9, 2022.
- 611 M. Lindauer, F. Karl, A. Klier, J. Moosbauer, A. Tornede, A. Müller, F. Hutter, M. Feurer, and  
612 B. Bischl. Position: A call to action for a human-centered automl paradigm. In Salakhutdinov  
613 et al. (2024).
- 614 A. Makarova, H. Shen, V. Perrone, A. Klein, J.B. Faddoul, A. Krause, M. Seeger, and C. Archam-  
615 beau. Automatic termination for hyperparameter optimization. In Guyon et al. (2022).
- 616
- 617 N. Mallik, C. Hvarfner, E. Bergman, D. Stoll, M. Janowski, M. Lindauer, L. Nardi, and F. Hut-  
618 ter. PriorBand: Practical hyperparameter optimization in the age of deep learning. In A. Oh,  
619 T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Proceedings of the 37th*  
620 *International Conference on Advances in Neural Information Processing Systems (NeurIPS’23)*.  
621 Curran Associates, 2023.
- 622 J. Močkus. On bayesian methods for seeking the extremum. In G. Marchuk (ed.), *Optimization*  
623 *Techniques IFIP Technical Conference 1974*, pp. 400–404. Springer, Springer, 1975.
- 624
- 625 J. Moosbauer, J. Herbinger, G. Casalicchio, M. Lindauer, and B. Bischl. Explaining hyperparameter  
626 optimization via partial dependence plots. In M. Ranzato, A. Beygelzimer, K. Nguyen, P. Liang,  
627 J. Vaughan, and Y. Dauphin (eds.), *Proceedings of the 35th International Conference on Ad-*  
628 *vances in Neural Information Processing Systems (NeurIPS’21)*, pp. 2280–2291. Curran Asso-  
629 ciates, 2021.
- 630 S. Müller, M. Feurer, N. Hollmann, and F. Hutter. PFNs4BO: In-Context Learning for Bayesian  
631 Optimization. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett (eds.),  
632 *Proceedings of the 40th International Conference on Machine Learning (ICML’23)*, volume 202  
633 of *Proceedings of Machine Learning Research*. PMLR, 2023.
- 634 R. Niwa, S. Yoshida, Y. Koyama, and Y. Ushiku. Cooperative design optimization through natural  
635 language interaction. *arXiv:2508.16077 [cs.HC]*, 2025.
- 636
- 637 N. Cesa-Bianchi P. Auer and P. Fischer. Finite-time analysis of the multiarmed bandit problem.  
638 *Mach. Learn.*, 47(2-3):235–256, 2002.
- 639 L. Papenmeier, N. Cheng, S. Becker, and L. Nardi. Exploring exploration in bayesian optimization.  
640 *arXiv:2502.08208 [cs.LG]*, 2025.
- 641
- 642 V. Perrone, H. Shen, M. Seeger, C. Archambeau, and R. Jenatton. Learning search spaces for  
643 bayesian optimization: Another view of hyperparameter transfer learning. In H. Wallach,  
644 H. Larochelle, A. Beygelzimer, F. d’Alche Buc, E. Fox, and R. Garnett (eds.), *Proceedings*  
645 *of the 33rd International Conference on Advances in Neural Information Processing Systems*  
646 *(NeurIPS’19)*, pp. 12751–12761. Curran Associates, 2019.
- 647 F. Pfisterer, L. Schneider, J. Moosbauer, M. Binder, and B. Bischl. YAHPO Gym – an efficient multi-  
objective multi-fidelity benchmark for hyperparameter optimization. In Guyon et al. (2022).

- 648 S. Pineda Arango, F. Ferreira, Kadra A., Hutter F., and Grabocka J. Quick-tune: Quickly learning  
649 which pretrained model to finetune and how. In *The Twelfth International Conference on Learning*  
650 *Representations (ICLR'24)* ICL (2024). Published online: `iclr.cc`.
- 651 H. Rakotoarison, S. Adriaensen, N. Mallik, S. Garibov, E. Bergman, and F. Hutter. In-context freeze-  
652 thaw bayesian optimization for hyperparameter optimization. In Salakhutdinov et al. (2024).
- 653 A. Ramachandran, S. Gupta, S. Rana, C. Li, and S. Venkatesh. Incorporating expert prior in  
654 Bayesian optimisation via space warping. *Knowledge-Based Systems*, 195, 2020.
- 655 C. Rasmussen and C. Williams. *Gaussian Processes for Machine Learning*. The MIT Press, 2006.
- 656 R. Salakhutdinov, Z. Kolter, K. Heller, A. Weller, N. Oliver, J. Scarlett, and F. Berkenkamp (eds.).  
657 *Proceedings of the 41st International Conference on Machine Learning (ICML'24)*, volume 251  
658 of *Proceedings of Machine Learning Research*, 2024. PMLR.
- 659 R. Sass, E. Bergman, A. Biedenkapp, F. Hutter, and M. Lindauer. Deepcave: An interactive analysis  
660 tool for automated machine learning. In M. Mutny, I. Bogunovic, W. Neiswanger, S. Ermon,  
661 Y. Yue, and A. Krause (eds.), *ICML Adaptive Experimental Design and Active Learning in the*  
662 *Real World (ReALML Workshop 2022)*, 2022.
- 663 S. Segel, H. Graf, A. Tornede, B. Bischl, and M. Lindauer. Symbolic explanations for hyperparame-  
664 ter optimization. In A. Faust, C. White, F. Hutter, R. Garnett, and J. Gardner (eds.), *Proceedings of*  
665 *the Second International Conference on Automated Machine Learning*. Proceedings of Machine  
666 Learning Research, 2023.
- 667 J. Seng, F. Ventola, Z. Yu, and K. Kersting. Hyperparameter optimization via interacting with  
668 probabilistic circuits. In R. Garnett, C. Doerr, J. van Rijn, and L. Akoglu (eds.), *Proceedings*  
669 *of the Third International Conference on Automated Machine Learning*. Proceedings of Machine  
670 Learning Research, 2025.
- 671 B. Shahriari, K. Swersky, Z. Wang, R. Adams, and N. de Freitas. Taking the human out of the loop:  
672 A review of Bayesian optimization. *Proceedings of the IEEE*, 104(1):148–175, 2016.
- 673 J. Snoek, H. Larochelle, and R. Adams. Practical Bayesian optimization of machine learning al-  
674 gorithms. In P. Bartlett, F. Pereira, C. Burges, L. Bottou, and K. Weinberger (eds.), *Proceed-*  
675 *ings of the 26th International Conference on Advances in Neural Information Processing Systems*  
676 *(NeurIPS'12)*, pp. 2960–2968. Curran Associates, 2012.
- 677 A. Souza, L. Nardi, L. Oliveira, K. Olukotun, M. Lindauer, and F. Hutter. Bayesian optimization  
678 with a prior for the optimum. In N. Oliver, F. Pérez-Cruz, S. Kramer, J. Read, and J. A. Lozano  
679 (eds.), *Machine Learning and Knowledge Discovery in Databases. Research Track*, volume 12975  
680 of *Lecture Notes in Artificial Intelligence*, pp. 265–296. Springer-Verlag, 2021a.
- 681 A. Souza, L. Nardi, L. Oliveira, K. Olukotun, M. Lindauer, and F. Hutter. Bayesian optimization  
682 with a prior for the optimum. In N. Oliver, F. Pérez-Cruz, S. Kramer, J. Read, and J. Lozano  
683 (eds.), *Machine Learning and Knowledge Discovery in Databases (ECML/PKDD'21)*, volume  
684 12975 of *Lecture Notes in Computer Science*, pp. 265–296. Springer, 2021b.
- 685 N. Srinivas, A. Krause, S. Kakade, and M. Seeger. Information-theoretic regret bounds for gaussian  
686 process optimization in the bandit setting. *IEEE Trans. Inf. Theory*, 58(5):3250–3265, 2012.
- 687 S. Studer, T. Binh Bui, C. Drescher, A. Hanuschkin, L. Winkler, S. Peters, and K. Müller. Towards  
688 CRISP-ML(Q): A machine learning process model with quality assurance methodology. *Mach.*  
689 *Learn. Knowl. Extr.*, 2021.
- 690 K. Swersky, J. Snoek, and R. Adams. Multi-task Bayesian optimization. In C. Burges, L. Bottou,  
691 M. Welling, Z. Ghahramani, and K. Weinberger (eds.), *Proceedings of the 27th International*  
692 *Conference on Advances in Neural Information Processing Systems (NeurIPS'13)*, pp. 2004–  
693 2012. Curran Associates, 2013.
- 694 T. Tornede, A. Tornede, J. Hanselle, F. Mohr, M. Wever, and E. Hüllermeier. Towards green au-  
695 tomated machine learning: Status quo and future directions. *Journal of Artificial Intelligence*  
696 *Research*, 77:427–457, 2023.

- 702 R. Turner, D. Eriksson, M. McCourt, J. Kiili, E. Laaksonen, Z. Xu, and I. Guyon. Bayesian opti-  
703 mization is superior to random search for machine learning hyperparameter tuning: Analysis of  
704 the Black-Box Optimization Challenge 2020. In H. Escalante and K. Hofmann (eds.), *Proceed-*  
705 *ings of the Neural Information Processing Systems Track Competition and Demonstration*, pp.  
706 3–26. Curran Associates, 2021.
- 707 K. Van der Blom, A. Serban, H. Hoos, and J. Visser. Automl adoption in ml software. In *8th ICML*  
708 *Workshop on Automated Machine Learning (AutoML)*, 2021.
- 709 J. van Rijn and F. Hutter. Hyperparameter importance across datasets. In Y. Guo and F. Farooq (eds.),  
710 *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and*  
711 *Data Mining (KDD’18)*, pp. 2367–2376. ACM Press, 2018.
- 712 A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. Gomez, L. Kaiser, and I. Polosukhin.  
713 Attention is all you need. In I. Guyon, U. von Luxburg, S. Bengio, H. Wallach, R. Fergus,  
714 S. Vishwanathan, and R. Garnett (eds.), *Proceedings of the 31st International Conference on*  
715 *Advances in Neural Information Processing Systems (NeurIPS’17)*. Curran Associates, Inc., 2017.
- 716 Q. Wang, Y. Ming, Z. Jin, Q. Shen, D. Liu, M. Smith, K. Veeramachaneni, and H. Qu. Atmseer:  
717 Increasing transparency and controllability in automated machine learning. In *Proceedings of the*  
718 *2019 CHI Conference on Human Factors in Computing Systems (CHI’19)*, pp. 1–12. ACM Press,  
719 2019.
- 720 Z. Wang, G. Dahl, K. Swersky, C. Lee, Z. Mariet, Z. Nado, J. Gilmer, J. Snoek, and Z. Ghahramani.  
721 Pre-trained Gaussian processes for Bayesian optimization. *J. Mach. Learn. Res.*, 2024.
- 722 W. Xu, M. Adachi, C. N. Jones, and M. A. Osborne. Principled bayesian optimisation in collabora-  
723 tion with human experts. *arXiv:2410.10452 [cs.LG]*, 2024.
- 724 M. Zöllner, W. Titov, T. Schlegel, and M. Huber. Xautoml: A visual analytics tool for understanding  
725 and validating automated machine learning. *ACM Trans. Interact. Intell. Syst.*, 13(4):28:1–28:39,  
726 2023.
- 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

756	ORGANIZATION OF TECHNICAL APPENDICES AND SUPPLEMENTARY	
757	MATERIAL	
758		
759	The appendix is split into four main parts. In Appendix A, we prove our theoretical guarantees.	
760	In Appendix B, we provide a detailed discussion of our prior construction and selection scheme.	
761	In Appendix C, we discuss the experimental setup. Lastly, in Appendix D, we provide additional	
762	evaluation results.	
763		
764	<b>A Theoretical Guarantees</b>	<b>16</b>
765		
766	A.1 Theorem - Almost Sure Convergence of DynaBO . . . . .	16
767	A.2 Corollary - Robustness to Misleading Priors . . . . .	17
768	A.3 Theorem - Acceleration of Convergence with Informative Priors . . . . .	18
769		
770		
771	<b>B Prior Handling</b>	<b>19</b>
772		
773	B.1 Prior Construction and Selection . . . . .	19
774	B.2 Facilitating Prior Behavior . . . . .	19
775	B.3 Further Details on the Rejection Criterion . . . . .	21
776		
777		
778	<b>C Detailed Experimental Setup</b>	<b>21</b>
779		
780	C.1 Summary of the Benchmark . . . . .	21
781	C.2 Competitor Setup . . . . .	21
782	C.3 <a href="#">Further Details of Prior Rejection</a> . . . . .	21
783		
784	<b>D Additional Empirical Results</b>	<b>22</b>
785		
786	D.1 Gaussian Process Main Results . . . . .	22
787	D.2 Random Forest Main Results . . . . .	23
788	D.3 Random Prior Location . . . . .	24
789	D.4 Increased Budget Results . . . . .	25
790	D.5 <a href="#">Naïve Dynamic Extension for <math>\pi</math>BO</a> . . . . .	26
791	D.6 Sensitivity Analysis of the Prior Rejection Criterion . . . . .	27
792	D.7 <a href="#">Prior Decay Ablation</a> . . . . .	29
793		
794		
795		
796	<b>E Declaration of LLM Usage</b>	<b>30</b>
797		
798		
799		
800		
801		
802		
803		
804		
805		
806		
807		
808		
809		

## 810 A THEORETICAL GUARANTEES

811 For all the subsequent proofs, we make the following assumptions:

812 **A1 - Objective Function Regularity:** The objective function  $f_{\text{obj}} : \Lambda \rightarrow \mathbb{R}$  is bounded and  
813 Lipschitz-continuous on a compact domain. Additionally,  $f_{\text{obj}}$  is assumed to be a realization from a  
814 Gaussian process prior with mean  $m : \Lambda \rightarrow \mathbb{R}$  and positive definite kernel  $k(\lambda, \lambda')$ .  
815  
816

817 W.l.o.g., we minimize the objective function. However, for the theoretical analysis, we consider the  
818 maximization of the equivalent utility function  $f(\lambda) := C - f_{\text{obj}}(\lambda)$ , where the constant  $C$  is chosen  
819 to ensure  $f(\lambda) > 0$  for all  $\lambda \in \Lambda$ . This affine transformation preserves the global optimum and  
820 guarantees the strict positivity required for multiplicative re-weighting via the sum of priors.

821 **A2 - Finiteness of User Priors:** Let  $\{(t^{(m)}, \pi^{(m)})\}_{m=1}^M$  be a finite series of user-specified prior  
822 functions  $\pi^{(m)} : \Lambda \rightarrow (0, 1]$  with iteration indices  $t^{(m)} \in \mathbb{N}$  such that  $t^{(1)} < \dots < t^{(M)}$ .  
823

824 **A3 - Vanishing Influence of Priors:** There exists  $\beta \in \mathbb{R}$  such that for all indices  $m$  and  $\lambda \in \Lambda$

$$825 \lim_{t \rightarrow \infty} \pi^{(m)}(\lambda)^{\beta/(t-t^{(m)})} = 1.$$

826 That is, the multiplicative influence of the prior function on the acquisition criterion diminishes to  
827 unity with increasing number of iterations.

828 Additionally, we assume Upper Confidence Bound (UCB) as an acquisition function for maximiza-  
829 tion of the utility function.

### 830 A.1 THEOREM - ALMOST SURE CONVERGENCE OF DYNABO

831 Given Assumptions A1, A2 and A3, the DynaBO algorithm converges almost surely to the global  
832 optimum of  $f$ ; meaning, the sequence of points  $(\lambda_t)_{t \geq 1}$  selected by DynaBO satisfies

$$833 f(\lambda_t) \xrightarrow{\text{a.s.}} f(\lambda^*).$$

834 *Proof.* The proof follows standard convergence results for Bayesian optimization (BO) (Srinivas  
835 et al., 2012) augmented by the dynamic influence of priors.

836 Let  $\lambda^* \in \arg \max_{\lambda \in \Lambda} f(\lambda)$  be a global maximizer of  $f$ . At iteration  $t > t^{(M)}$ , the DynaBO  
837 acquisition function is defined as

$$838 \alpha_{\text{dyna}}(\lambda, t) := \alpha(\lambda, t) \sum_{m=1}^M \pi^{(m)}(\lambda)^{\beta/(t-t^{(m)})},$$

839 where  $\alpha(\lambda, t)$  denotes a standard GP-UCB-acquisition function.

840 To prove convergence, we must show that the instantaneous regret  $r_t := f(\lambda^*) - f(\lambda_t)$  converges  
841 to zero almost surely. We prove this by analyzing the event  $\mathcal{E}_t = \{\lambda_t \in D_\epsilon\}$  at any time  $t$ , where  
842  $D_\epsilon = \{\lambda \in \Lambda \mid f(\lambda) \leq f(\lambda^*) - \epsilon\}$  is a sub-optimal region. First, the UCB parameter  $\beta_{\text{UCB}_t}$   
843 is chosen such that the probability of the confidence bounds failing,  $\mathbb{P}(\mathcal{F}_t)$ , is summable. By the  
844 Borel-Cantelli lemma, this occurs only finitely many times, almost surely.

845 Second, we analyze the event  $\mathcal{E}_t$  on the high-probability event  $\mathcal{F}_t^C$ , complement of  $\mathcal{F}_t$ , where the  
846 bounds hold. If the algorithm selects a suboptimal sample  $\lambda_t \in D_\epsilon$ , it must be that  $\alpha_{\text{dyna}}(\lambda_t, t-1) \geq$   
847  $\alpha_{\text{dyna}}(\lambda^*, t-1)$ . Let  $P(\lambda, t) := \sum_{m=1}^M \pi^{(m)}(\lambda)^{\beta/(t-t^{(m)})}$ . By definition:

$$848 \alpha(\lambda_t, t-1) \cdot P(\lambda_t, t-1) \geq \alpha(\lambda^*, t-1) \cdot P(\lambda^*, t-1).$$

849 Since the bounds hold,  $\alpha(\lambda, t-1)$  satisfies  $\alpha(\lambda_t, t-1) \leq f(\lambda_t) + 2\beta_t^{1/2} \sigma_{t-1}(\lambda_t)$  and  $\alpha(\lambda^*, t-1) \geq$   
850  $f(\lambda^*)$ . Substituting these gives:

$$851 (f(\lambda_t) + 2\beta_t^{1/2} \sigma_{t-1}(\lambda_t)) \cdot P(\lambda_t, t-1) \geq f(\lambda^*) \cdot P(\lambda^*, t-1).$$

By Assumption A3 and the finiteness of  $M$ ,  $\lim_{t \rightarrow \infty} P(\lambda, t) = 1$ . This convergence is uniform on the compact set  $\Lambda$ , by Dini's Theorem (Jost, 2005), given continuity of  $\pi^{(m)}$  per the prior construction as a normal distribution. Therefore, for any  $\delta > 0$ , there exists a  $T$  such that for all  $t > T$ ,  $P(\lambda, t) \in (1 - \delta, 1 + \delta)$  for all  $\lambda \in \Lambda$ . For  $t > T$ , we can bound the terms  $P(\lambda_t, t - 1) \leq (1 + \delta)$  and  $P(\lambda^*, t - 1) \geq (1 - \delta)$ :

$$(f(\lambda_t) + 2\beta_t^{1/2}\sigma_{t-1}(\lambda_t))(1 + \delta) \geq f(\lambda^*)(1 - \delta).$$

As  $f(\lambda_t) \leq f(\lambda^*) - \epsilon$  (since  $\lambda_t \in D_\epsilon$ ), we have:

$$(f(\lambda^*) - \epsilon + 2\beta_t^{1/2}\sigma_{t-1}(\lambda_t))(1 + \delta) \geq f(\lambda^*)(1 - \delta).$$

Rearranging for the variance term (and assuming  $f, \epsilon > 0$  for simplicity):

$$\begin{aligned} 2\beta_t^{1/2}\sigma_{t-1}(\lambda_t)(1 + \delta) &\geq f(\lambda^*)(1 - \delta) - (f(\lambda^*) - \epsilon)(1 + \delta) \\ \Leftrightarrow 2\beta_t^{1/2}\sigma_{t-1}(\lambda_t)(1 + \delta) &\geq \epsilon(1 + \delta) - 2\delta f(\lambda^*). \end{aligned}$$

As this must hold for any  $\delta > 0$  (by picking  $t$  large enough), we can let  $\delta \rightarrow 0$ , which implies that for  $t$  sufficiently large, the condition for sampling in  $D_\epsilon$  requires:

$$2\beta_t^{1/2}\sigma_{t-1}(\lambda_t) \geq \epsilon.$$

We now show that this condition can only be met a finite number of times. Let  $t_n$  be the time of the  $n$ -th sample drawn from the sub-optimal region  $D_\epsilon$ . The UCB parameter  $\beta_{t_n}$  grows as  $O(\log t_n)$  (Srinivas et al., 2012). Crucially, Assumption A1 guarantees that the posterior variance at that point,  $\sigma_{t_n-1}(\lambda_{t_n})$ , converges to 0 as  $n \rightarrow \infty$  (Srinivas et al., 2012). Because the variance  $\sigma_{t_n-1}$  converges to 0 faster than the  $\beta_{t_n}^{1/2}$  term grows (Srinivas et al., 2012), their product must also converge to zero:

$$\lim_{n \rightarrow \infty} 2\beta_{t_n}^{1/2}\sigma_{t_n-1}(\lambda_{t_n}) = 0.$$

This directly implies that for any  $\epsilon > 0$ , the condition  $2\beta_t^{1/2}\sigma_{t-1}(\lambda_t) \geq \epsilon$  must eventually be violated. In other words, there exists a finite sample count  $N_\epsilon$  after which the necessary condition for sampling in  $D_\epsilon$  is permanently false.

In conclusion, the bounds fail only finitely often almost surely, and any sub-optimal region  $D_\epsilon$  is sampled only a finite number of times (a.s.). Therefore, the event  $\mathcal{E}_t$  occurs only finitely many times, almost surely. Since this holds for any  $\epsilon > 0$ , it follows that  $r_t \xrightarrow{\text{a.s.}} 0$  and

$$f(\lambda_t) \xrightarrow{\text{a.s.}} f(\lambda^*).$$

□

## A.2 COROLLARY - ROBUSTNESS TO MISLEADING PRIORS

Under Assumptions A1–A3, the asymptotic performance of DynaBO is robust to any finite set of misleading or incorrect priors; that is, the presence of such priors does not degrade convergence compared to standard BO. Formally, let

$$f_t^* := \max_{1 \leq i \leq t} f(\lambda_i) \quad \text{and} \quad f_{t,BO}^* := \max_{1 \leq i \leq t} f(\lambda_i^{BO})$$

be the maximum function values found by DynaBO and standard BO, respectively, at time  $t$ . Then,

$$\limsup_{t \rightarrow \infty} (f_{t,BO}^* - f_t^*) \leq 0.$$

*Proof.* By Assumption A3, analogously to Theorem A.1, the DynaBO acquisition function

$$\alpha_{\text{dyna}}(\lambda, t) = \alpha(\lambda, t) \sum_{m=1}^M \pi^{(m)}(\lambda)^{\beta/(t-t^{(m)})}$$

converges uniformly over  $\Lambda$  to the standard BO acquisition function:

$$\sup_{\lambda \in \Lambda} |\alpha(\lambda, t) - \alpha_{\text{dyna}}(\lambda, t)| \rightarrow 0 \quad \text{as } t \rightarrow \infty.$$

Hence, the sequence of query points  $\{\lambda_t\}$  selected by DynaBO converges in distribution to those of standard BO,  $\{\lambda_t^{\text{BO}}\}$ . By Assumption 1, the Lipschitz-continuity of  $f$ , the Continuous Mapping Theorem implies

$$f(\lambda_t) \xrightarrow{d} f(\lambda_t^{\text{BO}}).$$

Since  $\{f_t^*\}_{t \geq 1}$  and  $\{f_{t, \text{BO}}^*\}_{t \geq 1}$  are bounded non-decreasing sequences, their limits exist. Therefore,

$$\limsup_{t \rightarrow \infty} (f_{t, \text{BO}}^* - f_t^*) \leq 0,$$

establishing the asymptotic robustness of DynaBO to misleading priors.  $\square$

### A.3 THEOREM - ACCELERATION OF CONVERGENCE WITH INFORMATIVE PRIORS

Under Assumptions A1-A3, suppose there exists a user prior  $\pi^{(m)}$  such that the prior density  $\pi^{(m)}(\lambda)$  places a substantial mass in a neighborhood  $U_\epsilon(\lambda^*)$  containing the global optimum  $\lambda^*$ :

$$\int_{U_\epsilon(\lambda^*)} \pi^{(m)}(\lambda) d\lambda \geq 1 - \delta$$

for small  $\delta > 0$ . Then, the DynaBO algorithm achieves an improved upper bound on cumulative regret over  $T$  iterations, specifically

$$R_T := \sum_{t=1}^T (f(\lambda^*) - f(\lambda_t)) = O(\sqrt{T\beta_{\text{UCB}_T}\gamma_T(U_\epsilon)}).$$

Here,  $\gamma_T(U_\epsilon)$  denotes the maximum information gain restricted to the neighborhood  $U_\epsilon(\lambda^*)$ , satisfying  $\gamma_T(U_\epsilon) < \gamma_T(\Lambda)$ . This implies asymptotically faster convergence rates than standard BO without informative priors.

*Proof.* Define the index sets

$$\mathcal{I}_{U_\epsilon} := \{1 \leq t \leq T \mid \lambda_t \in U_\epsilon\} \quad \text{and} \quad \mathcal{I}_{\Lambda \setminus U_\epsilon} := \{1 \leq t \leq T \mid \lambda_t \notin U_\epsilon\}.$$

Similarly to Theorem A.1, BO regret bounds using the GP kernel and mutual information gain apply to  $U_\epsilon$  (Srinivas et al., 2012):

$$\sum_{t \in \mathcal{I}_{U_\epsilon}} (f(\lambda^*) - f(\lambda_t)) \leq C\sqrt{|I_{U_\epsilon}| \beta_{\text{UCB}|I_{U_\epsilon}|} \gamma_{|I_{U_\epsilon}|}(U_\epsilon)} \leq C\sqrt{T\beta_{\text{UCB}_T}\gamma_T(U_\epsilon)},$$

where  $C$  is kernel-dependent.

Now, let  $B := \max_{\lambda \in \Lambda} (f(\lambda^*) - f(\lambda))$  be the maximum instantaneous regret. Then,

$$\sum_{t \in \mathcal{I}_{\Lambda \setminus U_\epsilon}} (f(\lambda^*) - f(\lambda_t)) \leq B|\mathcal{I}_{\Lambda \setminus U_\epsilon}|.$$

The region  $\Lambda \setminus U_\epsilon$  is sub-optimal. It is a result of UCB analysis (P. Auer & Fischer, 2002) that the expected number of times a sub-optimal arm is sampled,  $\mathbb{E}[|\mathcal{I}_{\Lambda \setminus U_\epsilon}|]$ , is bounded logarithmically in time, i.e.,  $O(\log T)$ . Since  $O(\log T)$  grows strictly slower than  $T$ , it is  $o(T)$ . Therefore,

$$\mathbb{E} \left[ \sum_{t \in \mathcal{I}_{\Lambda \setminus U_\epsilon}} (f(\lambda^*) - f(\lambda_t)) \right] \leq B\mathbb{E}[|\mathcal{I}_{\Lambda \setminus U_\epsilon}|] = o(T).$$

By linearity of expectation,

$$\mathbb{E}[R_T] = \mathbb{E} \left[ \sum_{t=1}^T (f(\lambda^*) - f(\lambda_t)) \right] = \mathbb{E} \left[ \sum_{t \in \mathcal{I}_{U_\epsilon}} (f(\lambda^*) - f(\lambda_t)) \right] + \mathbb{E} \left[ \sum_{t \in \mathcal{I}_{\Lambda \setminus U_\epsilon}} (f(\lambda^*) - f(\lambda_t)) \right].$$

Plugging in the previous results yields

$$\mathbb{E}[R_T] \leq C\sqrt{T\beta_{\text{UCB}_T}\gamma_T(U_\epsilon)} + o(T).$$

As for the asymptotic implications: Since  $\gamma_T(U_\epsilon) < \gamma_T(\Lambda)$ , the cumulative regret bound improves over standard BO without informative priors when  $\delta$  is sufficiently small.  $\square$

## 972 B PRIOR HANDLING

### 973 B.1 PRIOR CONSTRUCTION AND SELECTION

974 As mentioned in Section 2, given a learner (or learning algorithm)  $A$ , dataset  $D = (D_T, D_V)$ , and  
 975 a loss function  $\ell$ , HPO aims to find a well-performing hyperparameter configuration  $\hat{\lambda}$ . A learning  
 976 algorithm  $A$ , configured with  $\hat{\lambda}$ , and applied to  $D$ , results in a hypothesis  $h_{\hat{\lambda}, D}$  minimizing the loss  
 977 function  $\ell$ .

978 As discussed in Section 6.1, we construct priors based on data collected through Bayesian optimiza-  
 979 tion runs. These data generation runs are conducted as follows:

- 980 1. Generate prior data: For each learner  $A$ , dataset  $D$  combination, execute explorative  
 981 Bayesian optimization runs with both the more greedy Expected Improvement (EI) and  
 982 more explorative Lower Confidence Bounds (LCB) (Papenmeier et al., 2025). For each ac-  
 983 quisition function, run 10 seeds for a budget of 5,000 iterations. Then, for every algorithm,  
 984 dataset combination, concatenate the lists of preliminary incumbents and assemble a joint  
 985 list sorted by losses  $\ell_\lambda$ :

$$986 I_{A,D} = [(\lambda_1, \ell_{\lambda_1}), (\lambda_2, \ell_{\lambda_2}), \dots, (\lambda_n, \ell_{\lambda_n})].$$

987 Note the abuse of notation: We denote with  $\ell_\lambda$  the loss on the validation set as introduced  
 988 in Equation (1):

$$989 \ell_\lambda = \left[ \frac{1}{|D_V|} \sum_{(x,y) \in D_V} \ell(y, h_{\lambda, D_T}(x)) \right].$$

- 990 2. To ensure that also non-well-performing areas of the configuration space are covered,  $I_{A,D}$   
 991 is supplemented with  $n$  non-incumbent configurations and their loss.
- 992 3. Due to structured configuration spaces, some hyperparameters may not be active. In the  
 993 case of our experiments, this only occurs for numeric hyperparameters. These values are  
 994 filled with  $-1$ .
- 995 4. Create clusters of configurations of preliminary incumbents in the configuration space:  
 996 For each learner  $A$ , dataset  $D$  combination, cluster the incumbent configurations into 100  
 997 clusters

$$998 C_{A,D} = \{c_1, c_2, \dots, c_{100}\} \quad \text{with} \quad c_i = \{(\lambda_{c_i^1}, \ell_{\lambda_{c_i^1}}), (\lambda_{c_i^2}, \ell_{\lambda_{c_i^2}}), \dots\}$$

1000 using Agglomerative Clustering with Gower’s Distance (Gower, 1971) and Ward Linkage.  
 1001 For each cluster, compute a centroid  $\bar{c}_i$  and the median performance  $\bar{\ell}_{c_i}$ .

1002 During optimization of  $A$  on dataset  $D$ , priors are generated dynamically.

- 1003 1. Sample prior configuration  $\lambda^+$  as described in Section 6.1.
- 1004 2. Build prior: We hypothesize that with each prior provided to DYNABO, the confidence of a  
 1005 user would grow. In our synthetic prior generation, we therefore build the  $k$ -th prior  $\pi^k$  as  
 1006 follows: For each numerical hyperparameter  $\lambda_j$  with lower bounds  $\lambda_1^l, \lambda_2^l, \dots, \lambda_d^l$  and upper  
 1007 bounds  $\lambda_1^u, \lambda_2^u, \dots, \lambda_d^u$ , set

$$1008 \pi^m = [\mu_j, \sigma_j]_{j=1}^d = \left[ \left( \lambda_j^+, \frac{|\lambda_j^u - \lambda_j^l|}{k \cdot 5} \right) \right]_{j=1}^d. \quad (4)$$

- 1009 3. For each categorical hyperparameter  $\lambda_j = \lambda_j^+$ .

1010 As mentioned in Section 6.2, we provide four priors for evaluations on YAHPO Gym, and four  
 1011 priors  $\pi^1, \pi^2, \pi^3, \pi^4$  for evaluations on PD1, respectively.

### 1012 B.2 FACILITATING PRIOR BEHAVIOR

1013 **Numerical Stability** To ensure numerical stability and to maintain the impact of the initial ac-  
 1014 quisition function, we clip priors at  $1e-12$  and thereby ensure that priors take values in  $(0, 1]$ . Further-  
 1015 more, due to the decaying mechanism, priors converge to 1 with increasing  $t$ .

**Adapting Candidate Sampling** Furthermore, to facilitate sampling candidate configurations with a high prior impact, we adapt SMAC3’s (Lindauer et al., 2022) Local and Random Search candidate generation. If no prior is provided, the standard SMAC acquisition function samples 5000 random configurations, ranks them according to the acquisition function, and selects the 10 best-ranked configurations as starting points for local search among other candidates.

Suppose a finite sequence of user-specified priors  $\{\pi^{(m)}\}_{m=1}^M$  is provided at times  $\{t^{(m)}\}_{m=1}^M$ , with  $t^{(1)} < \dots < t^{(M)} \leq T$ , then the random configurations are replaced as follows:

- At iteration  $t$  prior  $\pi^{(m)}$  is associated with a weight  $\omega_m = e^{-0.126 \cdot (t - t^{(m)})}$ .
- If  $\sum_{m=1}^M \omega_m \leq 0.9$ , each prior is used to sample  $\lfloor \omega_m + 0.5 \rfloor$  configurations. The rest are sampled uniformly at random.
- If  $\sum_{m=1}^M \omega_m > 0.9$ , each prior is used to sample  $\left\lfloor \frac{\omega_m}{\sum_{j=1}^M \omega_j} + 0.5 \right\rfloor$  configurations. The rest are sampled uniformly at random.

### B.3 FURTHER DETAILS ON THE REJECTION CRITERION

To apply our prior rejection scheme

$$\mathbb{E}_{\lambda \sim \mathcal{N}_{\pi^{(m)}}} [\alpha_{\hat{f}}(\lambda)] - \mathbb{E}_{\lambda^+ \sim \mathcal{N}_{\hat{\lambda}^+}} [\alpha_{\hat{f}}(\lambda^+)] \geq \tau . \quad (5)$$

On categorical hyperparameters, we conduct the following adaptations: For the configurations sampled according to the prior  $\lambda \sim \mathcal{N}_{\pi^{(m)}}$ , categorical values are sampled according to the provided weights. For the configurations sampled in the area around the incumbent  $\lambda^+ \sim \mathcal{N}_{\hat{\lambda}^+}$ , the incumbents configuration is utilized.

## C DETAILED EXPERIMENTAL SETUP

### C.1 SUMMARY OF THE BENCHMARK

Rather than training many models with different hyperparameter configurations, we use surrogate models provided by (Pfisterer et al., 2022) for XGBoost (`xgboost`) (Chen & Guestrin, 2016) and multi-layer perceptrons (`lcbench`) dubbed traditional machine learning. For complex architectures, we utilize surrogates trained for Mallik et al. (2023) based on data collected by Wang et al. (2024). We consider a wide ResNet (He et al., 2016) (`widernet`) trained on CIFAR10 (Krizhevsky, 2009), a ResNet (He et al., 2016) (`resnet`) trained on ImageNet (Deng et al., 2009), a transformer (Vaswani et al., 2017) (`transf`) trained on translatewmt (Bojar et al., 2015), and a transformer trained on WMT15 (`xformer`) (Bojar et al., 2015). A learner, searchspace, and dataset overview is provided in Table 1.

Table 1: An overview of the evaluated *scenarios*, each with the considered configuration *configuration space* type and number of *datasets*, with which each scenario was evaluated.

Scenario	Configuration Space	# Datasets
rbv2_xgboost	14D: Mixed	119
lcbench	7D: Numeric	34
cifar100_wideresnet_2048	4D: Numeric	1
imagenet_resnet_512	4D: Numeric	1
lm1b_transformer_2048	4D: Numeric	1
translatewmt_xformer_64	4D: Numeric	1

Our experiments are scheduled, and the results are logged in a MySQL database using the PyExperimenter library (Tornede et al., 2023).

### C.2 COMPETITOR SETUP

Our experiments are built on top of SMAC3 (Lindauer et al., 2024), for vanilla BO,  $\pi$ BO, and DynaBO. We utilize the Hyperparameter Optimization Facade, but deactivate its standard log transformations. We also refit the surrogate after every evaluated configuration.

### C.3 IMPLEMENTATION DETAILS OF PRIOR REJECTION

Following the principle of optimism in the face of uncertainty, we assess the potential of both regions in terms of their lower confidence bounds (LCBs) (Agrawal, 1995):

$$LCB(\lambda) = (-1) \cdot (\mu(\lambda) - \kappa\sigma(\lambda)),$$

where  $\mu(\cdot)$  denotes the mean predicted with  $\hat{f}$  for  $\lambda$ , and  $\sigma(\cdot)$  the uncertainty in terms of standard deviation. This way, we interpret all uncertainty the surrogate model may have at a configuration  $\lambda$  as the potential for its performance to achieve an improvement at all. Intuitively, the LCB criterion allows us to quantify the explorative potential of the prior region as opposed to the exploitative potential close to the current incumbent, since uncertainty close to the incumbent is typically low.

Provided a prior  $\pi^{(m)}$  with mean  $\mu^{(m)}$ , and standard deviation  $\sigma^{(m)}$  as in Equation (4), our initialization of the rejection criterion detailed in Equation (3)

$$\mathbb{E}_{\lambda \in \mathcal{N}_{\pi^{(m)}}} [\alpha(\lambda)] - \mathbb{E}_{\lambda \in \mathcal{N}_{\hat{\lambda}}} [\alpha(\lambda)] \geq \tau \quad (6)$$

utilizes 500 configurations from normal distributions for both  $\mathcal{N}_{\pi^{(m)}} \sim (\mu^{(m)}, \sigma^{(m)})$ , and  $\mathcal{N}_{\hat{\lambda}} \sim (\hat{\lambda}, \mu^{(m)})$ . Our main experiments utilize  $\tau = -0.15$ .

## D ADDITIONAL EMPIRICAL RESULTS

Our additional experimental results focus on further validating our approach. Appendix D.1 contains the results on the PD1 benchmark (Wang et al., 2024) with gaussian processes as a surrogate. For an easier comparison, the corresponding random forest results are provided in Appendix D.2. Appendix D.3 contains the results of experiments conducted with randomly sampled prior locations. All experiments shown here are conducted on the PD1 benchmark.

### D.1 GAUSSIAN PROCESS MAIN RESULTS

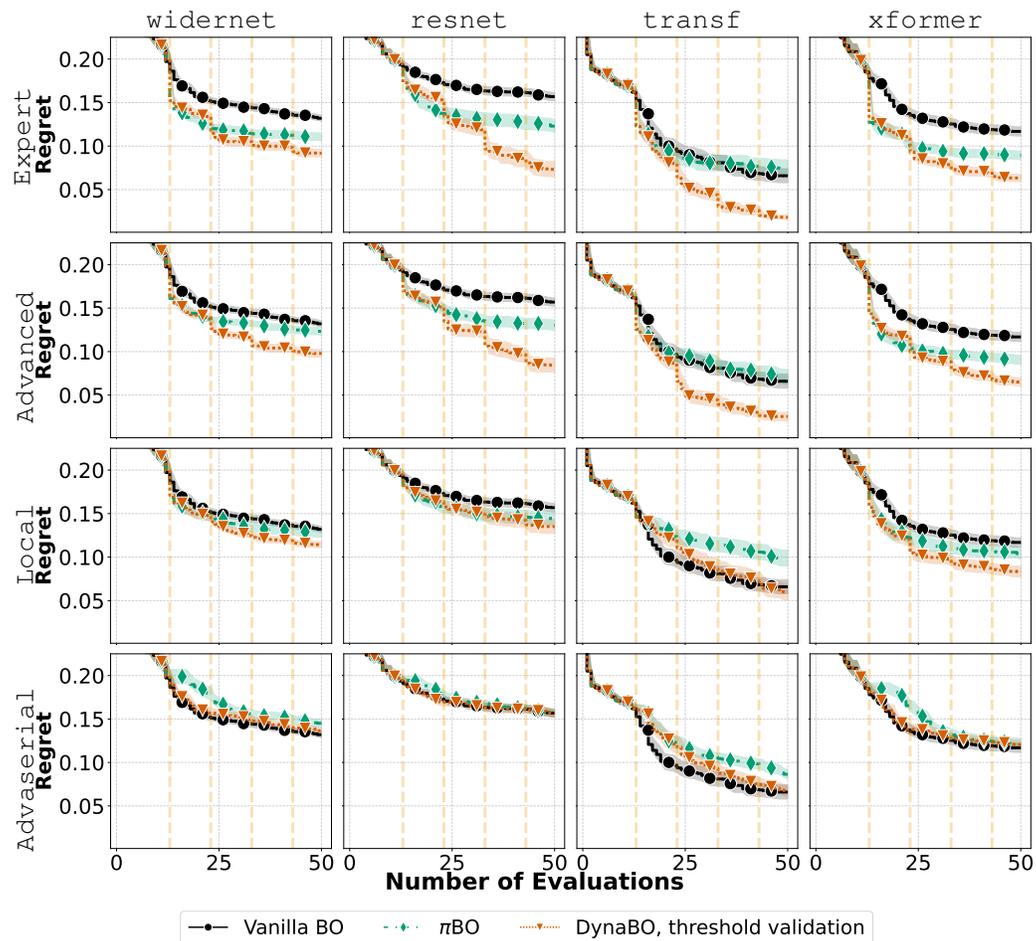


Figure 7: Mean regret for and PD1 using Expert, Advanced, Local, and Adversarial priors, with Gaussian processes as surrogate models. Priors are provided at vertical lines. The shaded areas visualize the standard error. The results indicate DynaBO outperforming  $\pi$ BO and remaining competitive with vanilla BO for adversarial priors.

## D.2 RANDOM FOREST MAIN RESULTS

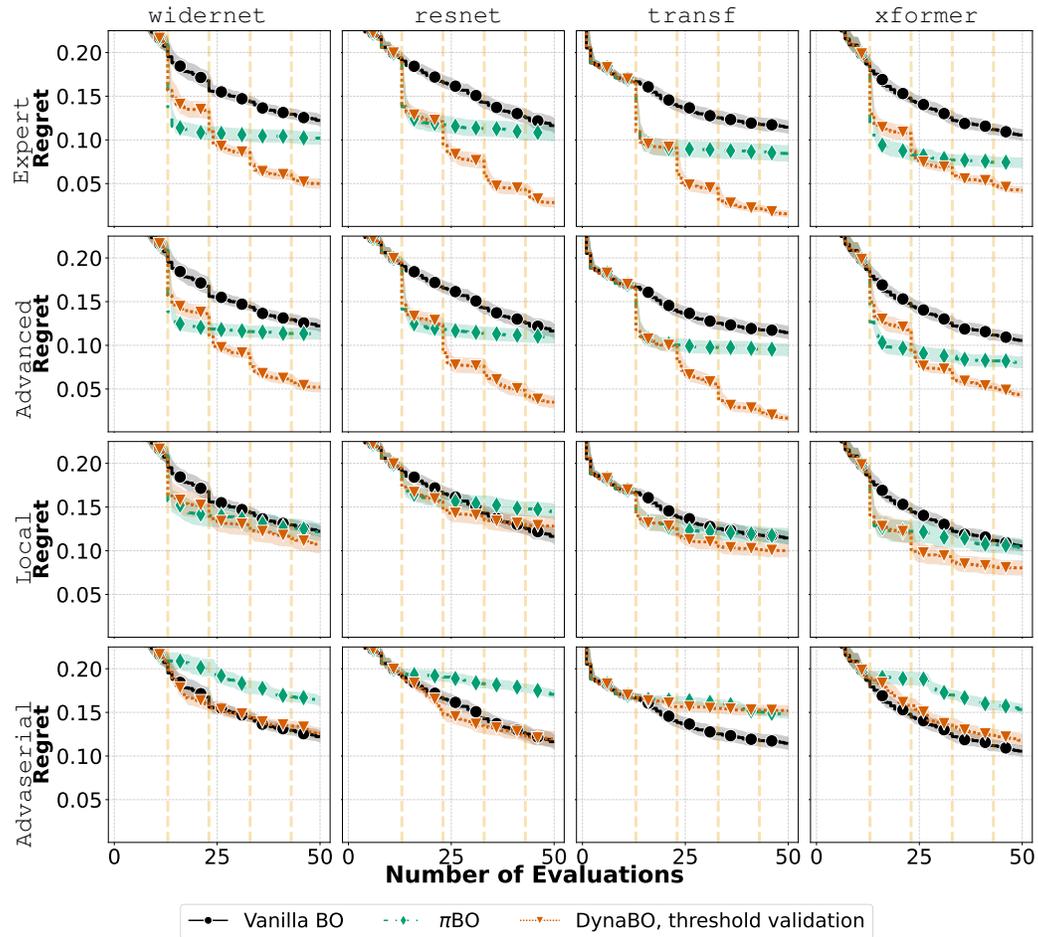


Figure 8: Mean regret for PD1 using Expert, Advanced, Local, and Adversarial priors, with random forests as surrogate models. Priors are provided at vertical lines. The shaded areas visualize the standard error. The results indicate DynaBO outperforming  $\pi$ BO and remaining competitive with vanilla BO for adversarial priors.

### D.3 RANDOM PRIOR LOCATION

For experiments with randomly chosen prior locations, we model user behavior as follows. Each user provides an initial prior at the start of the optimization. If the last prior was given at time  $t_i$ , a new prior is provided at time  $m$  with probability:

$$\mathbb{P}_{t_i}(\pi^{(m)}) = 1 - e^{-0.15 \cdot (m - t_i)}.$$

The resulting outcomes exhibit trends consistent with those observed in the main experiments.

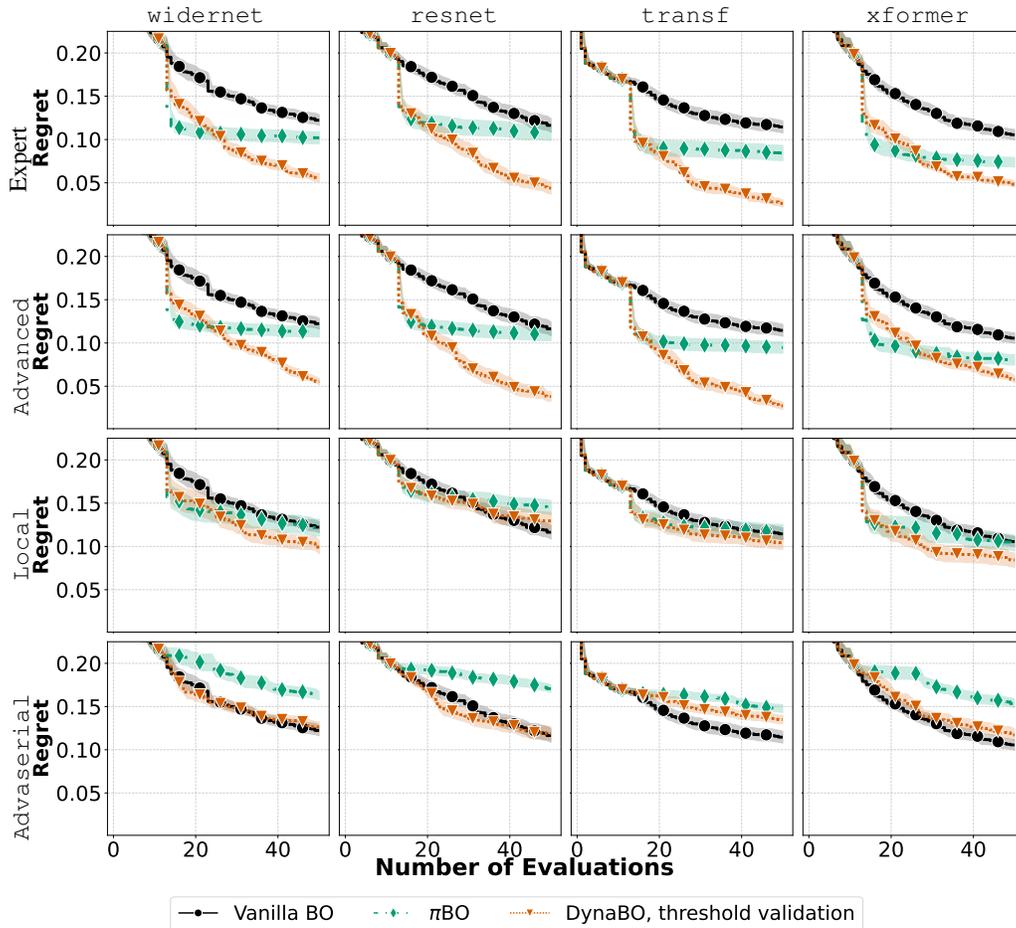


Figure 9: Mean regret for PD1 using Expert, Advanced, Local, and Adversarial priors, with random forests as surrogate models. The shaded areas visualize the standard error. The results indicate DynaBO outperforming  $\pi$ BO and remaining competitive with vanilla BO for adversarial priors.

## D.4 INCREASED BUDGET RESULTS

In our experiments with an increased budget, we introduce a prior every ten trials following the initial design. The indicate that random forests perform poorly without prior rejection, whereas Gaussian processes remain robust. Nevertheless, DynaBO with prior rejection consistently outperforms  $\pi$ BO on most scenarios.

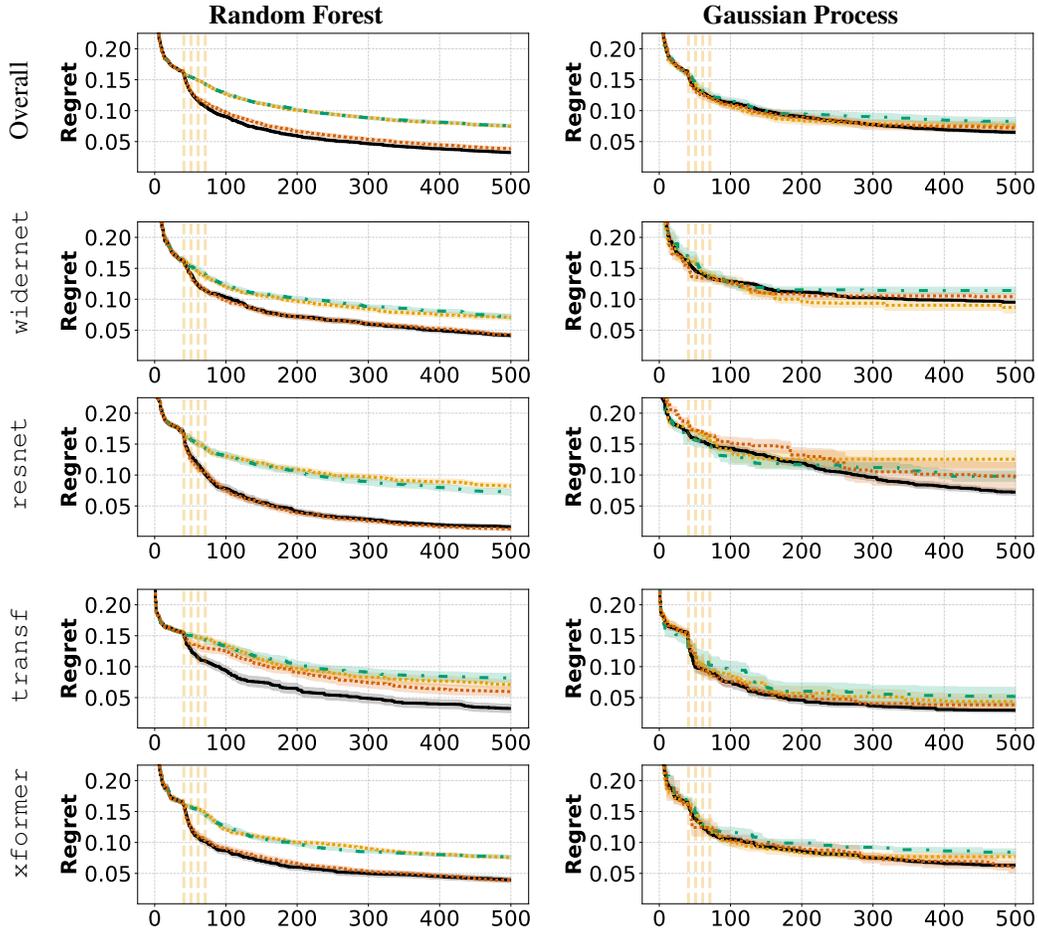


Table 2: Evaluation results across scenarios (rows) for two surrogate models (columns).

D.5 NAÏVE DYNAMIC EXTENSION FOR  $\pi$ BO

In this ablation, we compare a less sophisticated baseline of naively removing old priors to the proposed mechanism of DynaBO for summing priors. The results are visualized in Figure 10. We find that no substantial performance difference can be observed between the two methods in our standard evaluation setup. This means that there is no harm in continuing to use the old priors.

However, to evaluate whether DynaBO’s intuition of old information being useful holds in practice, one has to evaluate the impact of comparing positive and negative priors. To that end, we evaluate what happens if an expert, advanced, or local prior, each provided with a chance of  $1/3$  is followed by an adversarial prior. In this setup, the quality of the two methods differs significantly. For example, when positive priors are followed by negative priors, the results show a degraded performance, as can be seen in Figure 11. This result holds, even though the positive prior results in an immediate performance boost for both approaches.

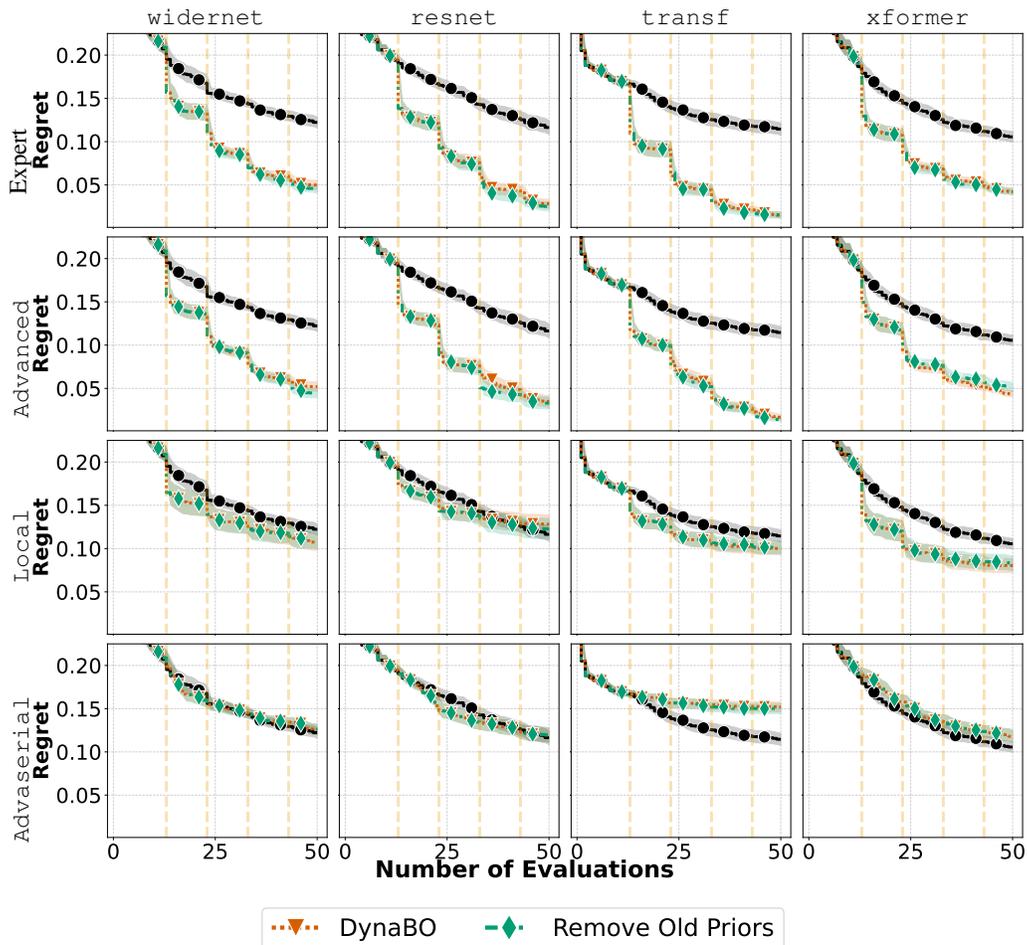


Figure 10: Mean regret for PD1 using Expert, Advanced, Local, and Adversarial priors, with random forests as surrogate models. The shaded areas visualize the standard error. The results indicate DynaBO outperforming  $\pi$ BO and remaining competitive with vanilla BO for adversarial priors.

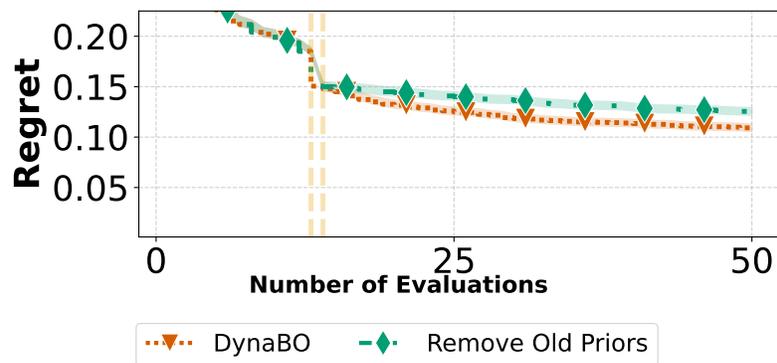


Figure 11: Investigation of helpful priors followed by adversarial priors.

## D.6 SENSITIVITY ANALYSIS OF THE PRIOR REJECTION CRITERION

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

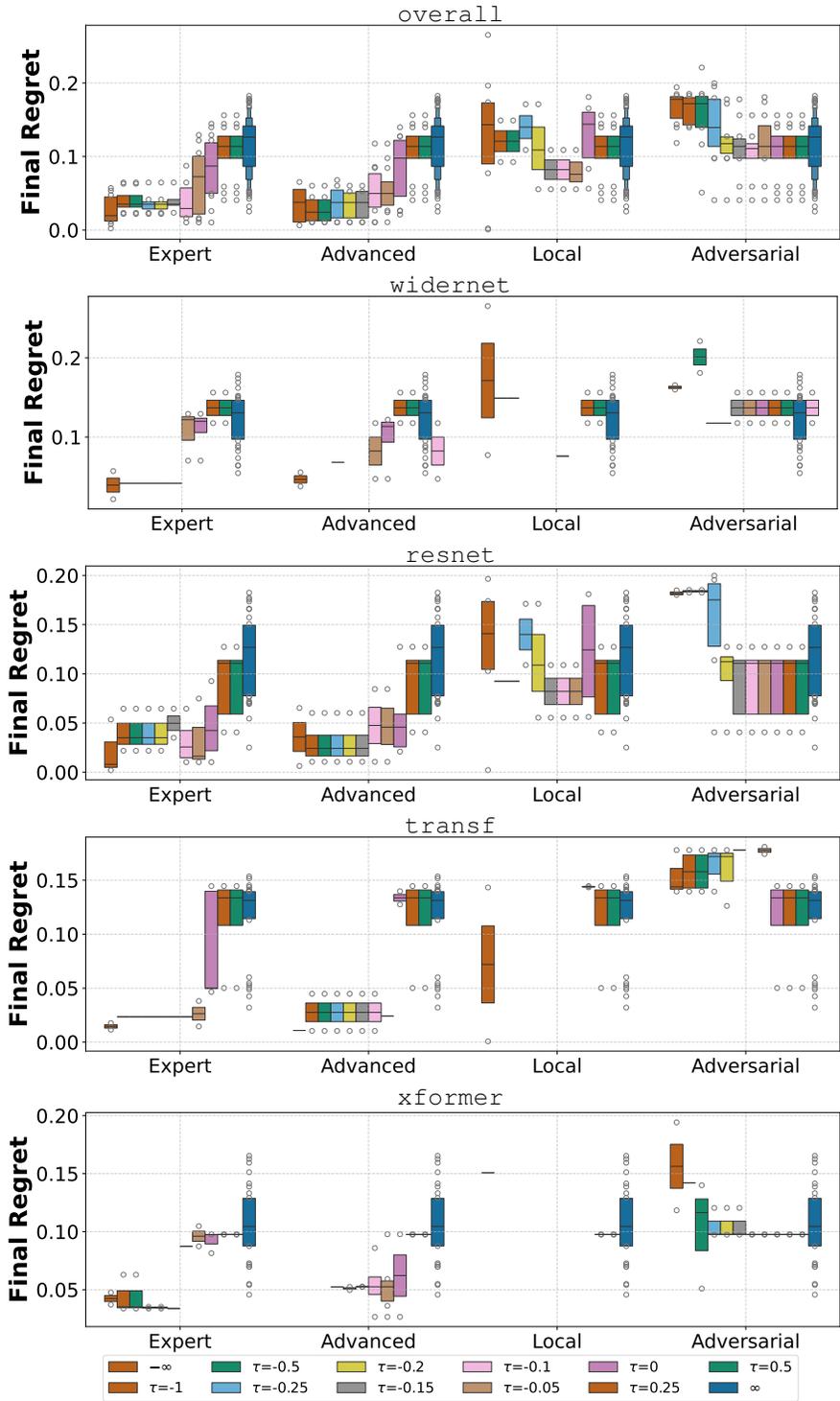


Figure 12: Sensitivity analysis of different thresholds  $\tau$ .  $\tau = -\infty$  accepts all, and  $\tau = \infty$  rejects all priors. *overall* contains the merged results from all scenarios.

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

## D.7 PRIOR DECAY ABLATION

Due to providing multiple priors during the experiment, we reinvestigated the speed at which priors are decayed. We therefore ablated the function  $\phi$ .

$$\alpha_{\hat{f}}^{\text{dyna}}(\lambda) := \alpha_{\hat{f}}(\lambda) \cdot \sum_{m=1}^M \pi^{(m)}(\lambda)^{\beta/\phi(t-t^{(m)})}$$

Our investigation reveals that the optimal decay rate varies according to the prior quality. In our paper, we utilize a linear decay, as it is a good trade-off between good and adversarial priors.

Table 3: Mean regret ( $\mu$ ) and standard error (SE) for each prior type across decay configurations. All values are rounded to three decimal places.

Config	Expert		Advanced		Local		Adversarial	
	$\mu$	SE	$\mu$	SE	$\mu$	SE	$\mu$	SE
Logarithmic Decay	0.027	0.001	0.033	0.001	0.105	0.002	0.177	0.001
Linear Decay	0.028	0.001	0.036	0.001	0.102	0.002	0.162	0.001
Quadratic Decay	0.034	0.001	0.041	0.001	0.103	0.002	0.149	0.001
Cubic Decay	0.038	0.001	0.042	0.001	0.104	0.002	0.146	0.001
To the Power of 4 Decay	0.039	0.001	0.046	0.001	0.101	0.002	0.145	0.001
To the Power of 5 Decay	0.039	0.001	0.045	0.001	0.103	0.001	0.145	0.001

1566 E DECLARATION OF LLM USAGE  
1567

1568 Throughout this submission, we made limited use of Large Language Models (LLMs) in the follow-  
1569 ing ways:  
1570

- 1571 • Code generation from specific instructions, primarily for producing plots and tables.
  - 1572 • Writing support, including translation and alternative phrasings.
  - 1573 • Assistance in locating related research.
- 1574

1575 All conceptual contributions, methodological developments, experimental designs, and analyses  
1576 were carried out solely by the authors.  
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