
FunBO: Discovering Acquisition Functions for Bayesian Optimization with FunSearch

Virginia Aglietti¹ Ira Ktena¹ Jessica Schrouff^{*1} Eleni Sgouritsa¹
Francisco J. R. Ruiz¹ Alan Malek¹ Alexis Bellot¹ Silvia Chiappa¹

Abstract

The sample efficiency of Bayesian optimization algorithms depends on carefully crafted acquisition functions (AFs) guiding the sequential collection of function evaluations. The best-performing AFs can vary significantly across optimization problems, often requiring ad-hoc and problem-specific choices. This work tackles the challenge of designing novel AFs that perform well across a variety of experimental settings. Based on FunSearch, a recent work using Large Language Models (LLMs) for discovery in mathematical sciences, we propose FunBO, an LLM-based method that can be used to learn new AFs written in computer code by leveraging access to a number of evaluations for a limited set of objective functions. We provide the analytic expression of all discovered AFs and evaluate them on various global optimization benchmarks and hyperparameter optimization tasks. We show how FunBO identifies AFs that generalize well both in *and* out of the training distribution of functions, thus outperforming established general-purpose AFs and achieving competitive performance against AFs that are customized to specific function types and are learned via transfer-learning algorithms.

1. Introduction

Bayesian optimization (BO) (Kushner, 1962; 1964; Mockus, 1974; Jones et al., 1998) is a methodology for optimizing complex and expensive-to-evaluate black-box functions that emerge in many scientific disciplines. BO has been used across a wide variety of applications ranging from hyperparameter tuning in machine learning (Bergstra et al., 2011; Snoek et al., 2012; Cho et al., 2020) to designing policies

in robotics (Calandra et al., 2016) and recommending new molecules in drug design (Korovina et al., 2020). Two main components lie at the heart of any BO algorithm: a surrogate model and an acquisition function (AF). The surrogate model expresses assumptions about the objective function, e.g., its smoothness, and it is often given by a Gaussian Process (GP) (Rasmussen & Williams, 2006). Based on the surrogate model, the AF determines the sequential collection of function evaluations by assigning a score to potential observation locations. BO’s success heavily depends on the AF’s ability to efficiently balance exploitation (i.e. assigning a high score to locations that are likely to yield optimal function values) and exploration (i.e. assigning a high score to regions with higher uncertainty about the objective function in order to inform future decisions), thus leading to the identification of the optimum with the minimum number of evaluations.

Existing AFs aim to provide either general-purpose optimization strategies or approaches tailored to specific objective types. For example, Expected Improvement (EI) (Šaltanis, 1971; Mockus, 1974), Upper Confidence Bound (UCB) (Lai & Robbins, 1985) and Probability of Improvement (PofI) (Kushner, 1962; 1964) are all widely adopted *general-purpose* AFs that can be used out-of-the-box across BO algorithms and objective functions. The performance of these AFs varies significantly across different types of black-box functions, making the AF choice an ad-hoc, empirically driven, decision. There exists an extensive literature on alternative AFs outperforming EI, UCB and PofI, for instance entropy-based (Wang & Jegelka, 2017) or knowledge-gradient (Šaltanis, 1971; Frazier et al., 2008) optimizers, see Garnett (2023, Chapter 7) for a review. However, while these functions are often interpretable, they are generally hard to implement and expensive to evaluate, partly defeating the purpose of replacing the expensive original optimization with the optimization of a much cheaper and faster to evaluate AF. In order to avoid the limitations of current AFs, several works have proposed self-adjusting the hyperparameters of known AFs in a data driven way throughout the optimization process (Benjamins et al., 2023; Ding et al., 2022), combining different AFs in a portfolio and selecting them via an online multi-armed bandit strategy (Hoffman

^{*}Now at GlaxoSmithKline. ¹Google DeepMind, London, UK. Correspondence to: Virginia Aglietti <aglietti@google.com>.

et al., 2011) or computing a Pareto front over different AFs (Lyu et al., 2018). Other prior works (Hsieh et al., 2021; Volpp et al., 2020; Wistuba & Grabocka, 2021) have instead proposed representing AFs via neural networks thus bypassing the need for an analytical representation and learning new AFs tailored to specific objectives by transferring information from a set of related training functions via, e.g., reinforcement learning or transformers. While such learned AFs can outperform general-purpose AFs, their generalization performance to objectives outside of the training distribution is often poor (see experimental section and discussion on generalization behavior in Volpp et al. (2020)). More recently, the concurrent work of Yao et al. (2024) investigated representing AFs in code for specific optimization settings where the experimentation budget is limited.

Defining methodologies that *automatically* identify new AFs capable of outperforming general-purpose *and* function-specific alternatives, both in and out of the training distribution, remains a significant and unaddressed challenge. In this work we tackle this challenge by considering AFs represented in computer code. Learning new AFs expressed in code presents three main difficulties: (i) the vast space of all possible programs makes exhaustive search infeasible, (ii) efficiently exploring a constrained space of possible programs requires scalable methods and (iii) there is no clear criterion for ensuring the validity and effectiveness of generated AFs.

Contributions. We overcome these difficulties by formulating the problem of learning novel AFs written in computer code as an algorithm discovery problem and address it by extending FunSearch (Romera-Paredes et al., 2023), a recently proposed algorithm solving open problems in mathematical sciences via LLMs. In particular, we introduce FunBO, a novel method that explores the large space of AFs written in computer code by taking an initial AF as input and, with a limited number of evaluations for a set of objective functions, iteratively modifying it to improve the performance of the resulting BO algorithm. We focus on Python programs but FunBO can be readily applied to languages supported by FunSearch.

Unlike existing algorithms, FunBO outputs code snippets corresponding to improved AFs. This approach offers several advantages: (i) *interpretability*, as the code-based AFs can be directly inspected, analyzed, and understood. This allows for a clear comprehension of the logic behind the optimization strategy employed by the discovered AF, which contrasts with neural network-based AFs for which it is often difficult to understand why they perform well or how they balance exploration and exploitation; (ii) *deployability and simplicity*. The AFs discovered by FunBO are output as concise code snippets (e.g., Python functions) that can be readily integrated into existing BO libraries and work-

flows with minimal overhead; (iii) *leveraging foundational knowledge* embedded in LLMs. Having been trained on vast corpora of code and text, these models possess knowledge about both programming constructs and existing BO strategies, which FunBO utilizes to efficiently explore the program space and construct effective heuristics.

We extensively test FunBO on a range of optimization problems including standard global optimization benchmarks and hyperparameter optimization (HPO) tasks. For each experiment, we report the explicit functional form of the discovered AFs and show that they generalize well to the optimization of functions both in and out of the training distribution, outperforming general-purpose AFs while comparing favorably to function-specific ones.

2. Preliminaries

We consider an expensive-to-evaluate black-box function $f : \mathcal{X} \rightarrow \mathbb{R}$ over the input space $\mathcal{X} \subseteq \mathbb{R}^d$ for which we aim to identify the global minimum $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$. We assume access to a set of auxiliary black-box and expensive-to-evaluate objective functions, $\mathcal{G} = \{g_j\}_{j=1}^J$, with $g_j : \mathcal{X}_j \rightarrow \mathbb{R}$, $\mathcal{X}_j \subseteq \mathbb{R}^{d_j}$ for $j = 1, \dots, J$, from which we can obtain a set of evaluations.

Bayesian optimization. BO seeks to identify \mathbf{x}^* with the smallest number T of sequential evaluations of f given N initial observations $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$, with $y_i = f(\mathbf{x}_i)$.¹ BO relies on a probabilistic surrogate model for f which in this work is set to a GP with prior distribution over any batch of input points $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ given by $p(f|\mathbf{X}) = \mathcal{N}(m(\mathbf{X}), K_\theta(\mathbf{X}, \mathbf{X}'))$ with prior mean $m(\mathbf{X})$ and kernel $K_\theta(\mathbf{X}, \mathbf{X}')$ with hyperparameters θ . The posterior distribution $p(f|\mathcal{D})$ is available in closed form via standard GP updates. At every step t in the optimization process, BO selects the next evaluation location by optimizing an AF $\alpha(\cdot|\mathcal{D}_t) : \mathcal{X} \rightarrow \mathbb{R}$, given the current posterior distribution $p(f|\mathcal{D}_t)$, with \mathcal{D}_t denoting the function evaluations collected up to trial t (including \mathcal{D}). A commonly used AF is the Expected Improvement (EI), which is defined as $\alpha_{\text{EI}}(\mathbf{x}|\mathcal{D}_t) = (y^* - m(\mathbf{x}|\mathcal{D}_t))\Phi(z) + \sigma(\mathbf{x}|\mathcal{D}_t)\phi(z)$, where y^* denotes the best function value observed in \mathcal{D}_t , also called incumbent, $z = (y^* - m(\mathbf{x}|\mathcal{D}_t))/\sigma(\mathbf{x}|\mathcal{D}_t)$, ϕ and Φ are the standard Normal density and distribution functions, and $m(\mathbf{x}|\mathcal{D}_t)$ and $\sigma(\mathbf{x}|\mathcal{D}_t)$ are the GP posterior mean and standard deviation computed at $\mathbf{x} \in \mathcal{X}$. Other general-purpose AFs proposed in the literature are: UCB ($\alpha_{\text{UCB}}(\mathbf{x}|\mathcal{D}_t) = m(\mathbf{x}|\mathcal{D}_t) - \beta\sigma(\mathbf{x}|\mathcal{D}_t)$ with hyperparameter β), PoFI ($\alpha_{\text{PoFI}}(\mathbf{x}|\mathcal{D}_t) = \Phi((y^* - m(\mathbf{x}|\mathcal{D}_t))/\sigma(\mathbf{x}|\mathcal{D}_t))$) and the posterior mean $\alpha_{\text{MEAN}}(\mathbf{x}|\mathcal{D}_t) = m(\mathbf{x}|\mathcal{D}_t)$ (denoted

¹We focus on noiseless observations but the method can be equivalently applied to noisy outcomes.

by MEAN hereinafter).²

Unlike general-purpose AFs, several works have proposed increasing the efficiency of BO for a specific optimization problem, say the optimization of f , by either adaptively selecting and/or adjusting known AFs in a data-driven manner (Benjamins et al., 2023) or by *learning* problem-specific AFs (Hsieh et al., 2021; Volpp et al., 2020; Wistuba & Grabocka, 2021). The learned AFs are trained on the set \mathcal{G} , whose functions are assumed to be drawn from the same distribution or function class associated with f , reflecting a meta-learning setup. “Function class” here refers to a set of functions with a shared structure and obtained by, e.g., applying scaling and translation transformations to their input and output values or evaluating the loss function of the same machine learning model, e.g., AdaBoost, on different data sets. For instance, Wistuba et al. (2018) learns an AF that is a weighted superposition of EIs by exploiting access to a sufficiently large dataset for functions in \mathcal{G} . Volpp et al. (2020) considered settings where the observations for functions in \mathcal{G} are limited and proposed MetaBO, a reinforcement learning based algorithm that learns a specialized neural AF, i.e., a neural network representing the AF. The neural AF takes as inputs a set of potential locations (with a given d), the posterior mean and variance at those points, the trial t and the budget T and is trained using a proximal policy optimization algorithm (Schulman et al., 2017). Similarly, Hsieh et al. (2021) proposed FSAF, an AF obtained via few-shot adaptation of a learned AF using a small number of function instances in \mathcal{G} .

Note that, while general-purpose AFs are used to optimize objectives across function classes, learned AFs aim to achieve high performance for the single function class to which f and \mathcal{G} belong. See Section A for an extensive discussion of related work.

FunSearch. FunSearch (Romera-Paredes et al., 2023) is a recently proposed evolutionary algorithm for *searching* in the *functional* space by combining a pre-trained LLM used for generating new computer programs with an efficient evaluator, which guards against hallucinations and scores fitness. An example problem that FunSearch tackles is the online bin packing problem (Coffman et al., 1984), where a set of items of various sizes arriving online needs to be packed into the smallest possible number of fixed size bins.

A set of heuristics has been designed for deciding which bin to assign an incoming item to, e.g., “first fit.” FunSearch aims to discover new heuristics that improve on existing ones by taking as inputs: (i) the computer code of an `evolve` function $h(\cdot)$ representing the initial heuristic to be improved by the LLM, e.g., “first fit” and (ii) an `evaluate` function $e(h, \cdot)$, also written in computer code, specifying

²We focus on AFs that can be evaluated in closed form given the posterior parameters of a GP surrogate model and exclude those whose computation involves approximations.

the problem at hand (also called “problem specification”) and scoring each $h(\cdot)$ according to a predefined performance metric, e.g., the number of bins used in $h(\cdot)$. The inputs of both $h(\cdot)$ (denoted by h hereinafter) and $e(h, \cdot)$ (denoted by e hereinafter), are problem specific. A description of h ’s inputs is provided in the function’s docstring³ together with an explanation of how the function itself is used within e .

Given these initial components, FunSearch prompts an LLM to propose an improved h , scores the proposals on a set of inputs, e.g., on different bin-packing instances, and adds them to a programs database. The programs database stores correct h functions⁴ together with their respective scores. In order to encourage diversity of programs and enable exploration of different solutions, a population-based approach inspired by genetic algorithms (Tanese, 1989) is adopted for the programs database (DB). At a subsequent step, functions in the database are sampled to create a new prompt, the LLM’s proposals are scored and stored again. The process repeats for $\tau = 1, \dots, T$ until a time budget T is reached and the heuristic with the highest score on a set of inputs is returned.

3. FunBO

FunBO is a FunSearch-based method for discovering novel AFs that increase BO efficiency by exploiting the set of auxiliary objectives \mathcal{G} . In particular, FunBO (i) uses the same prompt and DB structure as FunSearch, but (ii) proposes a new problem specification by viewing the learning of AFs as an algorithm discovery problem, and (iii) introduces a novel initialization and evaluation pipeline that is used within the FunSearch structure. FunBO does not make assumptions about similarities between f and \mathcal{G} , nor does it assume access to a large dataset for each function in \mathcal{G} . Therefore, FunBO can be used to discover both general-purpose and function-specific AFs as well as to adapt AFs via few-shots. FunBO leverages the LLMs’ ability to generate executable code to make the search for novel AFs automatic and scalable, potentially leveraging the extensive LLMs’ knowledge of BO and AFs while delivering more interpretable AFs than those represented by neural networks. Furthermore, while FunSearch was only applied to problems that required evolving functions with simple inputs (integers, floats or short tuples; with only one application taking as input a single array), FunBO explores a significantly more complex function space where programs take as inputs multiple arrays. This demonstrates how the same formulation can be applied to problems of increasing complexity as long as an appropriate scoring mechanism is identified.

³wide variety on Python programs.

⁴The definition of a correct function is also problem specific. For instance, a program can be considered correct if it compiles.

Inputs: $\mathcal{G}_{Tr}, \mathcal{G}_V, N_{DB}, B, \mathcal{T}$

Setup: Initialize h (Fig. 6), e (Fig. 7-8) and DB with N_{DB} islands. Assign h to each island.

while $\tau < \mathcal{T}$ **do**

1. Sample two programs from DB and create prompt (Fig. 2, right)
2. Get a batch of B samples from the LLM
3. For each correct h^τ in the batch compute $s_{h^\tau}(\mathcal{G}_{Tr})$
4. Add correct h^τ to DB and update it (see Appendix C)
5. Update step $\tau = \tau + 1$

end

Output: Return h in DB with score in the top 20th percentile for \mathcal{G}_{Tr} and highest score on \mathcal{G}_V .

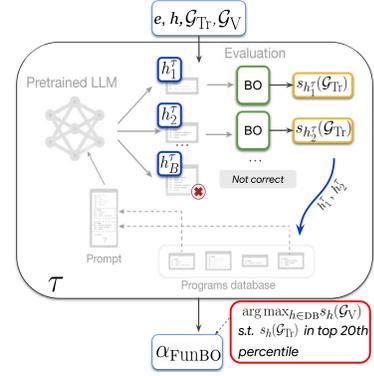


Figure 1. *Left:* The FunBO algorithm. *Right:* Graphical representation of FunBO. The different FunBO component w.r.t. FunSearch (Romera-Paredes et al., 2023, Fig. 1) are highlighted in color.

Method overview. FunBO sequentially prompts an LLM to improve an **initial AF** expressed in code so as to enhance the performance of the corresponding BO algorithm when optimizing objectives in \mathcal{G} . At every step τ of FunBO, an LLM’s **prompt** is created by including the code for two AF instances generated and stored in a **programs database** (DB) at previous iterations. With this prompt, a number (B) of alternative AFs are sampled from the LLM and are evaluated based on their average performance on a subset $\mathcal{G}_{Tr} \subseteq \mathcal{G}$, which acts as training dataset. The **evaluation** process for an AF, say h^τ at step τ , on \mathcal{G}_{Tr} gives a numeric score $s_{h^\tau}(\mathcal{G}_{Tr})$ that is used to store programs in DB and sample them for subsequent prompts. The “process” of prompt creation, LLM sampling, and AF scoring and storing repeats until a time budget \mathcal{T} is reached. Out of the top performing⁵ AFs on \mathcal{G}_{Tr} , the algorithm returns the AF performing the best, on average, in the optimization of $\mathcal{G}_V = \mathcal{G} \setminus \mathcal{G}_{Tr}$, which acts as a validation dataset. When no validation functions are used ($\mathcal{G} = \mathcal{G}_{Tr}$), the AF with the highest average performance on \mathcal{G}_{Tr} is returned. Each FunBO component highlighted in bold is described below in more detail, along with the complete algorithm and graphical representation in Fig. 1. We denote the AF returned by FunBO as α_{FunBO} .

Initial AF. FunBO’s initial program h determines the input variables that can be used to generate alternative AFs while imposing a prior on the programs the LLM will generate at successive steps. For these reasons it is important for guiding the search process effectively. We consider an initial acquisition function that takes the functional form of the EI and has as inputs the union of the inputs given to EI, UCB and PofI (see code in Fig. 6). The AF returns an integer representing the index of the point in a vector of potential locations that should be selected for the next function evaluation. All programs generated by the LLM share the same inputs and output, but vary in their implementation, which

defines different optimization strategies, see for instance the AF generated for one of our experiments in Fig. 3 (left).⁶

Prompt. At every algorithm iteration, a prompt is constructed by sampling two AFS, h_i and h_j , previously generated and stored in DB. h_i and h_j are sampled from DB in a way that favors higher scoring and shorter programs (see paragraph below for more details) and are sorted in the prompt in ascending order based on their scores $s_{h_i}(\mathcal{G}_{Tr})$ and $s_{h_j}(\mathcal{G}_{Tr})$, see the prompt skeleton⁷ in Fig. 2 (bottom). The LLM is then asked to generate a new AF representing an improved version of the last, higher scoring, program.

Evaluation. As expected, the evaluation protocol is critical for the discovery of appropriate AFS. Indeed, a scoring mechanism that captures small improvements in the proposed AF is needed to steer the LLM toward promising regions of the function space. Our novel evaluation setup, unlike the one used in FunSearch, entails performing a full BO loop to evaluate program fitness. In particular, each function generated by the LLM is (i) checked to verify it is correct, i.e., it compiles and returns a numerical output; (ii) scored based on the average performance of a BO algorithm using h^τ as an AF on \mathcal{G}_{Tr} . Evaluation is performed by running a full BO loop with h^τ for each function $g_j \in \mathcal{G}_{Tr}$ and computing a score that contains two terms: a term that rewards AFS finding values close to the true optimum, and a term that rewards AFS finding the optimum in fewer evaluations (often called

⁶We explored using a random selection of initial points as an alternative to EI. However, this approach did not yield good results as using a random selection was incentivizing the generation of functions with a stochastic output, for which convergence results are not reproducible.

⁷When $\tau = 1$, only the initial program will be available in DB thus the prompt in Fig. 2 will be simplified by removing `acquisition_function_v1` and replacing `v_2` with `v_1`.

⁵In this work we consider the programs with a score in the top 20th percentile.

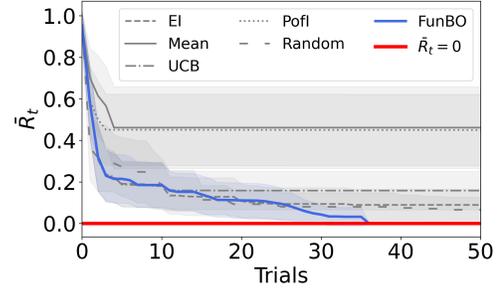
```

"""Improve Bayesian Optimization by discovering a new
acquisition function."""

def acquisition_function_v0(predictive_mean, predictive_var,
incumbent, beta=1.0):
    """Returns the index of the point to collect ...
    (Full docstring in Fig. 8)"""
    # Code for lowest-scoring sampled AF.
    return ...

def acquisition_function_v1(predictive_mean, predictive_var,
incumbent, beta=1.0):
    """Improved version of 'acquisition_function_v0'."""
    # Code for highest-scoring sampled AF.
    return ...

def acquisition_function_v2(predictive_mean, predictive_var,
incumbent, beta=1.0):
    """Improved version of the previous
    'acquisition_function'."""
    
```



```

def acquisition_function(predictive_mean,
predictive_var, incumbent, beta=1.0):
    """Returns the index of the point to collect
    ... (Full docstring in Fig. 8)."""
    z = (incumbent - predictive_mean) /
    np.sqrt(predictive_var)
    predictive_std = np.sqrt(predictive_var)
    vals = (incumbent - predictive_mean)
    * norm.cdf(z) + predictive_std * norm.pdf(z)
    return np.argmax(vals)
    
```

Figure 2. *Left*: FunBO prompt includes two previously generated AFs which are sampled from DB and are sorted in ascending order based on the score achieved on \mathcal{G}_{Tr} . The LLM generates a third AF, `acquisition_function_v2`, representing an improved version of the highest scoring program. *Right*: OOD-Bench. Code for α_{FunBO} (bottom) and average BO performance when using general purpose AFs and α_{FunBO} (top). Shaded area gives \pm std. dev./2. The red line gives $\bar{R}_t = 0$, i.e. zero average regret.

trials). Specifically, we use the score $s_{h^\tau}(\mathcal{G}_{Tr})$ defined as :

$$\frac{1}{|\mathcal{G}_{Tr}|} \sum_{j=1}^J \left[\left(1 - \frac{g_j(\mathbf{x}_{j,h^\tau}^*) - y_j^*}{g_j(\mathbf{x}_{j,t=0}^*) - y_j^*} \right) + \left(1 - \frac{T_{h_j^\tau}}{T} \right) \right] \quad (1)$$

where, for each g_j , y_j^* is the ground truth optimal value, $\mathbf{x}_{j,t=0}^*$ gives the optimum found at $t = 0$ ⁸, \mathbf{x}_{j,h^τ}^* is the found optimum for g_j with h^τ and $T_{h_j^\tau}$ gives the number of trials out of T that h^τ selected before reaching y_j^* (if the optimum for function g_j was not found, then $T_{h_j^\tau} = T$ to indicate that all available trials have been used). The first term in the square brackets of Eq. (1) quantifies the discrepancy between the function values at the returned optimum and the true optimum. This term becomes zero when \mathbf{x}_{j,h^τ}^* equals $\mathbf{x}_{j,t=0}^*$, indicating a failure to explore the search space. Conversely, if h^τ successfully identifies the true optimum, such that $g_j(\mathbf{x}_{j,h^\tau}^*) = y_j^*$, this term reaches its maximum value of one. The second term in Eq. (1) captures how quickly h^τ identifies y_j^* . When $T_{h_j^\tau} = T$, indicating the algorithm has not converged, this term becomes zero, and the score is solely determined by the discrepancy between the discovered and true optimum. If, instead, the algorithm reaches the global optimum, this term represents the proportion of trials, out of the total budget T , needed to do so. Code for

⁸This is the input value corresponding to the optimal value found among the set of initial observations \mathcal{D} . We assume this to be different from the ground truth optimum.

the evaluation process is presented in Appendix B.

Programs database. Similar to FunSearch, scored AFs are added to DB, which keeps a population of correct programs following an island model (Tanese, 1989). DB is initialized with a number N_{DB} of islands that evolve independently. Sampling of h_i and h_j from DB is done by first uniformly sampling an island and, within that island, sampling programs by favouring those that are shorter and higher scoring. A new program generated when using h_i and h_j in the prompt is added to the same island and, within that, to a cluster of programs performing similarly on \mathcal{G}_{Tr} , see Appendix C for more details.

4. Experiments

Our experiments explore FunBO’s ability to generate novel and efficient AFs across a wide variety of settings. In particular, we demonstrate its potential to generate AFs that generalize well to the optimization of functions both in distribution (ID, i.e. within function classes) and out of distribution (OOD, i.e. across function classes) by running three different types of experiments:

1. OOD-Bench tests generalization across function classes by running FunBO with \mathcal{G} containing different standard global optimization benchmarks and testing on a set \mathcal{F}

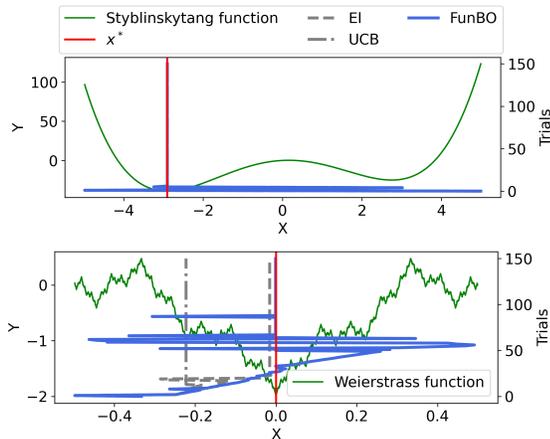


Figure 3. OOD-Bench. Different AFs trading-off exploration and exploitation for two one-dimensional objective functions (green lines). Blue and gray trajectories track the points queried by α_{FunBO} , EI and UCB over 150 trials (right y -axis).

that similarly comprises diverse functions in terms of smoothness, input ranges and dimensionality and output magnitudes. We do not scale the output values nor normalize the input domains to facilitate learning, but rather use the objective functions as available in standard BO packages out-of-the-box. In this case \mathcal{G} and \mathcal{F} do not share any particular structure, thus the generated AFs are closer to general-purpose AFs.

2. ID-Bench, HPO-ID and GPS-ID test FunBO-generated AFs within function classes for standard global optimization benchmarks, HPO tasks, and general function classes, respectively. As this setting is closer to the one considered by meta-learning approaches introduced in Section 2, we compare FunBO against MetaBO (Volpp et al., 2020),⁹ the state-of-the-art transfer AF.
3. FEW-SHOT demonstrates how FunBO can be used in the context of few-shot fast adaptation of an AF. In this case, the AF is learned using a general function class as \mathcal{G} and is then tuned, using a very small (5) number of examples, to optimize a specific synthetic function. We compare our approach to Hsieh et al. (2021),¹⁰ the most relevant few-shot learning method.

We report all results in terms of normalized average simple regret on a test set, \bar{R}_t , as a function of the trial t . For an objective function f , this is defined as $R_t = f(\mathbf{x}_t^*) - y^*$ where y^* is the ground truth optimal value and \mathbf{x}_t^* is the best selected point within the data collected up to t . As \mathcal{F}

⁹We used the author-provided implementation at github.com/boschresearch/MetaBO.

¹⁰We used the author-provided implementation at github.com/pinghsieh/FSAF.

might include functions with different scales, we normalize the regret values to be in $[0, 1]$ before averaging them. Note that the mean and standard deviations shown in all regret plots represent the performance variation across the set of test functions, not across multiple independent runs of a BO algorithm on a single function instance. For instance, for OOD-Bench, the mean and standard deviation in Fig. 2 (right, top) are over the 9 distinct test functions.

To isolate the effects of different AFs, we employ the same experimental setting across all methods in terms of: (i) the number of trials T ; (ii) the hyperparameters of the GP surrogate models (tuned offline); (iii) the evaluation grid for the AF, which is set to be a Sobol grid (Sobol', 1967) on the input space; and (iv) the initial design, which includes the input point yielding the maximum function value on the grid. Note that we use a GP model with zero mean function and RBF kernel across experiments. Therefore, the discovered AFs are conditioned on this choice of surrogate model. We acknowledge that while this setup ensures a consistent comparison across all AFs by removing confounding effects (e.g., from the quality of hyperparameter optimization routines), it might deviate from a fully realistic deployment scenario. For this reason, in Section D.6, we included a comparison of the AFs found by FunBO for OOD-Bench when evaluated using the standard BO pipeline available in BoTorch, which employs default settings for aspects such as AF optimization, GP hyperparameter optimization, and the random selection of initial points.

All experiments are conducted using FunSearch with default hyperparameters¹¹ unless otherwise stated. We employ Codey, an LLM fine-tuned on a large code corpus and based on the PaLM model family (Google-PaLM-2-Team, 2023), to generate AFs.¹²

OOD-Bench. We test FunBO’s capabilities to generate AFs that perform well *across* function classes by including the one-dimensional functions Ackley, Levy, and Schwefel in \mathcal{G}_{Tr} and using the one-dimensional Rosenbrock function for \mathcal{G}_{V} . We test the resulting α_{FunBO} on nine very different objective functions: Sphere ($d = 1$), Styblinski-Tang ($d = 1$), Weierstrass ($d = 1$), Beale ($d = 2$), Branin ($d = 2$), Michalewicz ($d = 2$), Goldstein-Price ($d = 2$) and Hartmann with both $d = 3$ and $d = 6$. We do not compare against MetaBO as (i) it was developed for settings in which the functions in \mathcal{G} and \mathcal{F} belong to the same class and, (ii) the neural AF is trained with evaluation points of a given dimension, thus it cannot be deployed for the optimization of functions across different d . For completeness, we re-

¹¹See code at github.com/google-deepmind/funsearch.

¹²Codey is publicly accessible via its API (Vertex AI, 2023). For AF sampling, we used 5 Codey instances running on tensor processing units on a computing cluster. For scoring, we used 100 CPU evaluators per LLM instance.

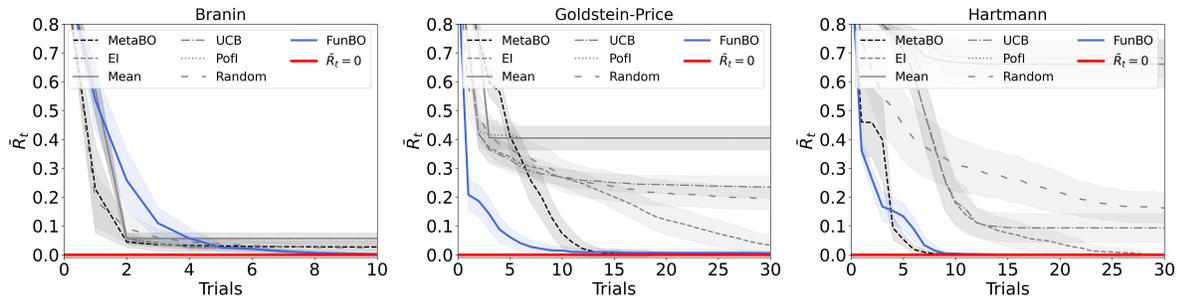


Figure 4. ID-Bench. Average BO performance when using general-purpose AFs (gray lines), the AF learned by MetaBO (black dashed line) and α_{FunBO} (blue line) on 100 function instances. Shaded area gives \pm std. dev./2. The red line represents $\bar{R}_t = 0$, i.e. zero average regret.

port a comparison with a dimensionality-agnostic version of MetaBO in Appendix D.1 (Fig. 9) together with all experimental details, e.g., input ranges and hyperparameters.

AF interpretation: In this experiment, α_{FunBO} (Fig. 2, right bottom plot) represents a combination of EI and UCB which, due to the `beta*predictive_std` term, is more exploratory than EI but, considering the incumbent value, still factors in the expected magnitude of the improvement and reduces to EI when `beta=0`. This determines the way α_{FunBO} trades-off exploration and exploitation which can be visualized by looking at the "exploration path", i.e., the sequence of \mathbf{x} values selected over t , as shown in the right plots of Fig. 3 (t measured on the secondary y-axis). For objective functions that are smooth, for example Styblinski-Tang, the exploration path of α_{FunBO} matches those of EI and UCB (top plot, note that trajectories are overlapping). In this scenario, all AFs exhibit similar behavior, converging to \mathbf{x}^* (red vertical line) with fewer than 25 trials. When instead the objective function has a lot of local optima (bottom plot) as in Weierstrass, both EI and UCB get stuck after a few trials while FunBO continues to explore the search space eventually converging to \mathbf{x}^* . Notice how in this plot the convergence paths of all AFs differ and only the blue line aligns with the red line, i.e., converges to \mathbf{x}^* , after a few trials.

Using α_{FunBO} to optimize the nine functions in \mathcal{F} leads to a fast and accurate convergence to the global optima (Fig. 2, right top plot). This is confirmed when extending the test set to include 50 scaled and translated instances of the functions in \mathcal{F} (Fig. 9, right). Finally, note how the top right plot in Fig. 2 shows a surprisingly good performance of random search. This is due to random search performing competitively on functions with numerous local optima, which are generally harder to optimize. Aggregating performance across all functions in \mathcal{F} highlights that *no single known general-purpose AF consistently outperforms the others*. This aligns with the well-established understanding that the effectiveness of AFs can vary significantly across black-box functions and is consistent with findings reported

in the literature (Perrone et al., 2019; Li et al., 2018).

ID-Bench. Next we evaluate FunBO capabilities to generate AFs that perform well *within* function classes using Branin, Goldstein-Price and Hartmann ($d = 3$). For each of these three functions, we train both FunBO and MetaBO with $|\mathcal{G}| = 25$ instances of the original function obtained by scaling and translating it with values in $[0.9, 1.1]$ and $[-0.1, 0.1]^d$ respectively.¹³ For FunBO we randomly assign 5 functions in \mathcal{G} to \mathcal{G}_V and keep the rest in \mathcal{G}_{Tr} . We test the performance of the learned AFs on another 100 instances of the same function, with randomly sampled values of scale and translation from the same ranges. We additionally compare against a BO algorithm that uses EI, UCB, PofI, MEAN or a random selection of points. All hyper-parameter settings for this experiment are provided in Appendix D.2. Across all objective functions, α_{FunBO} leads to a convergence performance that outperforms general-purpose AFs (Fig. 4). More importantly, despite using the same inputs of EI or UCB, FunBO is able to reach performances that are comparable or superior to those of AFs that are parameterized by neural networks and use additional inputs (Fig. 4). In terms of interpretability, notice how the AF for Goldstein-Price (Fig. 12) can be written as $\sigma^2(\mathbf{x}|\mathcal{D}_t)\Phi(\frac{y^* - \mu(\mathbf{x}|\mathcal{D}_t)}{\sigma(\mathbf{x}|\mathcal{D}_t)})$ thus giving a modified PofI where the probability of observing an improvement over the incumbent is multiplied by the predictive variance.

The AFs found in this experiment (code in Figs. 11-13) are "customized" to a given function class thus being closer, in spirit, to the transfer AF. However, in order to further validate the generalizability of α_{FunBO} found in OOD-Bench, we tested this AF across instances of Branin, Goldstein-Price and Hartmann (Fig. 10, green line). We found it to perform well against general-purpose AFs thus confirming the strong results observed in OOD-Bench while being, as expected, slower than AFs customized to a specific objective.

HPO-ID. We test FunBO on two HPO tasks where the goal

¹³Throughout the paper we adopt MetaBO's translation and scaling ranges.

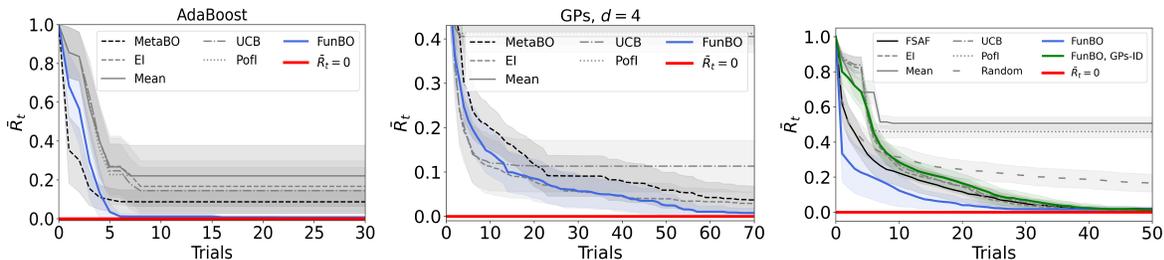


Figure 5. Average BO performance when using known general-purpose AFs (gray lines), the AF learned by MetaBO (black dashed line) and α_{FunBO} (blue line). Shaded area gives \pm std. dev./2. The red line represents $\bar{R}_t = 0$, i.e. zero average regret. *Left*: HPO-ID. *Middle*: GPs-ID with $d = 4$. *Right*: FEW-SHOT.

is to minimize the loss ($d = 2$) of an RBF-based SVM and an AdaBoost algorithm.¹⁴ As in ID-Bench, we test the ability to generate AFs that generalize well within function classes. Therefore, we train FunBO and MetaBO with losses computed on a random selection of 35 of the 50 available datasets and test on losses computed on the remaining 15 datasets. For FunBO we randomly assign 5 datasets to \mathcal{G}_V and keep the rest in \mathcal{G}_{Tr} . FunBO identifies AFs (code in Fig. 15-16) that outperform all other AFs in AdaBoost (Fig. 5, left) while performing similarly to general-purpose or meta-learned AFs for SVM (Fig. 14). Across the two tasks, α_{FunBO} found in OOD-Bench still outperforms general-purpose AFs while yielding slightly worse performance compared to MetaBO and FunBO-customized AFs (Fig. 14, green lines).

GPs-ID. Similar results are obtained for general function classes whose members do not exhibit any particular shared structure. We let \mathcal{G}_{Tr} include 25 functions sampled from a GP prior with $d = 3$, RBF kernel and length-scale drawn uniformly from $[0.05, 0.5]$. We test the found AF on 100 other GP samples defined both for $d = 3$ and $d = 4$ and length-scale values sampled similarly. As done by Volpp et al. (2020), we consider a dimensionality-agnostic version of MetaBO that allows deploying the function learned from $d = 3$ functions on $d = 4$ objectives. We found α_{FunBO} to outperform all other AFs (code in Fig. 18) in $d = 4$ (Fig. 5, right) while matching EI and outperforming MetaBO in $d = 3$ (Fig. 17, left). In terms of interpretability, note how the AF found in this case can be simplified to be written as $(\text{EI} ** 2) / (1 + (z/\text{beta}) ** 2 * \text{sqrt}(\text{var})) ** 2$. This function calculates the standard EI, squares it, and then divides it by a penalty term that increases with both the standardized improvement $z = (\text{incumbent} - \text{mean}) / \text{std}$ and the uncertainty $\text{sqrt}(\text{var})$. The squaring of the EI nonlinearly amplifies regions with high EI values relative to

¹⁴We use precomputed loss values across 50 datasets given as part of the HyLAP project. For SVM, the two hyperparameters are the RBF kernel parameter and the penalty parameter while for AdaBoost they correspond to the number of product terms and the number of iterations.

those with low or moderate EI values. This increases the "peakiness" of the AF, leading to stronger exploitation of the most promising point(s). The term $(1 + (z/\text{beta}) ** 2 * \text{sqrt}(\text{var})) ** 2$ penalizes points more heavily if they have a very high EI (z is large and positive) and/or high uncertainty ($\text{sqrt}(\text{var})$ is large). This might act as a regularizer against over-optimism or excessive jumps into highly uncertain areas, even if they look very promising according to EI. It's a novel way of balancing exploration and exploitation discovered by FunBO.

FEW-SHOT. We conclude our experimental analysis by demonstrating how FunBO can be used in the context of few-shot adaptation. In this setting, we aim to learn an AF for a specific function class by "adapting" an initial AF with a small number of instances from the target class.

We consider Ackley ($d = 2$) as the objective function and compare against FSAF (Hsieh et al., 2021), the closest few-shot adaptation method for BO. FSAF trains the initial AF with a set of GPs, adapts it using 5 instances of scaled and translated Ackley functions, then tests the adapted AF on 100 additional Ackley instances, generated in the same manner. Note that FSAF uses a wide variety of GP functions with different kernels and various hyperparameters for training the initial AF. On the contrary, FunBO few-shot adaptation is performed by setting the initial h function to the one found in GPs-ID (Fig. 5, right plot, green line) using 25 GPs with RBF kernel, and including the 5 instances of Ackley used by FSAF in \mathcal{G}_{Tr} . Despite the limited training set, FunBO adapts very quickly to the new function instances, identifying an AF (code in Fig. 19) that outperforms both general purpose AFs and FSAF (Fig. 5, right).

5. Conclusions and Discussion

We tackled the problem of discovering novel, well-performing AFs for BO through FunBO, a FunSearch-based algorithm that explores the space of AFs by letting an LLM iteratively modify the AF expression in native computer code to improve the efficiency of the corresponding BO algorithm.

We have shown across a variety of settings that FunBO learns AFs that generalize well within and across function classes while being easily adaptable to specific objective functions of interest with only a few training examples.

Limitations. FunBO inherits the strengths of FunSearch along with some of its inherent constraints. While FunSearch allows finding programs that are concise and interpretable, it works best for programs that can be quickly evaluated and for which the score provides an accurate quantification of the improvement achieved. Therefore, a potential limitation of FunBO is the computational overhead associated with running a full BO loop for each candidate AF on every function in the set \mathcal{G} . This limits the scalability of FunBO for larger sets \mathcal{G} and hinders its application to more complex optimization problems, such as those with multiple objectives. However, several practical scenarios exist where the functions in \mathcal{G} are relatively inexpensive to evaluate, thus lowering the evaluation cost incurred during FunBO’s discovery phase. For instance, \mathcal{G} could consist of: (i) faster, lower-fidelity simulators or simplified analytical models related to the expensive target problem f ; (ii) cheap-to-evaluate surrogate models learned from previous related tasks; or (iii) hyperparameter optimization losses on smaller datasets where training and evaluation are faster than on the final, large target dataset. It is important to note that the computational cost incurred during this one-time search phase for a novel AF can be offset by its subsequent performance during online deployment, where the discovered AF may significantly reduce the number of costly function evaluations needed on the actual target problems.

Furthermore, the simple metric considered in this work in Eq. (1), only captures the distance from the true optimum and the number of trials needed to identify it. More research needs to be done to understand if a metric that better characterizes the convergence path for a given AF can improve FunBO’s performance. In addition, each FunBO experiment shown in this work required obtaining a large number of LLM samples. This means that the overall cost of experiments, which depends on the LLM used as well as the algorithm’s implementation (e.g. single-threaded or distributed, as originally proposed by FunSearch), can be high. Finally, as reported by Romera-Paredes et al. (2023), the variance in the quality of the AF found by FunBO is high. This is due to the randomness in both the LLM sampling and the evolutionary procedure. While we were able to reproduce the results shown for ID-Bench, HPO-ID and GPS-ID with different FunBO experiments, finding AFs that perform well across function classes required multiple FunBO runs.

Future work. This work opens up several promising avenues for future research. While our focus here was on the simplest single-output BO algorithm with a GP surrogate model, FunBO can be extended to learn AFs for various

adaptations of this problem, e.g., constrained optimization, noisy evaluations, multi-fidelity settings or alternative surrogate models. In addition, FunBO can be used to search in the space of functions with different inputs thus potentially discovering, e.g., non-myopic AFs. Our method is inherently flexible and can accommodate such extensions, which we view as natural follow-up work. Additionally, FunBO demonstrates the potential to harness the power of LLMs while maintaining the interpretability of AFs expressed in code. This opens an exciting avenue for exploring how and what assumptions can be encoded within AFs, based on the desired program characteristics and prior knowledge about the objective function. Finally, the discovered AFs might have intrinsic value, independently of how they were discovered. Future work could focus on more extensively testing their properties and identify those that can be added to the standard suite of AFs available in BO packages.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Benjamins, C., Raponi, E., Jankovic, A., Doerr, C., and Lindauer, M. Self-adjusting weighted expected improvement for Bayesian optimization. In *International Conference on Automated Machine Learning*, 2023.
- Bergstra, J., Bardenet, R., Bengio, Y., and Kégl, B. Algorithms for hyper-parameter optimization. In *Advances in Neural Information Processing Systems*, 2011.
- Calandra, R., Seyfarth, A., Peters, J., and Deisenroth, M. P. Bayesian optimization for learning gaits under uncertainty: An experimental comparison on a dynamic bipedal walker. *Annals of Mathematics and Artificial Intelligence*, 76:5–23, 2016.
- Chen, A., Dohan, D., and So, D. Evoprompting: Language models for code-level neural architecture search. In *Advances in Neural Information Processing Systems*, 2024a.
- Chen, L., Chen, J., Goldstein, T., Huang, H., and Zhou, T. InstructZero: Efficient instruction optimization for black-box large language models. In *International Conference on Machine Learning*, pp. 6503–6518, 2024b.
- Chen, Y., Hoffman, M. W., Colmenarejo, S. G., Denil, M., Lillicrap, T. P., Botvinick, M., and Freitas, N. Learning to learn without gradient descent by gradient descent. In *International Conference on Machine Learning*, pp. 748–756, 2017.

- Chen, Y., Song, X., Lee, C., Wang, Z., Zhang, R., Dohan, D., Kawakami, K., Kochanski, G., Doucet, A., Ranzato, M., Perel, S., and de Freitas, N. Towards learning universal hyperparameter optimizers with transformers. In *Advances in Neural Information Processing Systems*, pp. 32053–32068, 2022.
- Cheng, J., Liu, X., Zheng, K., Ke, P., Wang, H., Dong, Y., Tang, J., and Huang, M. Black-box prompt optimization: Aligning large language models without model training. In *Annual Meeting of the Association for Computational Linguistics*, 2024.
- Cho, H., Kim, Y., Lee, E., Choi, D., Lee, Y., and Rhee, W. Basic enhancement strategies when using Bayesian optimization for hyperparameter tuning of deep neural networks. *IEEE Access*, 8:52588–52608, 2020.
- Coffman, E. G., Garey, M. R., and Johnson, D. S. Approximation algorithms for bin-packing — an updated survey. In Ausiello, G., Lucertini, M., and Serafini, P. (eds.), *Algorithm Design for Computer System Design*, pp. 49–106. Springer, 1984.
- Ding, Q., Kang, Y., Liu, Y.-W., Lee, T. C. M., Hsieh, C.-J., and Sharpnack, J. Syndicated bandits: A framework for auto tuning hyper-parameters in contextual bandit algorithms. In *Advances in Neural Information Processing Systems*, pp. 1170–1181, 2022.
- Fernando, C., Banarse, D. S., Michalewski, H., Osindero, S., and Rocktäschel, T. Promptbreeder: Self-referential self-improvement via prompt evolution. In *International Conference on Machine Learning*, pp. 13481–13544, 2024.
- Feurer, M., Letham, B., and Bakshy, E. Scalable meta-learning for Bayesian optimization using ranking-weighted Gaussian process ensembles. In *AutoML Workshop at International Conference on Machine Learning*, 2018.
- Frazier, P. I., Powell, W. B., and Dayanik, S. A knowledge-gradient policy for sequential information collection. *SIAM Journal on Control and Optimization*, 47(5):2410–2439, 2008.
- Garnett, R. *Bayesian Optimization*. Cambridge University Press, 2023.
- Google-PaLM-2-Team. PaLM 2 Technical Report. *arXiv preprint arXiv:2305.10403*, 2023.
- Guo, Q., Wang, R., Guo, J., Li, B., Song, K., Tan, X., Liu, G., Bian, J., and Yang, Y. EvoPrompt: Connecting LLMs with evolutionary algorithms yields powerful prompt optimizers. In *International Conference on Learning Representations*, 2024.
- Hoffman, M., Brochu, E., De Freitas, N., et al. Portfolio allocation for Bayesian optimization. In *Conference on Uncertainty in Artificial Intelligence*, pp. 327–336, 2011.
- Hsieh, B.-J., Hsieh, P.-C., and Liu, X. Reinforced few-shot acquisition function learning for Bayesian optimization. In *Advances in Neural Information Processing Systems*, pp. 7718–7731, 2021.
- Jiang, A. Q., Li, W., Tworkowski, S., Czechowski, K., Odrzygóźdź, T., Miłoś, P., Wu, Y., and Jamnik, M. Thor: Wielding hammers to integrate language models and automated theorem provers. In *Advances in Neural Information Processing Systems*, pp. 8360–8373, 2022.
- Jones, D. R., Schonlau, M., and Welch, W. J. Efficient global optimization of expensive black-box functions. *Journal of Global Optimization*, 13:455–492, 1998.
- Korovina, K., Xu, S., Kandasamy, K., Neiswanger, W., Poczoz, B., Schneider, J., and Xing, E. Chembo: Bayesian optimization of small organic molecules with synthesizable recommendations. In *International Conference on Artificial Intelligence and Statistics*, pp. 3393–3403, 2020.
- Koza, J. R. Genetic programming as a means for programming computers by natural selection. *Statistics and computing*, 4:87–112, 1994.
- Kristiadi, A., Strieth-Kalthoff, F., Skreta, M., Poupart, P., Aspuru-Guzik, A., and Pleiss, G. A sober look at LLMs for material discovery: Are they actually good for Bayesian optimization over molecules? In *International Conference on Machine Learning*, pp. 25603–25622, 2024.
- Kushner, H. J. A versatile stochastic model of a function of unknown and time varying form. *Journal of Mathematical Analysis and Applications*, 5(1):150–167, 1962.
- Kushner, H. J. A new method of locating the maximum point of an arbitrary multipeak curve in the presence of noise. *Journal Basic Engineering*, 86(1):97–106, 1964.
- Lai, T. L. and Robbins, H. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6(1): 4–22, 1985.
- Lehman, J., Gordon, J., Jain, S., Ndousse, K., Yeh, C., and Stanley, K. O. Evolution through large models. In Banzhaf, W., Machado, P., and Zhang, M. (eds.), *Handbook of evolutionary machine learning*, pp. 331–366. Springer, 2023.
- Li, L., Jamieson, K., DeSalvo, G., Rostamizadeh, A., and Talwalkar, A. Hyperband: A novel bandit-based approach to hyperparameter optimization. *Journal of Machine Learning Research*, 18(185):1–52, 2018.

- Liu, F., Xialiang, T., Yuan, M., Lin, X., Luo, F., Wang, Z., Lu, Z., and Zhang, Q. Evolution of heuristics: Towards efficient automatic algorithm design using large language model. In *International Conference on Machine Learning*, pp. 32201–32223, 2024a.
- Liu, T., Astorga, N., Seedat, N., and van der Schaar, M. Large language models to enhance Bayesian optimization. In *International Conference on Learning Representations*, 2024b.
- Lyu, W., Yang, F., Yan, C., Zhou, D., and Zeng, X. Batch Bayesian optimization via multi-objective acquisition ensemble for automated analog circuit design. In *International Conference on Machine Learning*, pp. 3306–3314, 2018.
- Maraval, A., Zimmer, M., Grosnit, A., and Bou Ammar, H. End-to-end meta-Bayesian optimisation with transformer neural processes. In *Advances in Neural Information Processing Systems*, 2024.
- Meyerson, E., Nelson, M. J., Bradley, H., Gaier, A., Moradi, A., Hoover, A. K., and Lehman, J. Language model crossover: Variation through few-shot prompting. *ACM Transactions on Evolutionary Learning and Optimization*, 4(4), 2024.
- Mockus, J. On Bayesian methods for seeking the extremum. *Proceedings of the IFIP Technical Conference*, pp. 400–404, 1974.
- Müller, S., Feurer, M., Hollmann, N., and Hutter, F. In-context learning for bayesian optimization. In *International Conference on Machine Learning*, pp. 25444–25470, 2023.
- Nasir, M. U., Earle, S., Togelius, J., James, S., and Cleghorn, C. Llmatic: Neural architecture search via large language models and quality diversity optimization. In *Genetic and Evolutionary Computation Conference*, pp. 1110–1118, 2024.
- Perrone, V., Shen, H., Seeger, M. W., Archambeau, C., and Jenatton, R. Learning search spaces for Bayesian optimization: Another view of hyperparameter transfer learning. In *Advances in neural information processing systems*, 2019.
- Polu, S. and Sutskever, I. Generative language modeling for automated theorem proving. *arXiv preprint arXiv:2009.03393*, 2020.
- Ramos, M. C., Michtavy, S. S., Porosoff, M. D., and White, A. D. Bayesian optimization of catalysts with in-context learning. *arXiv preprint arXiv:2304.05341*, 2023.
- Ranković, B. and Schwaller, P. Bochemian: Large language model embeddings for Bayesian optimization of chemical reactions. In *Adaptive Experimental Design and Active Learning in the Real World Workshop at Advances in Neural Information Processing Systems*, 2023.
- Rasmussen, C. E. and Williams, C. K. *Gaussian Processes for Machine Learning*. MIT Press Cambridge, MA, 2006.
- Romera-Paredes, B., Barekatin, M., Novikov, A., Balog, M., Kumar, M. P., Dupont, E., Ruiz, F. J., Ellenberg, J. S., Wang, P., Fawzi, O., et al. Mathematical discoveries from program search with large language models. *Nature*, pp. 1–3, 2023.
- Šaltanis, R. P. One method of multiextremum optimization. *Avtomatika i Vychislitel'naya Tekhnika (Automatic Control and Computer Sciences)*, 5(3):33–38, 1971.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Snoek, J., Larochelle, H., and Adams, R. P. Practical Bayesian optimization of machine learning algorithms. In *Advances in Neural Information Processing Systems*, 2012.
- Sobol', I. M. On the distribution of points in a cube and the approximate evaluation of integrals. *Zhurnal Vychislitel'noi Matematiki i Matematicheskoi Fiziki*, 7, 1967.
- Song, L., Gao, C., Xue, K., Wu, C., Li, D., Hao, J., Zhang, Z., and Qian, C. Reinforced in-context black-box optimization. *arXiv preprint arXiv:2402.17423*, 2024a.
- Song, X., Tian, Y., Lange, R. T., Lee, C., Tang, Y., and Chen, Y. Position: leverage foundational models for black-box optimization. In *International Conference on Machine Learning*, 2024b.
- Sun, T., Shao, Y., Qian, H., Huang, X., and Qiu, X. Black-box tuning for language-model-as-a-service. In *International Conference on Machine Learning*, pp. 20841–20855, 2022.
- Swersky, K., Snoek, J., and Adams, R. P. Multi-task Bayesian optimization. In *Advances in Neural Information Processing Systems*, 2013.
- Tanese, R. *Distributed Genetic Algorithms for Function Optimization*. PhD thesis, University of Michigan, 1989.
- Tribes, C., Benarroch-Lelong, S., Lu, P., and Kobzyev, I. Hyperparameter optimization for large language model instruction-tuning. In *AAAI Conference on Artificial Intelligence*, 2024.

- Vertex AI, G. C. Code models overview. 2023. URL <https://cloud.google.com/vertex-ai/docs/generative-ai/code/code-models-overview>.
- Volpp, M., Fröhlich, L. P., Fischer, K., Doerr, A., Falkner, S., Hutter, F., and Daniel, C. Meta-learning acquisition functions for transfer learning in Bayesian optimization. In *International Conference on Learning Representations*, 2020.
- Wang, C., Liu, X., and Awadallah, A. H. Cost-effective hyperparameter optimization for large language model generation inference. In *International Conference on Automated Machine Learning*, pp. 21–1, 2023.
- Wang, Z. and Jegelka, S. Max-value entropy search for efficient Bayesian optimization. In *International Conference on Machine Learning*, pp. 3627–3635, 2017.
- Wistuba, M. and Grabocka, J. Few-shot Bayesian optimization with deep kernel surrogates. In *International Conference on Learning Representations*, 2021.
- Wistuba, M., Schilling, N., and Schmidt-Thieme, L. Scalable Gaussian process-based transfer surrogates for hyperparameter optimization. *Machine Learning*, 107(1): 43–78, 2018.
- Yang, C., Wang, X., Lu, Y., Liu, H., Le, Q. V., Zhou, D., and Chen, X. Large language models as optimizers. In *International Conference on Learning Representations*, 2024.
- Yao, Y., Liu, F., Cheng, J., and Zhang, Q. Evolve cost-aware acquisition functions using large language models. In Affenzeller, M., Winkler, S. M., Kononova, A. V., Trautmann, H., Tušar, T., Machado, P., and Bäck, T. (eds.), *Parallel Problem Solving from Nature – PPSN XVIII*, pp. 374–390. Springer Nature Switzerland, 2024.
- Yogatama, D. and Mann, G. Efficient transfer learning method for automatic hyperparameter tuning. In *International Conference on Artificial Intelligence and Statistics*, pp. 1077–1085, 2014.
- Zhang, M. R., Desai, N., Bae, J., Lorraine, J., and Ba, J. Using large language models for hyperparameter optimization. In *Foundation Models for Decision Making Workshop at Advances in Neural Information Processing Systems*, 2023.
- Zheng, M., Su, X., You, S., Wang, F., Qian, C., Xu, C., and Albanie, S. Can GPT-4 perform neural architecture search? *arXiv preprint arXiv:2304.10970*, 2023.

A. Related work

LLMs as mutation operators. FunBO expands FunSearch (Romera-Paredes et al., 2023), an evolutionary algorithm pairing an LLM with an evaluator to solve open problems in mathematics and algorithm design. The idea of using LLMs as mutation operators paired with a scoring mechanism has been explored to create a self-improvement loop (Lehman et al., 2023), to optimize code for robotic simulations, to optimize prompts (Guo et al., 2024), or to evolve stable diffusion images with simple genetic algorithms (Meyerson et al., 2024). Other works explore the use of LLMs to search over neural network architectures described with Python code (Nasir et al., 2024; Zheng et al., 2023; Chen et al., 2024a), to find formal proofs for automatic theorem proving (Polu & Sutskever, 2020; Jiang et al., 2022) or to automatically design heuristics (Liu et al., 2024a). Symbolic regression, similarly to FunBO, also searches for mathematical expressions, which are often represented as trees. However, while FunBO employs an LLM for both generation and mutation, symbolic regression typically evolves these expressions using genetic operators (Koza, 1994). In addition, FunBO operates directly on code representations and uses full BO performance as its fitness metric. Therefore, while both methodologies aim to find functional forms, their search space representations and primary evolutionary operators are different.

Meta-learning for BO. Our work is also related to the literature on meta-learning for BO. In this realm, several studies have focused on meta-learning an accurate surrogate model for the objective function by exploiting observations from related functions, for instance by using standard multi-task GPs (Swersky et al., 2013; Yogatama & Mann, 2014), ensembles of GP models (Feurer et al., 2018; Wistuba et al., 2018; Wistuba & Grabocka, 2021) or neural processes that are trained to approximate the posterior predictive distribution through in-context learning on any prior distribution that can be efficiently sampled from (Müller et al., 2023). Others have focused on meta-learning general-purpose optimizers by using recurrent neural networks with access to gradient information (Chen et al., 2017) or transformers (Chen et al., 2022; Song et al., 2024a). Note that, while meta-learned surrogate models *explicitly* learn structure from past functions by observing data-points for each of them, methods that meta-learn AFs via \mathcal{G} *implicitly* learn similarities between these objectives by observing the optimization pattern achieved by previously sampled AFs for each objective function in \mathcal{G} . Interestingly, the most significant performance gains observed for the approach proposed by Chen et al. (2022) stem from using a standard AF (EI) on top of the transformer architecture for output predictions. This confirms the continued importance of AFs as crucial components in BO, even when combined with transformer-based approaches, and highlights the importance of a method such as FunBO that can be seamlessly integrated with these newer architectures, potentially leading to further improvements in performance. More relevant to our work are studies focusing on transferring information from related tasks by learning novel AFs that more efficiently solve the classic exploration-exploitation trade-off in BO algorithms (Volpp et al., 2020; Hsieh et al., 2021; Maraval et al., 2024). In contrast to prior works in this literature, FunBO produces AFs that are more interpretable, simpler and cheaper to deploy than neural network-based AFs and generalize not only within specific function classes but also across different classes.

LLMs and black-box optimization. Several works have investigated the use of LLMs to solve black-box optimization problems. Song et al. (2024b) discussed different ways in which foundational language models can revolutionize optimization, from harnessing the vast wealth of information encapsulated in free-form text to enriching task comprehension. Both Liu et al. (2024b) and Yang et al. (2024) framed optimization problems in natural language and asked LLMs to iteratively propose promising solutions and/or evaluate them. Similarly, Ramos et al. (2023) replaced surrogate modeling with LLMs within a BO algorithm targeted at catalyst or molecule optimization. Kristiadi et al. (2024) also investigated whether LLMs can accelerate BO in the molecular space while Ranković & Schwaller (2023) explored the integration of LLMs with BO in the domain of chemical reaction optimization. Other works have focused on exploiting black-box methods for prompt optimization (Sun et al., 2022; Chen et al., 2024b; Cheng et al., 2024; Fernando et al., 2024), solving HPO tasks with LLMs (Zhang et al., 2023) or identifying optimal LLM hyperparameter settings via black-box optimization approaches (Wang et al., 2023; Tribes et al., 2024). Concurrent to our work, Yao et al. (2024) propose using an LLM coupled with an evolutionary procedure to find cost-aware AFs.

AFs representations Works proposing new meta-learned or general-purpose AFs can also be classified based on the representation used for the AFs. Unlike general-purpose AFs, for which an analytical representation is available, recent works have explored representing AFs via neural networks or code. Among the works using neural networks, Volpp et al. (2020) proposed a *neural* AF that is a MLP with ReLU activations while Chen et al. (2022) and Maraval et al. (2024) jointly trained surrogate models and AFs via transformers or neural processes. Instead, the recent work by Yao et al. (2024) represents AFs for settings with limited experimentation budgets in code.

```

import numpy as np
from scipy.stats import norm

def acquisition_function(predictive_mean, predictive_var, incumbent, beta=1.0):
    """Returns the index of the point to collect in a vector of eval points.

    Given the posterior mean and posterior variance of a GP model for the objective function, this function computes
    a heuristic and finds its optimum. The next function evaluation will be placed at the point corresponding to the
    selected index in a vector of possible eval points.

    Args:
        predictive_mean: an array of shape [num_points, dim] containing the predicted mean values for the GP model on
            the objective function for 'num_points' points of dimensionality 'dim'.
        predictive_var: an array of shape [num_points, dim] containing the predicted variance values for the GP model on
            the objective function for 'num_points' points of dimensionality 'dim'.
        incumbent: current optimum value of objective function observed.
        beta: a possible hyperparameter to construct the heuristic.

    Returns:
        An integer representing the index of the point in the array of shape [num_points, dim]
        that needs to be selected for function evaluation.
    """
    z = (incumbent - predictive_mean) / np.sqrt(predictive_var)
    predictive_std = np.sqrt(predictive_var)
    vals = (incumbent - predictive_mean) * norm.cdf(z) + predictive_std * norm.pdf(z)
    return np.argmax(vals)

```

Figure 6. Python code for FunBO’s initial h function with full docstring.

B. Code for FunBO components

Fig. 6 gives the Python code for the initial acquisition function used by FunBO, including the full docstring. The docstring describes the inputs of the function and the way in which the function itself is used within the evaluate function e . Evaluation of the functions generated by FunBO is done by first running a full BO loop (see Fig. 7 for Python code) and then, based on its output (the initial optimal input value, the true optimum, the found optimum and the percentage of steps taken before finding the latter), computing the score as in the Python code of Fig. 8. Note how the latter captures how accurately and quickly a BO algorithm using the proposed AF finds the true optimum.

C. Programs Database

The DB structure matches the one proposed by FunSearch (Romera-Paredes et al., 2023). We discuss it here for completeness. A multiple-deme model (Tanese, 1989) is employed to preserve and encourage diversity in the generated programs. Specifically, the program population in DB is divided into N_{DB} islands, each initialized with the given initial h and evolved independently. Within each island, programs are clustered based on their scores on the functions in \mathcal{G}_{Tr} , with AFs having the same scores grouped together. Sampling from DB involves first uniformly selecting an island and then sampling two AFs from it. Within the chosen island, a cluster is sampled, favoring those with higher scores, followed by sampling a program within that cluster, favoring shorter ones. The newly generated AF is added to the same island associated with the instances in the prompt, but to a cluster reflecting its scores on \mathcal{G}_{Tr} . Every 4 hours, all programs from the $N_{DB}/2$ islands with the lowest-scoring best AF are discarded. These islands are then reseeded with a single program from the surviving islands. This procedure eliminates underperforming AFs, creating space for more promising programs. See the Methods section in Romera-Paredes et al. (2023) for further details.

```

"""Evaluate an AF with a full BO loop for the objective f."""
import GPy
import numpy as np
import utils

def run_bo(f, # objective function to minimize
          acquisition_function, # h given by LLM
          num_eval_points = 1000, num_trials = 30):
    """Run a BO loop and return the minimum value found and the percentage of trials required to reach it."""

    # Get evaluation points for AF. get_eval_points() returns a given number of points on a
    # Sobol grid on the f's input space
    eval_points = utils.get_eval_points(f, num_eval_points)

    # Get the initial point with get_initial_design(). This is set to be the point giving the
    # maximum (worst) function evaluation among eval_points
    initial_x, initial_y = utils.get_initial_design(f)

    # Initialize GP hyper-parameters. We pre-compute the RBF kernel hyper-parameters
    # for each given f. These are returned by get_hyperparameters().
    hp = utils.get_hyperparameters(f)

    # Initialize kernel and model.
    kernel = GPy.kern.RBF(input_dim=initial_x.shape[1], variance=hp['variance'],
                          lengthscale=hp['lengthscale'], ARD=hp['ard'])
    model = GPy.models.GPRegression(initial_x, initial_y, kernel)

    # Get initial predictive mean and var.
    predictive_mean, predictive_var = model.predict(eval_points)

    # Get initial optimum value.
    found_min = initial_min_y = float(np.min(model.Y))

    # Get true optimum value.
    true_min = np.min(f(eval_points))

    # Optimization loop.
    for _ in range(num_trials):
        selected_idx = acquisition_function(predictive_mean, predictive_var, found_min) # Get index for new point using AF.
        new_input = eval_points[selected_idx, :] # Get new point.
        new_output = f(new_input) # Evaluate new point.
        # Append to dataset.
        model.set_XY(np.concatenate((model.X, new_input), axis=0), np.concatenate((model.Y, new_output), axis=0))
        # Get updated mean and var
        predictive_mean, predictive_var = model.predict(eval_points)
        found_min = float(np.min(model.Y)) # Get current optimum value.

    # Get percentage of trials (note that we remove the number of given points in the initial design) needed
    # to identify the optimum.
    percentage_steps_before_converging = (np.argmax(model.Y) - initial_x.shape[0]) / (
        num_trials) if found_min == true_min else 1.0
    return (found_min, true_min, initial_min_y, percentage_steps_before_converging)

```

Figure 7. Python code for the first part of e used in FunBO. This function runs a full BO loop with a given number of trials and points on a Sobol grid to assess how efficiently a given AF allows optimizing f .

```

"""Score an AF given the output of run_bo()."""

import numpy as np

def score(found_min, true_min, initial_min_y, percentage_steps_before_converging):
    """Compute a score based on the output of run_bo()."""

    # Get score based on how close the found optimum is to the true one (first term
    # in Eq. (1)).
    score_min_reached = 1.0 - np.abs(found_min - true_min) / (initial_min_y - true_min)

    # Get score based on how the percentage of trials needed to identify the true
    # optimum (second term in Eq. (1)).
    score_steps_needed = 1.0 - percentage_steps_before_converging

    return score_min_reached + score_steps_needed

```

Figure 8. Python code for the second part of e used in FunBO. Based on the output of `run_bo()`, this function computes a score capturing how accurately and quickly an AF allows identifying the true optimum.

D. Experimental details

In this section, we provide the experimental details for all our experiments. We run FunBO with $\mathcal{T} = 48\text{hrs}$, $B = 12$ and $N_{\text{DB}} = 10$. To isolate the effect of using different AFs and eliminate confounding factors related to AF maximization or surrogate models, we maximized all AFs on a fixed Sobol grid (of size N_{SG}) over each function’s input space. We also ensured the same initial design across all methods (including the point with the highest/worst function value on the Sobol grid) and used consistent GP hyperparameters which were tuned offline and fixed. In particular, we use a GP model with zero mean function and RBF kernel defined as $K_{\theta}(\mathbf{X}, \mathbf{X}') = \sigma_f^2 \exp(-\|\mathbf{X} - \mathbf{X}'\|^2 / 2\ell^2)$ with $\theta = (\ell, \sigma_f^2)$ where ℓ and σ_f^2 are the length-scale and kernel variance respectively. The Gaussian likelihood noise σ^2 is set to $1e - 5$ unless otherwise stated. We set $T = 30$ for all experiments apart from HPO-ID and GPs-ID for which we use $T = 20$ to ensure faster evaluations of generated AFs. We used the MetaBO implementation provided by the authors at github.com/boschresearch/MetaBO, retaining default parameters except for removing the local maximization of AFs and ensuring consistency in the initial design. We followed the same procedure for FSAF, using code available at github.com/pinghsieh/FSAF. We ran UCB with $\beta = 1$. Experiment-specific settings are detailed below.

D.1. OOD-Bench

The parameter configurations adopted for each objective function used in this experiment, either in \mathcal{G} or in \mathcal{F} , are given in Table 1. Notice that for Hartmann with $d = 3$ we use an ARD kernel. Scaled and translated functions are obtained with translations sampled uniformly in $[-0.1, 0.1]^d$ and scalings sampled uniformly in $[0.9, 1.1]$. Fig. 9 gives the results achieved by α_{FunBO} (blue line) and a dimensionality-agnostic version of MetaBO that does not take the possible evaluation points as input to the neural AF. This allows the neural AF to be trained on one-dimensional functions and be used to optimize functions across input dimensions.

D.2. ID-Bench

The parameter configurations for Branin, Goldstein-Price and Hartmann are given in Table 2. For this experiment, we adopt the parameters used by Volpp et al. (2020) thus optimize the functions in the unit-hypercube and use ARD RBF kernels. Fig. 10 gives the results achieved by α_{FunBO} (blue line) and the AF found by FunBO for OOD-Bench (green). The Python code for the found AFs is given in Figs. 11-13.

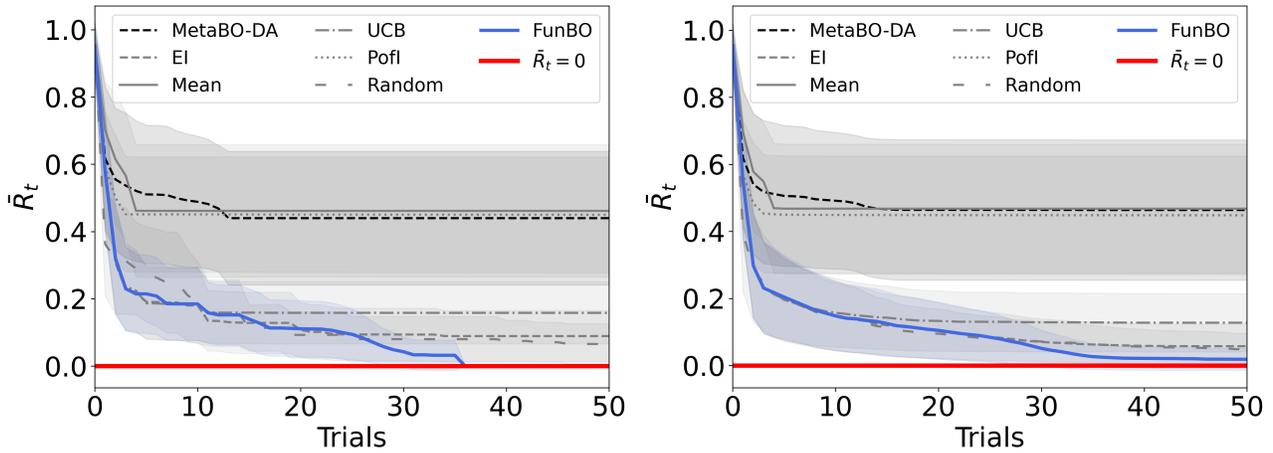


Figure 9. OOD-Bench. Average BO performance when using known general-purpose AFs (gray lines with different patterns), the AF learned by a dimensionality-agnostic version of MetaBO (MetaBO-DA, black dashed line) and α_{FunBO} (blue line). Shaded area gives \pm std. dev./2. The red line represents $\bar{R}_t = 0$, i.e., zero average regret. *Left*: \mathcal{F} includes nine different synthetic functions. *Right*: Extended test set including, for each function in \mathcal{F} , 50 randomly scaled and translated instances.

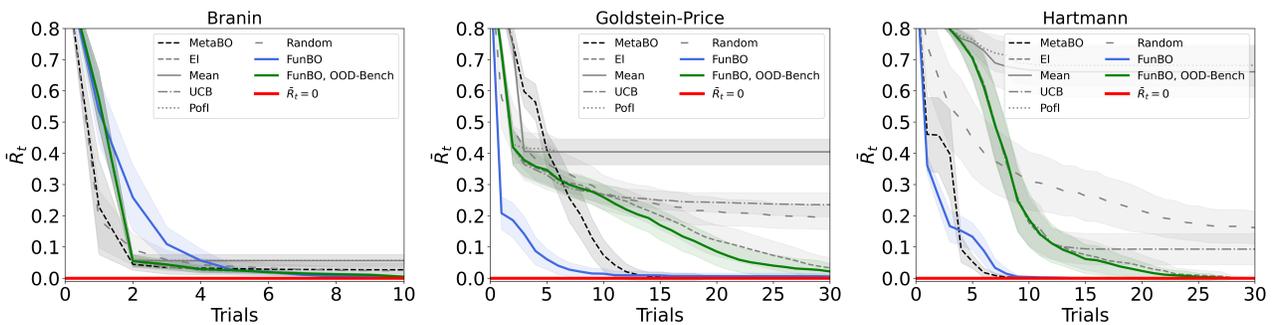


Figure 10. ID-Bench. Average BO performance when using known general-purpose AFs (gray lines with different patterns), α_{FunBO} found in OOD-Bench (green line), the AF learned by MetaBO (black dashed line) and α_{FunBO} (blue line) on 100 instances of Branin, Goldstein-Price and Hartmann. Shaded area gives \pm std. dev./2. The red line represents $\bar{R}_t = 0$, i.e., zero average regret.

```

import numpy as np
from scipy.stats import norm

def acquisition_function(predictive_mean, predictive_var, incumbent, beta=1.0):
    """Returns the index of the point to collect ..."""
    y_pred = predictive_mean + 2 * predictive_var
    diff_current_best_y_pred = incumbent - y_pred
    bound_standard_deviation = np.maximum(np.sqrt(predictive_var), 1e-15)
    z = diff_current_best_y_pred / bound_standard_deviation
    vals = (diff_current_best_y_pred * norm.cdf(z)
            + np.sqrt(predictive_var) * norm.cdf(z + 0.5)
            + (norm.cdf(z) - norm.cdf(z + 0.5)) * predictive_var / 2)
    a = np.maximum(diff_current_best_y_pred, incumbent)
    alpha = diff_current_best_y_pred if incumbent > 0.0 else -np.inf
    alpha = np.maximum(alpha, 0.) * (-alpha + 0.5 * a) - y_pred
    y_vals = np.absolute(alpha + a + np.abs(y_pred)) * (a >= 0.)
    for y_val in y_vals:
        idx = np.argmax(vals - (y_val - y_pred) / bound_standard_deviation)
        vals[idx] = 0
    return np.argmax(vals)

```

Figure 11. ID-Bench. Python code for α_{FunBO} for Branin. The BO performance corresponding to this AF is given in Fig. 4 (left).

```

import numpy as np
from scipy.stats import norm

def acquisition_function(predictive_mean, predictive_var, incumbent, beta=1.0):
    """Returns the index of the point to collect ..."""
    shape, dim = predictive_mean.shape
    best_score = 0.0
    g_i = 0.0

    predictive_var[(shape-10)//2] *= dim
    predictive_var[~ np.isfinite(predictive_var)] = 1.0

    for i in range(predictive_mean.shape[0]):
        curr_z = (incumbent - predictive_mean[i]) / np.sqrt(predictive_var[i])
        new_score = predictive_var[i] * norm.cdf(curr_z, 0.5)
        if new_score > best_score:
            best_score = new_score
            g_i = i
    return g_i

```

Figure 12. ID-Bench. Python code for α_{FunBO} for Goldstein-Price. The BO performance corresponding to this AF is given in Fig. 4 (middle).

Table 1. Parameters used for OOD-Bench.

	d	\mathcal{X}	N_{SG}	ℓ	σ_f^2	σ^2
Ackley	1	$[-4, 4]$	1000	0.21	28.19	$1e-5$
Levy	1	$[-10, 10]$	1000	1.05	83.32	$1e-5$
Schwefel	1	$[-500, 500]$	1000	18.46	76868.65	$1e-5$
Rosenbrock	1	$[-5, 10]$	1000	1.20	87328.20	$1e-5$
Sphere	1	$[-5, 5]$	1000	18.46	924202.43	$1e-5$
Styblinski-Tang	1	$[-5, 5]$	1000	7.34	119522207.86	$1e-5$
Weierstrass	1	$[-0.5, 0.5]$	1000	0.01	0.39	$1e-5$
Beale	2	$[-4, 5]^2$	10000	0.46	546837.32	$1e-5$
Branin	2	$[-5, 10] \times [0, 15]$	10000	4.65	155233.52	$1e-5$
Michalewicz	2	$[0, \pi]^2$	10000	0.22	0.10	$1e-5$
Goldstein-Price	2	$[-2, 2]^2$	10000	0.27	117903.96	$1e-5$
Hartmann-3	3	$[0, 1]^3$	1728	$[0.716, 0.298, 0.186]$	0.83	$1.688e-11$
Hartmann-6	6	$[0, 1]^6$	729	1.0	1.0	$1e-5$

Table 2. Parameters used for ID-Bench.

	d	\mathcal{X}	N_{SG}	ℓ	σ_f^2	σ^2
Branin	2	$[0, 1]^2$	961	$[0.235, 0.578]$	2.0	$8.9e-16$
Goldstein-Price	2	$[0, 1]^2$	961	$[0.130, 0.07]$	0.616	$1e-6$
Hartmann-3	3	$[0, 1]^3$	1728	$[0.716, 0.298, 0.186]$	0.83	$1.688e-11$

D.3. HPO-ID

For this experiment, we adopt the GP hyperparameters used by Volpp et al. (2020). From the training datasets used in MetaBO, we assign “bands”, “wine”, “coil2000”, “winequality-red” and “titanic” for Adaboost, and “bands”, “breast-cancer”, “banana”, “yeast” and “vehicle” for SVM to \mathcal{G}_V . We keep the rest in \mathcal{G}_{Tr} . Fig. 14 gives the results achieved by α_{FunBO} (blue lines) and the AF found by FunBO for OOD-Bench (green lines). The Python code for the found AFs is given in Figs. 15-16.

D.4. GPS-ID

The functions included in \mathcal{G} and \mathcal{F} are sampled from a GP prior with RBF kernel and length-scale values drawn uniformly from $[0.05, 0.5]$. The functions are optimized in the input space $[0, 1]^3$ with $N_{SG} = 1728$ points. In terms of GP hyperparameters, we set $\sigma_f^2 = 1.0$, $\sigma^2 = 1e-20$ and use the length-scale value used to sample each function as ℓ . Fig. 17 gives the results achieved by α_{FunBO} and the AF found by FunBO for OOD-Bench. The Python code for α_{FunBO} is given in Fig. 18.

D.5. FEW-SHOT

For this experiment, the 5 Ackley functions used to “adapt” the initial AF are obtained by scaling and translating the output and inputs values with translations and scalings uniformly sampled in $[-0.1, 0.1]^d$ and $[0.9, 1.1]$ respectively. The test set includes 100 instances of Ackley similarly obtained with scale and translation values in $[0.7, 1.3]$ and $[-0.3, 0.3]^d$ respectively. Furthermore, we consider $[0, 1]^2$ as input space and use $N_{SB} = 1000$. The GP hyperparameters are set to $\ell = [0.07, 0.018]$ (ARD kernel), $\sigma_f^2 = 1.0$ and $\sigma^2 = 8.9e-16$. Python code for α_{FunBO} is given in Fig. 19.

D.6. OOD-Bench with BoTorch eval pipeline

To ensure a consistent comparison across all AFs, in the manuscript, we have presented results obtained by testing AFs using a fixed Sobol grid, fixed hyperparameters (tuned offline), and a fixed initial design that includes the input point yielding the maximum function value on the grid. As this setup might deviate from a fully realistic deployment scenario, in this section,

```

import numpy as np
from scipy.stats import norm
from scipy import stats

def acquisition_function(predictive_mean, predictive_var, incumbent, beta=1.0):
    """Returns the index of the point to collect ..."""
    diff_current_best_mean = incumbent - predictive_mean
    standard_deviation = np.sqrt(predictive_var)
    z = diff_current_best_mean / standard_deviation
    vals = diff_current_best_mean * norm.cdf(z)**3 + (
        norm.cdf(z)**2 + norm.cdf(z) + 1) * norm.pdf(z)
    index = np.argmax(stats.truncnorm.cdf(vals, a=-0.1, b=0.1))
    return index

```

Figure 13. ID-Bench. Python code for α_{FunBO} for Hartmann. The BO performance corresponding to this AF is given in Fig. 4 (right).

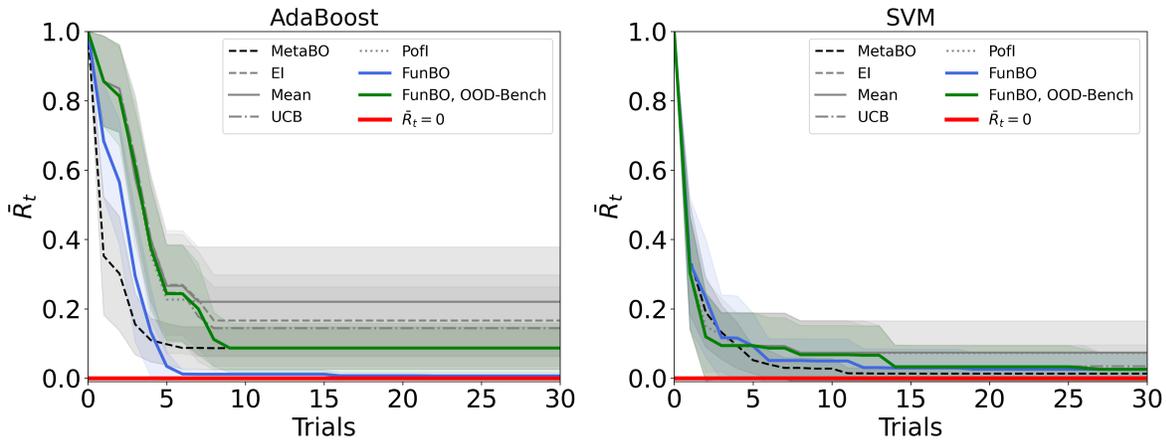


Figure 14. HPO-ID. Average BO performance when using known general-purpose AFS (gray lines with different patterns), a meta-learned AF by MetaBO (black dashed line), α_{FunBO} found in OOD-Bench (green lines) and α_{FunBO} (blue lines). Shaded area gives \pm std. dev./2. The red line represents $\bar{R}_t = 0$, i.e., zero average regret.

```

import numpy as np
from scipy.stats import norm

def acquisition_function(predictive_mean, predictive_var, incumbent, beta=1.0):
    """Returns the index of the point to collect ..."""
    c1 = np.exp(-beta)
    c2 = 2.0 * beta * np.exp(-beta)
    alpha = np.sqrt(2.0) * beta * np.sqrt(predictive_var)
    z = (incumbent - predictive_mean) / alpha
    vals = -abs(c1 * np.exp(- np.power(z, 2)) - 1.0 + c1 + incumbent
        ) + 2.0 * beta * np.power(z+c2, 2)
    vals -= np.log(np.power(alpha, 2))
    vals[np.argmin(vals)] = 1.0
    return np.argmin(vals)

```

Figure 15. HPO-ID. Python code for α_{FunBO} for AdaBoost. The BO performance corresponding to this AF is given in Fig. 5 (left).

```

import numpy as np
from scipy.stats import norm

def acquisition_function(predictive_mean, predictive_var, incumbent, beta=1.0):
    """Returns the index of the point to collect ..."""
    z = (incumbent - predictive_mean) / np.sqrt(predictive_var)
    vals = (incumbent - predictive_mean) * norm.cdf(z)
        + np.sqrt(predictive_var) * norm.pdf(z)
    t0_val = norm(loc=incumbent, scale=np.sqrt(predictive_var)).pdf(incumbent)
    t1_val = z * norm.pdf(z)
    vals = ((vals * t1_val - t0_val) / (1 - 2 * t1_val)
        + t1_val*(vals/(1-2*t1_val))
        - vals/(1 - 2*t1_val)**2 + t1_val*(t1_val - z)/beta)
    return np.argmax(vals)

```

Figure 16. HPO-ID. Python code for α_{FunBO} for SVM. The BO performance corresponding to this AF is given in Fig. 14 (right).

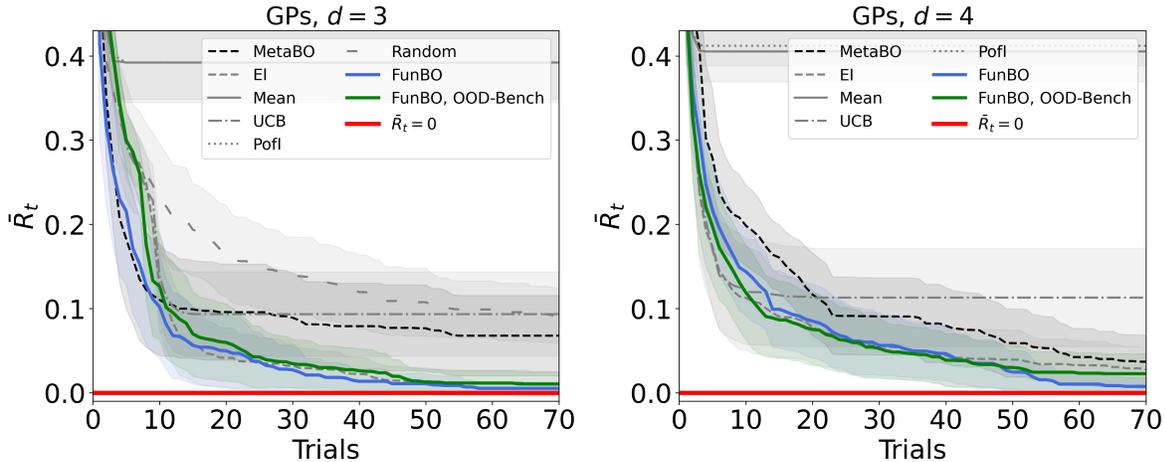


Figure 17. Average BO performance when using known general-purpose AFs (gray lines with different patterns), the AF learned by MetaBO (black dashed line), α_{FunBO} found in OOD-Bench (green lines) and α_{FunBO} (blue lines). Shaded area gives \pm std. dev./2. The red line represents $\bar{R}_t = 0$, i.e. zero average regret. *Left*: GPs-ID. \mathcal{F} includes functions with $d = 3$. *Right*: \mathcal{F} includes functions with $d = 4$.

```

import numpy as np
from scipy.stats import norm

def acquisition_function(predictive_mean, predictive_var, incumbent, beta = 1.0):
    """Returns the index of the point to collect ..."""
    z = (incumbent - predictive_mean) / np.sqrt(predictive_var)
    vals = ((incumbent - predictive_mean) * norm.cdf(z)
        + np.sqrt(predictive_var) * norm.pdf(z))**2
    vals = vals / (1 + (z / beta)**2 * np.sqrt(predictive_var))**2
    return np.argmax(vals)

```

Figure 18. GPs-ID. Python code for α_{FunBO} . The BO performance corresponding to this AF is given in Fig. 5 (right).

```

import numpy as np
from scipy.stats import norm

def acquisition_function(predictive_mean, predictive_var, incumbent, beta=1.0):
    """Returns the index of the point to collect ..."""
    num_points, _ = predictive_mean.shape
    a = 10
    z = (predictive_mean + 0.000001 - incumbent) / np.sqrt(predictive_var)
    vals = 1 / ((1 + (z / beta)**2 * np.sqrt(a * predictive_var + 0.00001)) **2)
    beta_sqrt_p_z = np.sqrt(beta) * z
    vals *= (1 + (z / beta)**2)*predictive_var/(
        (1+ (beta_sqrt_p_z / np.sqrt(predictive_var))**2 * predictive_var) * (
            1+(beta_sqrt_p_z / np.sqrt(predictive_var))**2))
    vals += (1 - beta_sqrt_p_z / np.sqrt(predictive_var))**2 * predictive_var / (
        1 + (beta_sqrt_p_z / np.sqrt(predictive_var))**2 * predictive_var)**2
    vals = (1 + (z / beta)**2) * vals - (1 - (z / beta)**2) * np.exp(- 1) ** 2
    vals = np.sqrt(a * predictive_var) * vals / np.sqrt(
        a * predictive_var + 0.00001)
    vals *= np.sqrt(np.sqrt(a * predictive_var) * predictive_var)
    vals *= predictive_var**2
    vals[:num_points // 2] = 0
    return np.argmax(vals)

```

Figure 19. FEW-SHOT. Python code for α_{FunBO} . The BO performance corresponding to this AF is given in Fig. 5 (right).

we repeat the evaluation of the AF found for OOD-Bench using the standard BoTorch evaluation pipeline, as demonstrated in the Colab notebooks [closed-loop tutorial](#) and [custom acquisition tutorial](#).

We adhere to the settings specified in these Colab notebooks, thus we optimize the AFs numerically at each trial, we fit the hyperparameters of the GP model at every trial using L-BFGS-B, and randomly select initial points. For each test function, we ran the algorithm with 10 different random initial designs and plot the results by averaging over the runs for each function in OOD-Bench. Note that we use the exact plot formatting from BoTorch shown in the linked Colabs (which differs from the other plots in the manuscript). Fig. 20 shows how, across all functions, FunBO performs either comparably or better than EI and UCB.

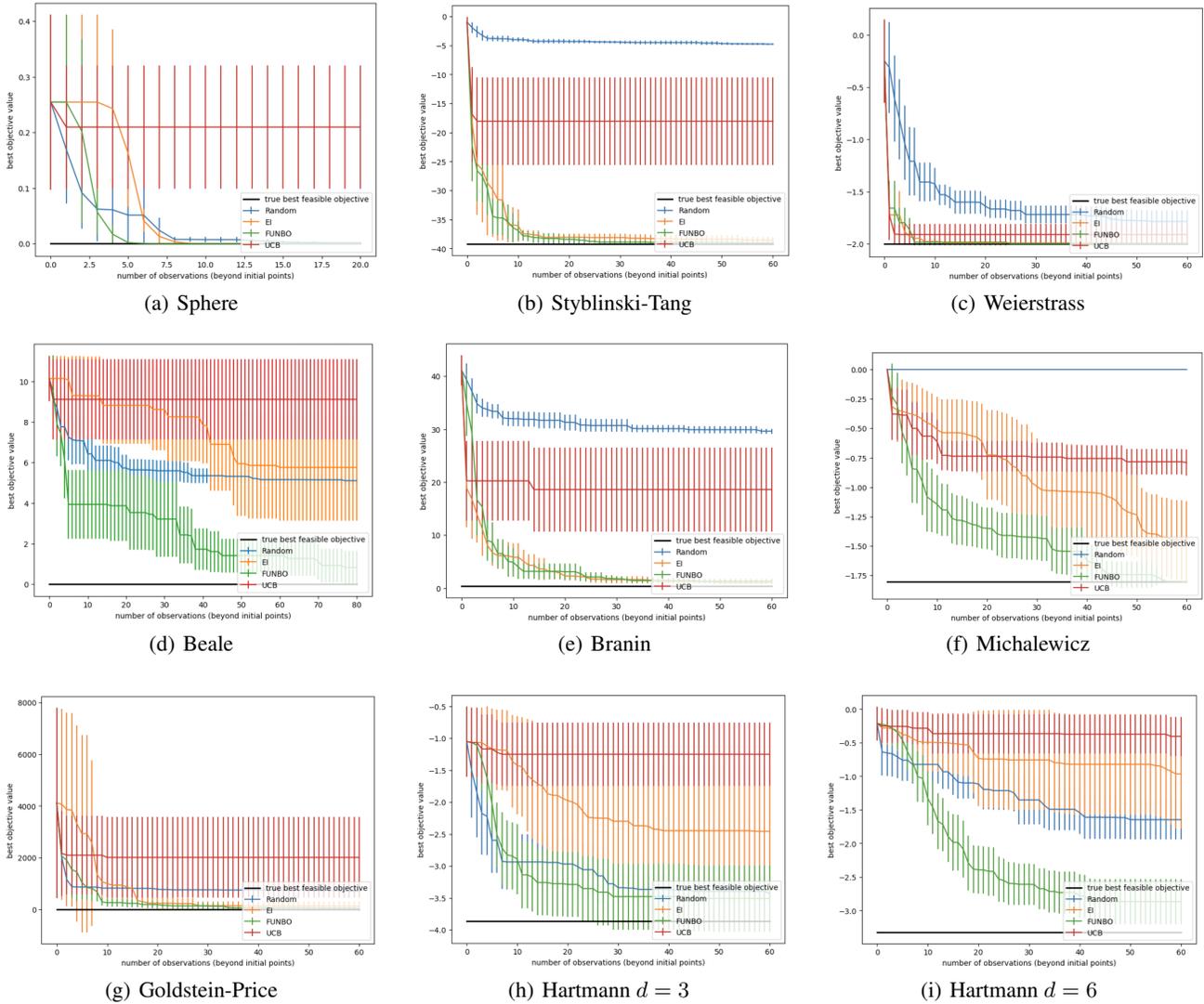


Figure 20. Performance comparison of FunBO (green) against EI (orange) and UCB (red) on the test functions for OOD-Bench when using the BoTorch evaluation pipeline where AFs are optimized numerically at each trial, GP hyperparameters are refitted at each trial using L-BFGS-B, and initial points are selected randomly. Results are averaged over 10 runs with different random initial designs for each function. Highlighted areas give \pm standard deviation over the random initial designs.