REINFORCED IN-CONTEXT BLACK-BOX OPTIMIZATION

Anonymous authors

Paper under double-blind review

Abstract

Black-Box Optimization (BBO) has found successful applications in many fields of science and engineering. Recently, there has been a growing interest in metalearning particular components of BBO algorithms to speed up optimization and get rid of tedious hand-crafted heuristics. As an extension, learning the entire algorithm from data requires the least labor from experts and can provide the most flexibility. In this paper, we propose RIBBO, a method to reinforce-learn a BBO algorithm from offline data in an end-to-end fashion. RIBBO employs expressive sequence models to learn the optimization histories produced by multiple behavior algorithms and tasks, leveraging the in-context learning ability of large models to extract task information and make decisions accordingly. Central to our method is to augment the optimization histories with *regret-to-go* tokens, which are designed to represent the performance of an algorithm based on cumulative regret over the future part of the histories. The integration of regret-to-go tokens enables RIBBO to automatically generate sequences of query points that are positive correlation to the user-desired regret, which is verified by its universally good empirical performance on diverse problems, including BBO benchmark functions, hyper-parameter optimization and robot control problems.

1 INTRODUCTION

029 030 031

003 004

010 011

012

013

014

015

016

017

018

019

021

025

026

027 028

Black-Box Optimization (BBO) Alarie et al. (2021); Audet & Hare (2017) refers to optimizing objective functions where neither analytic expressions nor derivatives of the objective are available. To solve BBO problems, we can only access the results of objective evaluation, which usually also incurs a high computational cost. Many fundamental problems in science and engineering involve optimization of expensive BBO functions, such as drug discovery Negoescu et al. (2011); Terayama et al. (2021), material design Frazier & Wang (2016); Gómez-Bombarelli et al. (2018), robot control Calandra et al. (2016); Chatzilygeroudis et al. (2019), and optimal experimental design Greenhill et al. (2020); Nguyen et al. (2023), just to name a few.

To date, a lot of BBO algorithms have been developed, among which the most prominent ones are
Bayesian Optimization (BO) Frazier (2018); Shahriari et al. (2016) and Evolutionary Algorithms
(EA) Back (1996); Zhou et al. (2019). Despite the advancements, these algorithms typically solve
BBO problems from scratch and rely on expert-derived heuristics. Consequently, they are often
hindered by slow convergence rates, and unable to leverage the inherent structures within the optimization problems Astudillo & Frazier (2021); Bai et al. (2023).

Recently, there has been a growing interest in meta-learning a particular component of the algorithms with previously collected data Arango et al. (2021); Feurer et al. (2021). Learning the component not only alleviates the need for the laborious design process of the domain experts, but also specifies the component with domain data to facilitate subsequent optimization. For example, some components in BO are proposed to be learned from data, including the surrogate model Müller et al. (2023); Perrone et al. (2018); Wang et al. (2021); Wistuba & Grabocka (2021), acquisition function Hsieh et al. (2021); Volpp et al. (2020), initialization strategy Feurer et al. (2015); Poloczek et al. (2016), and search space Li et al. (2022); Perrone & Shen (2019); some core evolutionary operations in EA have also been considered, e.g., learning the selection and mutation rate adaptation in genetic algorithm Lange et al. (2023a) or the update rules for evolution strategy Lange et al.

(2023b); the configuration of the algorithm can also be learned and dynamically adjusted throughout the optimization process Biedenkapp et al. (2020); Adriaensen et al. (2022).

There have also been some attempts to learn an entire algorithm in an End-to-End (E2E) fashion, 057 which requires almost no expert knowledge at all and provides the most flexibility across a broad range of BBO problems. However, existing practices require additional knowledge regarding the objective function during the training stage, e.g., the gradient information (often impractical for 060 BBO) Chen et al. (2017); TV et al. (2019) or online sampling from the objective function (often 061 very expensive) Maraval et al. (2023). Chen et al. (2022) proposed the OptFormer method to imitate 062 the behavior algorithms separately during training, presenting a challenge for the user to manually 063 specify which algorithm to execute during testing. Thus, these methods are less ideal for practical 064 scenarios where offline datasets are often available beforehand and a suitable algorithm for the given task has to be identified automatically without the involvement of domain experts. 065

066 In this paper, we introduce Reinforced In-context BBO (RIBBO), which learns a reinforced BBO 067 algorithm from offline datasets in an E2E fashion. RIBBO employs an expressive sequence model, 068 i.e., causal transformer, to fit the optimization histories in the offline datasets generated by executing 069 diverse behavior algorithms on multiple tasks. The sequence model is fed with previous query points and their function values, and trained to predict the distribution over the next query point. During 071 testing, the sequence model itself serves as a BBO algorithm by generating the next query points auto-regressively. Apart from this, RIBBO augments the optimization histories with regret-to-go 072 (RTG) tokens, which are calculated by summing up the regrets over the future part of the histories, 073 representing the future performance of an algorithm. A novel Hindsight Regret Relabelling (HRR) 074 strategy is proposed to update the RTG tokens during testing. By integrating the RTG tokens into 075 the modeling, RIBBO can automatically identify different algorithms, and generate sequences of 076 query points that are positive correlation to the user-desired regret. Such modeling enables RIBBO 077 to circumvent the impact of inferior data and further reinforce its performance on top of the behavior 078 algorithms. 079

We perform experiments on BBOB synthetic functions, hyper-parameter optimization and robot control problems by using some representatives of heuristic search, EA, and BO as behavior algorithms to generate the offline datasets. The results show that RIBBO can automatically generate sequences of query points related to the user-desired regret across diverse problems, and achieve good performance universally. Note that the best behavior algorithm depends on the problem at hand, and RIBBO can perform even better on some problems. Compared to the most related method OptFormer Chen et al. (2022), RIBBO also has clear advantage. In addition, we perform a series of experiments to analyze the influence of important components of RIBBO.

088

090

2 BACKGROUND

091 092 093

094 095

2.1 BLACK-BOX OPTIMIZATION

Let $f: \mathcal{X} \to \mathbb{R}$ be a black-box function, where $\mathcal{X} \subseteq \mathbb{R}^d$ is a *d*-dimensional search space. The goal of 096 BBO is to find an optimal solution $x^* \in \arg \max_{x \in \mathcal{X}} f(x)$, with the only permission of querying the 097 objective function value. Several classes of BBO algorithms have been proposed, e.g., BO Frazier 098 (2018); Shahriari et al. (2016) and EA Back (1996); Zhou et al. (2019). The basic framework of BO contains two critical components: a surrogate model, typically formalized as Gaussian Process 100 (GP) Rasmussen & Williams (2006), and an acquisition function Wilson et al. (2018), which are 101 used to model f and decide the next query point, respectively. EA is a class of heuristic optimization 102 algorithms inspired by natural evolution. It maintains a population of solutions and iterates through 103 mutation, crossover, and selection operations to find better solutions. 104

To evaluate the performance of BBO algorithms, regrets are often used. The instantaneous regret $r_t = f(x^*) - f(x_t)$ measures the gap of function values between an optimal solution x^* and the currently selected point x_t . The cumulative regret $\operatorname{Reg}_T = \sum_{i=1}^T r_i$ is the sum of instantaneous regrets in the first T iterations.

108 2.2 META-LEARNING IN BLACK-BOX OPTIMIZATION

Hand-crafted BBO algorithms usually require an expert to analyze the algorithms' behavior across a wide range of problems, a process that is both tedious and time-consuming. One solution is meta-learning Vilalta & Drissi (2002); Hospedales et al. (2021), which aims to exploit knowledge to improve the performance of learning algorithms given data from a collection of tasks. By parameter-izing a component of BBO algorithms or even an entire BBO algorithm that is traditionally manually designed, we can utilize historical data to incorporate the domain knowledge into the optimization, which may bring speedup.

117 Meta-learning particular components has been studied with different BBO algorithms. Metalearning in BO can be divided into four main categories according to "what to transfer" Bai et al. 118 (2023), including the design of the surrogate model, acquisition function, initialization strategy, and 119 search space. For surrogate model design, Wang et al. (2021) and Wistuba & Grabocka (2021) 120 parameterized the mean or kernel function of the GP model with Multi-Layer Perceptron (MLP), 121 while Perrone et al. (2018) and Müller et al. (2023) substituted GP with Bayesian linear regres-122 sion or neural process Garnelo et al. (2018); Müller et al. (2022). For acquisition function design, 123 MetaBO Volpp et al. (2020) uses Reinforcement Learning (RL) to meta-train an acquisition function 124 on a set of related tasks, and FSAF Hsieh et al. (2021) employs a Bayesian variant of deep Q-network 125 as a surrogate differentiable acquisition function trained by model-agnostic meta-learning Finn et al. 126 (2017). The remaining two categories focus on exploiting the previous good solutions to warm start 127 the optimization Feurer et al. (2015); Poloczek et al. (2016) or shrink the search space Perrone & 128 Shen (2019); Li et al. (2022). Meta-learning in EA usually focuses on learning specific evolutionary operations. For example, Lang et al. substituted core genetic operators, i.e., selection and muta-129 tion rate adaptation, with dot-product attention modules Lange et al. (2023a), and meta-learned a 130 self-attention-based architecture to discover effective and order-invariant update rules Lange et al. 131 (2023b). ALDes Zhao et al. (2024) introduces an auto-regressive learning-based approach to se-132 quentially generate components of meta-heuristic algorithms. Beyond that, Dynamic Algorithm 133 Configuration (DAC) Biedenkapp et al. (2020); Adriaensen et al. (2022) concentrates on learning 134 the configurations of algorithms, employing RL to dynamically adjust the configurations during the 135 optimization process. 136

Meta-learning entire algorithms has also been explored to obtain more flexible models. Early 137 works Chen et al. (2017); TV et al. (2019) use Recurrent Neural Network (RNN) to meta-learn 138 a BBO algorithm by optimizing the summed objective functions of some iterations. RNN uses its 139 memory state to store information about history and outputs the next query point. This work assumes 140 access to gradient information during the training phase, which is, however, usually impractical in 141 BBO problems. OptFormer Chen et al. (2022) uses a text-based transformer framework to learn an 142 algorithm, providing a universal E2E interface for BBO problems. It is trained to imitate different 143 BBO algorithms across a broad range of problems, which, however, presents a challenge for the user 144 to manually specify an algorithm for inference. Neural Acquisition Processes (NAP) Maraval et al. 145 (2023) uses transformer to meta-learn the surrogate model and acquisition function of BO jointly. Due to the lack of labeled acquisition data, NAP uses an online RL algorithm with a supervised 146 auxiliary loss for training, which requires online sampling from the expensive objective function 147 and lacks efficiency. Black-box Optimization NETworks (BONET) Krishnamoorthy et al. (2023) 148 employ a transformer model to fit regret-augmented trajectories in an offline BBO scenario, where 149 the training and testing data are from the same objective function, and a prefix sequence is required 150 to warm up the optimization before testing. OPT-GAN Lu et al. (2023) utilizes generative adversar-151 ial networks (GAN) to estimate the distribution of optimum gradually by exploration-exploitation 152 trade-off. Compared to the above state-of-the-art E2E methods, we consider the meta-BBO setting, 153 where the training datasets consist of diverse algorithms across different functions. Our approach of-154 fers the advantage of automatically identifying with RTG tokens and deploying the best-performing 155 algorithm without requiring the user to pre-specify which algorithm to use or to provide a prefix sequence during the testing phase. It utilizes a supervised learning loss for training on a fixed offline 156 dataset without the need for further interaction with the objective function. 157

158 159

- 2.3 DECISION TRANSFORMER
- 161 Transformer has emerged as a powerful architecture for sequence modeling tasks Khan et al. (2022); Wen et al. (2022); Wolf et al. (2020). A basic building block behind transformer is the self-attention



174 Figure 1: Illustration of RIBBO. Left: Data Generation. K existing BBO algorithms $\{A_j\}_{j=1}^{K}$ 175 and N BBO tasks $\{f_i\}_{i=1}^N$ are used to serve as the behavior algorithms and the training tasks, 176 respectively. The offline datasets $\{\mathcal{D}_{i,j}\}$ consist of the optimization histories $h_T = \{(x_t, y_t)\}_{t=1}^T$ 177 collected by executing each behavior algorithm A_j on each task f_i for T evaluation steps, which are 178 then augmented with the regret-to-go tokens R_t (calculated as the cumulative regret $\sum_{t'=t+1}^{T} (y^* - t)$ 179 $y_{t'}$) over the future optimization history) to generate the final dataset $\{\widehat{D}_{i,j}\}$ for training. *Right:* Training and Inference. Our model takes in triplets of (x_t, y_t, R_t) , embeds them into one token, 181 and outputs the distribution over the next query point x_{t+1} . During training, the ground-truth next 182 query point is used to minimize the loss in Eq. (4). During inference, the next query point x_{t+1} is 183 generated auto-regressively based on the current history \hat{h}_t .

185 mechanism Vaswani et al. (2017), which captures the correlation between tokens of any pair of 186 timesteps. As the scale of data and model increases, transformer has demonstrated the in-context 187 *learning* ability Brown et al. (2020), which refers to the capability of the model to infer the tasks 188 at hand based on the input contexts. Decision Transformer (DT) Chen et al. (2021) abstracts RL 189 as a sequence modeling problem, and introduces return-to-go tokens, representing the cumulative 190 rewards over future interactions. Conditioning on return-to-go tokens enables DT to correlate the 191 trajectories with their corresponding returns and generate future actions to achieve a user-specified 192 return. Inspired by DT, we will treat BBO tasks as a sequence modeling problem naturally, use a causal transformer for modeling, and train it by conditioning on future regrets. Such design is ex-193 pected to enable the learned model to distinguish algorithms with different performance and achieve 194 good performance with a user-specified low regret. 195

3 Method

196 197

205

This section presents Reinforced In-context Black-Box Optimization (RIBBO), which learns an enhanced BBO algorithm in an E2E fashion, as illustrated in Figure 1. We follow the task-distribution assumption, which is commonly adopted in meta-learning settings Finn et al. (2017); Hospedales et al. (2021); Zhou et al. (2023). Our goal is to learn a generalizable model \mathcal{M} capable of solving a wide range of BBO tasks, each associated with a BBO objective function f sampled from the task distribution $P(\mathcal{F})$, where \mathcal{F} denotes the function space.

Let [N] denote the integer set $\{1, 2, \dots, N\}$. During training, we usually access N source tasks 206 and each task corresponds to an objective function $f_i \sim P(\mathcal{F})$, where $i \in [N]$. Hereafter, we 207 use f_i to denote the task i if the context is clear. We assume that the information is available via 208 offline datasets $\mathcal{D}_{i,j}$, which are produced by executing a behavior algorithm \mathcal{A}_j on task f_i , where 209 $j \in [K]$ and $i \in [N]$. Each dataset $\mathcal{D}_{i,j} = \{h_T^{i,j,m}\}_{m=1}^M$ consists of M optimization histories 210 $h_T^{i,j,m} = \{(x_t, y_t)\}_{t=1}^T$, where x_t is the query point selected by A_j at iteration t, and $y_t = f_i(x_t)$ is 211 its objective value. If the context is clear, we will omit i, j, m and simply use h_T to denote a history 212 with length T. The initial history h_0 is defined as \emptyset . We impose no additional assumptions about 213 the behavior algorithms, allowing for a range of BBO algorithms, even random search. 214

With the datasets, we seek to learn a model $\mathcal{M}_{\theta}(x_t|h_{t-1})$, which is parameterized by θ and generates the next query point x_t by conditioning on the previous history h_{t-1} . As introduced in

Section 2.1, with a given budget T and the history h_T produced by an algorithm \mathcal{A} , we use the cumulative regret

220 221

$$\operatorname{Reg}_{T} = \sum_{t=1}^{T} (y^{*} - y_{t}) \tag{1}$$

as the evaluation metric, where y^* is the optimum value and $\{y_t\}_{t=1}^T$ are the function values in h_T .

222 3.1 METHOD OUTLINE

Given the current history h_{t-1} at iteration t, a BBO algorithm selects the next query point x_t , observes the function value $y_t = f_i(x_t)$, and updates the history $h_t = h_{t-1} \cup \{(x_t, y_t)\}$. Similar to the previous work Chen et al. (2017), we take this framework as a starting point and treat the learning of a universal BBO algorithm as learning a model \mathcal{M}_{θ} , which takes the preceding history h_{t-1} as input and outputs a distribution of the next query point x_t . The optimization histories in offline datasets provide a natural supervision for the learning process.

Suppose we have a set of histories $\{h_T\}$, generated by a single behavior algorithm \mathcal{A} on a single task f. By employing a causal transformer model \mathcal{M}_{θ} , we expect \mathcal{M}_{θ} to imitate \mathcal{A} and produce similar optimization history on f. In practice, we usually have datasets containing histories from multiple behavior algorithms $\{\mathcal{A}_j\}_{j=1}^K$ on multiple tasks $\{f_i\}_{i=1}^N$. To fit \mathcal{M}_{θ} , we use the negative log-likelihood loss

$$\mathcal{L}_{BC}(\boldsymbol{\theta}) = -\mathbb{E}_{\boldsymbol{h}_T \sim \mathcal{D}_{i,j}} \left[\sum_{t=1}^T \log \mathcal{M}_{\boldsymbol{\theta}}(\boldsymbol{x}_t | \boldsymbol{h}_{t-1}) \right].$$
(2)

To effectively minimize the loss, \mathcal{M}_{θ} needs to recognize both the task and the behavior algorithm in-context, and then imitate the optimization behavior of the corresponding behavior algorithm.

Nevertheless, naively imitating the offline datasets hinders the model since some inferior behavior algorithms may severely degenerate the model's performance. Inspired by DT Chen et al. (2021), we propose to augment the optimization history with Regret-To-Go (RTG) tokens R_t , defined as the sum of instantaneous regrets over the future history:

235

236 237

245 246

249

250

251

 $\hat{\boldsymbol{h}}_T = \{(\boldsymbol{x}_t, y_t, R_t)\}_{t=0}^T, \ R_t = \sum_{t'=t+1}^T (y^* - y_{t'}),$ (3)

where x_0 and y_0 are placeholders for padding, denoted as [PAD] in Figure 1(b), and $R_T = 0$. The augmented histories compose the augmented dataset $\hat{D}_{i,j}$, and the training objective of \mathcal{M}_{θ} becomes

 $\mathcal{L}_{\text{RIBBO}}(\boldsymbol{\theta}) = -\mathbb{E}_{\hat{\boldsymbol{h}}_T \sim \widehat{\mathcal{D}}_{i,j}} \left[\sum_{t=1}^T \log \mathcal{M}_{\boldsymbol{\theta}}(\boldsymbol{x}_t | \hat{\boldsymbol{h}}_{t-1}) \right].$ (4)

The integration of RTG tokens in the context brings identifiability of behavior algorithms, and the model \mathcal{M}_{θ} can effectively utilize them to make appropriate decisions. Furthermore, RTG tokens have a direct correlation with the metric of interest, i.e., cumulative regret Reg_T in Eq. (1). Conditioning on a lower RTG token provides a guidance to our model and reinforces \mathcal{M}_{θ} to exhibit superior performance. These advantages will be clearly shown by experiments in Section 4.4.

257 The resulting method RIBBO has implicitly utilized the in-context learning capacity of transformer 258 to guide the optimization with previous histories and the desired future regret as context. The incontext learning capacity of inferring the tasks at hand based on the input contexts has been observed 259 as the scale of data and model increases Kaplan et al. (2020). It has been explored to infer general 260 functional relationships as supervised learning or RL algorithms. For example, the model is expected 261 to predict accurately on the query input x_t by feeding the training dataset $\{(x_i, y_i)\}_{i=1}^{t-1}$ as the 262 context Guo et al. (2023); Hollmann et al. (2023); Li et al. (2023); Laskin et al. (2023) learned RL 263 algorithms using causal transformers. Here, we use it for BBO. 264

- 265 266
- 3.2 PRACTICAL IMPLEMENTATION
- 267 Next, we detail the model architecture, training, and inference of RIBBO.
- 269 Model Architecture. For the formalization of the model \mathcal{M}_{θ} , we adopt the commonly used GPT architecture Radford et al. (2018), which comprises a stack of causal attention blocks. Each block

Alg	orithm 1 Model Inference with HRR
Inp	ut : trained model \mathcal{M}_{θ} , budget T, optimum value y^*
Pro	cess:
1:	Initialize $\hat{h}_0 = \{(x_0, y_0, R_0)\}$, where x_0 and y_0 are placeholders for padding and $R_0 = 0$;
2:	for $t = 1, 2,, T$ do
3:	Generate the next query point $x_t \sim \mathcal{M}_{\theta}(\cdot \hat{h}_{t-1});$
4:	Evaluate x_t to obtain $y_t = f(x_t)$;
5:	Calculate the instantaneous regret $r = y^* - y_t$;
6:	Relabel $R_i \leftarrow R_i + r$, for each $(\boldsymbol{x}_i, y_i, R_i)$ in $\hat{\boldsymbol{h}}_{t-1}$;
7:	$\hat{h}_{t} = \hat{h}_{t-1} \cup \{(x_t, y_t, 0)\}$
8:	end for

is composed of an attention mechanism and a feed-forward network. We aggregate each triplet (x_i, y_i, R_i) using a two-layer MLP network. The output of \mathcal{M}_{θ} is a diagonal Gaussian distribution of the next query point. Note that previous works that adopt the sequence model as surrogate models Müller et al. (2022; 2023); Nguyen & Grover (2022) typically remove the positional encoding because the surrogate model should be invariant to the history order. On the contrary, our implementation preserves the positional encoding, naturally following the behavior of certain algorithms (e.g., BO or EA) and making it easier to learn from algorithms. Additionally, the positional encoding can help maintain the monotonically decreasing order of RTG tokens. More details about the architecture can be found in Appendix A.

Model Training. RTG tokens are calculated as outlined in Eq. (3) for the offline datasets before training. Since the calculation of regret requires the optimum value of task *i*, we use the bestobserved value y_{\max}^i as a proxy for the optimum. Let $\{\widehat{D}_{i,j}\}_{i \in [N], j \in [K]}$ denote the RTG augmented datasets with *N* tasks and *K* algorithms. During training, we sample a minibatch of consecutive subsequences of length $\tau < T$ uniformly from the augmented datasets. The training objective is to minimize the RIBBO loss in Eq. (4).

Model Inference. The model \mathcal{M}_{θ} generates the query points x_t auto-regressively during inference, which involves iteratively selecting a new query point x_t based on the current augmented history \hat{h}_{t-1} , evaluating the query as $y_t = f(x_t)$, and updating the history by $\hat{h}_t = \hat{h}_{t-1} \cup \{(x_t, y_t, R_t)\}$. A critical aspect of this process is how to specify the value of RTG (i.e., R_t) at every iteration t. Inspired by DT, a naive approach is to specify a desired performance as the initial RTG R_0 , and decrease it as $R_t = R_{t-1} - (y^* - y_t)$. However, this strategy has the risk of producing out-ofdistribution RTGs, since the values can fall below 0 due to an improperly selected R_0 .

Given the fact that RTGs are lower bounded by 0 and a value of 0 implies a good BBO algorithm 305 with low regret, we propose to set the immediate RTG as 0. Furthermore, we introduce a strategy 306 called Hindsight Regret Relabelling (HRR) to update previous RTGs based on the current sample 307 evaluations. The inference procedure with HRR is detailed in Algorithm 1. In line 1, the history 308 h_0 is initialized with padding placeholders x_0, y_0 and RTG $R_0 = 0$. At iteration t (i.e., lines 3–7), 309 the model \mathcal{M}_{θ} is fed with the augmented history \hat{h}_{t-1} to generate the next query point x_t in line 3, 310 followed by the evaluation procedure to obtain y_t in line 4. Then, the immediate RTG R_t is set to 0, 311 and we employ HRR to update previous RTG tokens in h_{t-1} , i.e., calculate the instantaneous regret 312 $r = y^* - y_t$ (line 5) and add r to every RTG token within h_{t-1} (line 6): 313

314

283

284

287

289

291

$$\forall 0 \le i < t, R_i \leftarrow R_i + (y^* - y_t). \tag{5}$$

315 Note that this relabelling process guarantees that $\forall 0 \leq i < t$, the RTG token $R_i = \sum_{t'=i+1}^{t} (y^* - i)^{t'}$ 316 $y_{t'}$, which can also be written as $\sum_{t'=i+1}^{T} (y^* - y_{t'})$, consistent with the definition in Eq. (3), because the immediate RTG $R_t = \sum_{t'=t+1}^{T} (y^* - y_{t'})$ is set to 0. In line 7, the history \hat{h}_t is 317 318 319 updated by expanding \hat{h}_{t-1} with $\{(x_t, y_t, 0)\}$, i.e., the current sampling and its immediate RTG 320 $R_t = 0$. The above process is repeated until reaching the budget T. Thus, we can find that HRR not 321 only exploits the full potential of \mathcal{M}_{θ} through using 0 as the immediate RTG and thereby demands the model to generate the most advantageous decisions, but also preserves the calculation of RTG 322 tokens following the same way as the training data, i.e., representing the cumulative regret over 323 future optimization history.

324 3.3 DATA GENERATION

Finally, we give some guidelines about data generation for using the proposed RIBBO method.

Data Collection. Given a set of tasks $\{f_i\}_{i=1}^N$ sampled from the task distribution $P(\mathcal{F})$, we can employ a diverse set of behavior algorithms for data collection. For example, we can select some representatives from different types of BBO algorithms, e.g., BO and EA. Datasets $\mathcal{D}_{i,j}$ are obtained by using each behavior algorithm to optimize each task with different random seeds. Each optimization history $h_T = \{(x_t, y_t)\}_{t=1}^T$ in $\mathcal{D}_{i,j}$ is then augmented with RTG tokens R_t , which is computed as in Eq. (3). The resulting histories $\hat{h}_T = \{(x_t, y_t, R_t)\}_{t=0}^T$ compose the final datasets $\hat{\mathcal{D}}_{i,j}$ for model training.

334 Data Normalization. To provide a unified interface and balance the statistic scales across tasks, it 335 is important to apply normalization to the inputs to our model. We normalize the point x by (x - x)336 $(x_{\rm min})/(x_{\rm max}-x_{\rm min})$, with $x_{\rm max}$ and $x_{\rm min}$ being the upper and lower bounds of the search space, 337 respectively. For the function value y, we apply random scaling akin to previous works Wistuba 338 & Grabocka (2021); Chen et al. (2022). That is, when sampling a history h_{τ} from the datasets 339 $\mathcal{D}_{i,j}$, we randomly sample the lower bound $l \sim \mathcal{U}(y_{\min}^i - \frac{s}{2}, y_{\min}^i + \frac{s}{2})$ and the upper bound $u \sim \mathcal{U}(y_{\max}^i - \frac{s}{2}, y_{\max}^i + \frac{s}{2})$, where \mathcal{U} stands for uniform distribution, y_{\min}^i, y_{\max}^i denote the observed 340 341 minimum and maximum values for f_i , and $s = y_{max}^i - y_{min}^i$; the values y_t in h_{τ} are then normalized 342 by $(y_t - l)/(u - l)$ for training. The RTG tokens are calculated accordingly with the normalized 343 values. The random normalization can make a model exhibit invariance across various scales of y. 344 For inference, the average values of the best-observed and worst-observed values across the training 345 tasks are used to normalize y.

346 347

348

4 EXPERIMENTS

In this section, we examine the performance of RIBBO on a wide range of tasks, including synthetic functions, Hyper-Parameter Optimization (HPO) and robot control problems. The model architecture and hyper-parameters are maintained consistently across these problems. We train our model using five distinct random seeds, ranging from 0 to 4, and each trained model is run five times independently during the execution phase. We will report the average performance and standard deviation. Details of the model hyper-parameters are given in Appendix A. The codes are provided in the supplementary.

356 357

358

4.1 EXPERIMENTAL SETUP

Benchmarks. We use BBO Benchmarks BBOB Elhara et al. (2019), HPO-B Arango et al. (2021), 359 and rover trajectory planning task Wang et al. (2018). The BBOB suite, a comprehensive and widely 360 used benchmark in the continuous domain, consists of 24 synthetic functions. For each function, a 361 series of linear and non-linear transformations are implemented on the search space to obtain a 362 distribution of functions with similar properties. According to the properties of these functions, they can be divided into 5 categories, and we select one from each category due to resource constraints, 364 including Greiwank Rosenbrock, Lunacek, Rastrigin, Rosenbrock, and Sharp Ridge. HPO-B is a 365 commonly used HPO benchmark and consists of a series of HPO problems. Each problem is to 366 optimize a machine learning model across various datasets, and an XGBoost model is provided as 367 the objective function for evaluation in a continuous space. We conduct experiments on two widely 368 used models, SVM and XGBoost, in the continue domain. For robot control optimization, we 369 perform experiments on rover trajectory planning task, which is a trajectory optimization problem to emulate rover navigation. Similar to Elhara et al. (2019); Volpp et al. (2020), we implement random 370 translations and scalings to the search space to construct a distribution of functions. For BBOB 371 and rover problems, we sample a set of functions from the task distribution as training and test 372 tasks, while for HPO-B, we use the meta-training/test task splits provided by the authors. Detailed 373 explanations of the benchmarks can be found in Appendix B.1. 374

375 Data. Similar to OptFormer Chen et al. (2022), we employ 7 behavior algorithms, i.e., Random
376 Search, Shuffled Grid Search, Hill Climbing, Regularized Evolution Real et al. (2019), Eagle Strat377 egy Yang & Deb (2010), CMA-ES Hansen (2016), and GP-EI Balandat et al. (2020), which are
representatives of heuristic search, EA, and BO, respectively. Datasets are generated by employing

each behavior algorithm to optimize various training functions sampled from the task distribution,
using different random seeds. For inference, new test functions are sampled from the same task
distribution to serve as the test set. Specifically, for the HPO-B problem, the meta-training/test splits
have been predefined by the authors and we adhere to this standard setup. Additional information
about the behavior algorithms and datasets can be found in Appendix B.2 and B.3.

384 4.2 BASELINES385

As RIBBO is an in-context E2E model, the most related baselines are those also training an E2E
model with offline datasets, including Behavior Cloning (BC) Bain & Sammut (1995) and OptFormer Chen et al. (2022). Their hyper-parameters are set as same as that of our model for fairness.
Note that the seven behavior algorithms used to generate datasets are also important baselines, and
included for comparison as well.

BC uses the same transformer architecture as RIBBO. The only difference is that we do not feed RTG tokens into the model of BC and train to minimize the BC loss in Eq. (2). When the solutions are generated auto-regressively, BC tends to imitate the average behavior of various behavior algorithms. Consequently, the inclusion of underperforming algorithms, e.g., Random Search and Shuffled Grid Search, may significantly degrade the performance. To mitigate this issue, we have also trained the model by excluding these underperforming algorithms, denoted as BC Filter.

397 **OptFormer** employs a transformer to imitate the behaviors of a set of algorithms and an algo-398 rithm identifier usually needs to be specified manually during inference for superior performance. 399 Its original implementation is built upon a text-based transformer with a large training scale. In 400 this paper, we re-implement a simplified version of OptFormer where we only retain the algorithm identifier within the metadata. The initial states, denoted as x_0 and y_0 , are used to distinguish be-401 tween algorithms. They are obtained by indexing the algorithm type through an embedding layer, 402 thereby aligning the initial states with the specific imitated algorithm. This enables the identifi-403 cation of distinct behavior algorithms within the simplified OptFormer. Further details about the 404 re-implementation can be found in Appendix C. 405

406

408

383

407 4.3 MAIN RESULTS

The results are shown in Figure 2. For the sake of clarity in visualization, we have omitted the 409 inclusion of Random Search and Shuffled Grid Search due to their poor performance from start to 410 finish. We can observe that RIBBO achieves superior or at least equivalent efficacy in comparison 411 to the best behavior algorithm on each problem except SVM and rover. This demonstrates the 412 versatility of RIBBO, while the most effective behavior algorithm depends upon the specific problem 413 at hand, e.g., the best behavior algorithms on Lunacek, Rastrigin and XGBoost are GP-EI, Eagle 414 Strategy and CMA-ES, respectively. Note that the good performance of RIBBO does not owe to 415 the memorization of optimal solutions, as the search space is transformed randomly, resulting in 416 variations in optimal solutions across different functions from the same distribution. It is because 417 RIBBO is capable of using RTG tokens to identify algorithms and reinforce the performance on top of the behavior algorithms, which will be clearly shown later. We can also observe that RIBBO 418 performs extremely well in the early stage, which draws the advantage from the HRR strategy, i.e., 419 employing 0 as the immediate RTG to generate the optimal potential solutions. 420

421 RIBBO does not perform well on the SVM problem, which may be due to the problem's low-422 dimensional nature (only three parameters) and its relative simplicity for optimization. Behavior algorithms can achieve good performance easily, while the complexity of RIBBO's training and 423 inference processes could instead result in the performance degradation. For the rover problem 424 where GP-EI performs the best, we collect less data from GP-EI than other behavior algorithms due 425 to the high time cost. This may limit RIBBO's capacity to leverage the high-quality data from GP-426 EI, given its small proportion relative to the data collected from other behavior algorithms. Despite 427 this, RIBBO is still the runner-up, significantly surpassing the other behavior algorithms. 428

Compared with BC and BC Filter, RIBBO performs consistently better except on the SVM problem. BC tends to imitate the average behavior of various algorithms, and its poor performance is due to the aggregation of behavior algorithms with inferior performance. BC Filter is generally better than BC, because the data from the two underperforming behavior algorithms, i.e., Random



Figure 2: Performance comparison among RIBBO, BC, BC Filter, OptFormer, and behavior algorithms on synthetic functions, HPO, and robot control problems. The *y*-axis is the normalized average objective value, and the length of vertical bars represents the standard deviation.

Search and Shuffled Grid Search, are manually excluded for the training of BC Filter. As introduced
before, OptFormer requires to manually specify which algorithm to execute. We have specified the
behavior algorithm Eagle Strategy in Figure 2, which obtains good overall performance on these
problems. It can be observed that OptFormer displays a close performance to Eagle Strategy, while
RIBBO performs better. More results about the imitation capacity of OptFormer can be found in
Appendix C.

Why Does RIBBO Behave Well? To better understand RIBBO, we train the model using only two behavior algorithms, Eagle Strategy and Random Search, which represent a good algorithm and an underperforming one, respectively. Figure 3(a) visualizes the contour lines of the 2D Branin function and the sampling points of RIBBO, Eagle Strategy, and Random Search, represented by red, orange, and gray points, respectively. The arrows are used to represent the optimization trajectory of RIBBO. Note that the two parameters of Branin have been scaled to [-1, 1] for better visualization. It can be observed that RIBBO makes a prediction preferring Eagle Strategy over Random Search, indicating its capability to automatically identify the quality of training data. Additionally, RIBBO achieves the exploration and exploitation trade-off capability using its knowledge about the task obtained during training, thus generating superior solutions over the ones in the training dataset.

Generalization. We also conduct experiments to examine the generalization capabilities of RIBBO. For this purpose, we train the model on all 24 BBOB synthetic functions simultaneously with the results shown in Figure 3(b). To aggregate results across functions with different output scaling, we normalize all the functions adhering to previous literature Turner et al. (2020); Arango et al. (2021); Chen et al. (2022). The results suggest that RIBBO demonstrates strong generalizing to a variety of functions with different properties. Further experiments are conducted to examine the cross-distribution generalization to unseen function distributions. The model is trained on 4 of the 5 cho-sen synthetic functions and tested with the remaining one. Note that each function (i.e., Greiwank Rosenbrock, Lunacek, Rastrigin, Rosenbrock, and Sharp Ridge) here actually represents a distribu-tion of functions with similar properties, and a set of functions is sampled from each distribution as introduced before. The results suggest that RIBBO has a strong generalizing ability to unseen function distributions. Due to the space limitation, the results and more details on cross-distribution generalization are deferred to Appendix D. The good generalization of RIBBO can be attributed to the paradigm of learning the entire algorithm, which can acquire general knowledge, such as exploration and exploitation trade-off from data, as observed in Chen et al. (2017). In contrast, such generalization may be limited if we learn surrogate models from data, because the function landscape inherent to surrogate models will contain only the knowledge of similar functions.



Figure 3: (a) Visualization of the contour lines of 2D Branin function and sampling points of RIBBO (red), Eagle Strategy (orange), and Random Search (gray), where the arrows represent the optimization trajectory of RIBBO. (b) Generalization by training the model across all 24 BBOB synthetic functions simultaneously. The results across functions with different output scaling are normalized to obtain the aggregate results. The legend shares with that of Figure 2. (c) Initial RTG R_0 's influence on performance. (d) RTG update strategy comparison between HRR and the naive strategy with various initial RTG R_0 .

502 4.4 ABLATION STUDIES

495

496

497

498

499

500

501

RIBBO augments the histories with RTG tokens, facilitating distinguishing algorithms and automatically generating algorithms with user-specified performance. Next, we will verify the effectiveness of RTG conditioning and HRR strategy.

Influence of Initial RTG Token R_0 . By incorporating RTG tokens, RIBBO is able to attend to 507 RTGs and generate optimization trajectories based on the specified initial RTG token R_0 . To val-508 idate this, we examine the performance of RIBBO with different values of R_0 , and the results are 509 presented in Figure 3(c). Here, the RTG values, instead of normalized objective values, are used as 510 the y-axis. We can observe that the cumulative regrets of the generated query sequence do correlate 511 with the specified RTG, indicating that RIBBO establishes the connection between regret and gen-512 eration. We also conduct experiments to explore the effect of varying the immediate RTG R_t , and 513 it is observed that setting to a value larger than 0 will decrease the performance and converge to a 514 worse value. Please see Appendix E. 515

Effectiveness of HRR. A key point of the inference procedure is how to update the value of RTG to-516 kens at each iteration. To assess the effectiveness of the proposed strategy HRR outlined in Eq. (5), 517 we compare it with the naive strategy, that sets an initial RTG token R_0 and decreases it by the 518 one-step regret after each iteration, a method that employs the same updating strategy mechanism as 519 DT. The results are shown in Figure 3(d). The naive strategy displays distinct behaviors depending 520 on the initial setting of R_0 . Specifically, when $R_0 = 0$, i.e., the lower bound of regret, the model 521 performs well initially. However, as the optimization progresses, the RTG tokens gradually decrease 522 to negative values, leading to poor performance since negative RTGs are out-of-distribution values. Using $R_0 = 5$ compromises the initial performance, as the model may not select the most aggressive 523 solutions with a high R_0 . However, a higher initial R_0 yields better convergence value since it pre-524 vents out-of-distribution RTGs in later stage. The proposed HRR strategy consistently outperforms 525 across the whole optimization stage, because setting the immediate RTG to 0 encourages the model 526 to make the most advantageous decisions at every iteration, while hindsight relabeling of previous 527 RTG tokens, as specified in Eq. (5), ensures that these values remain meaningful and feasible. 528

Further Studies. We also study the effects of the method to aggregate (x_i, y_i, R_i) tokens, the normalization method for y, the model size, and the sampled subsequence length τ . Please see Appendix F. More visualizations illustrating the effects of random transformations on the search space are detailed in Appendix G.

- 533 534
- 5 CONCLUSION

This paper proposes RIBBO, which employs a transformer architecture to learn a reinforced BBO algorithm from offline datasets in an E2E fashion. By incorporating RTG tokens into the optimization
histories, RIBBO can automatically generate optimization trajectories satisfying the user-desired regret. Comprehensive experiments on BBOB synthetic functions, HPO and robot control problems
show the versatility of RIBBO.

540 REFERENCES 541

549

551

560

561

- Steven Adriaensen, André Biedenkapp, Gresa Shala, Noor H. Awad, Theresa Eimer, Marius Lin-542 dauer, and Frank Hutter. Automated dynamic algorithm configuration. Journal of Artificial Intel-543 ligence Research, 75:1633–1699, 2022. 544
- Stéphane Alarie, Charles Audet, Aïmen E. Gheribi, Michael Kokkolaras, and Sébastien Le Digabel. 546 Two decades of blackbox optimization applications. EURO Journal on Computational Optimiza-547 tion, 9:100011, 2021.
- 548 Sebastian Pineda Arango, Hadi S Jomaa, Martin Wistuba, and Josif Grabocka. HPO-B: A large-scale reproducible benchmark for black-box HPO based on OpenML. arXiv preprint 550 arXiv:2106.06257, 2021.
- Raul Astudillo and Peter I. Frazier. Thinking inside the box: A tutorial on grey-box Bayesian 552 optimization. In Proceedings of the 51st Winter Simulation Conference (WSC'21), pp. 1–15, 553 Phoenix, AZ, 2021. 554
- 555 Charles Audet and Warren Hare. Derivative-Free and Blackbox Optimization. Springer, 2017. 556
- Thomas Back. Evolutionary Algorithms in Theory and Practice: Evolution Strategies, Evolutionary 557 Programming, Genetic Algorithms. Oxford University Press, 1996. 558
- 559 Tianyi Bai, Yang Li, Yu Shen, Xinyi Zhang, Wentao Zhang, and Bin Cui. Transfer learning for Bayesian optimization: A survey. arXiv preprint arXiv:2302.05927, 2023.
- Michael Bain and Claude Sammut. A framework for behavioural cloning. Machine Intelligence, 15: 562 103-129, 1995. 563
- 564 Maximilian Balandat, Brian Karrer, Daniel R. Jiang, Samuel Daulton, Benjamin Letham, An-565 drew Gordon Wilson, and Eytan Bakshy. BoTorch: A framework for efficient Monte-Carlo 566 Bayesian optimization. In Advances in Neural Information Processing Systems 33 (NeurIPS'20), 567 pp. 10113-10124, Virtual, 2020.
- 568 André Biedenkapp, H. Furkan Bozkurt, Theresa Eimer, Frank Hutter, and Marius Lindauer. Dy-569 namic algorithm configuration: Foundation of a new meta-algorithmic framework. In Proceed-570 ings of the 24th European Conference on Artificial Intelligence (ECAI'20), pp. 427–434, Santiago, 571 Spain, 2020. 572
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-573 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, 574 Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. 575 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, 576 Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, 577 Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In Advances in 578 Neural Information Processing Systems 33 (NeurIPS'20), pp. 1877–1901, Virtual, 2020. 579
- Roberto Calandra, André Seyfarth, Jan Peters, and Marc Peter Deisenroth. Bayesian optimization 580 for learning gaits under uncertainty - An experimental comparison on a dynamic bipedal walker. 581 Annals of Mathematics and Artificial Intelligence, 76(1-2):5–23, 2016. 582
- 583 Konstantinos Chatzilygeroudis, Vassilis Vassiliades, Freek Stulp, Sylvain Calinon, and Jean-584 Baptiste Mouret. A survey on policy search algorithms for learning robot controllers in a handful 585 of trials. IEEE Transactions on Robotics, 36(2):328-347, 2019.
- Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter 587 Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning 588 via sequence modeling. In Advances in Neural Information Processing Systems 34 (NeurIPS'21), 589 pp. 15084–15097, Virtual, 2021. 590
- Yutian Chen, Matthew W Hoffman, Sergio Gómez Colmenarejo, Misha Denil, Timothy P Lillicrap, Matt Botvinick, and Nando Freitas. Learning to learn without gradient descent by gradient de-592 scent. In Proceedings of the 34th International Conference on Machine Learning (ICML'17), pp. 748–756, Sydney, Australia, 2017.

613

618

620

621

622

623

624

628

629

630

594	Yutian Chen, Xingyou Song, Chansoo Lee, Zi Wang, Richard Zhang, David Dohan, Kazuya
595	Kawakami, Greg Kochanski, Arnaud Doucet, Marc'Aurelio Ranzato, Sagi Perel, and Nando
596	de Freitas. Towards learning universal hyperparameter optimizers with transformers. In Advances
597	in Neural Information Processing Systems 35 (NeurIPS'22), pp. 32053–32068, New Orleans, LA,
598	2022.
500	

- Ouassim Elhara, Konstantinos Varelas, Duc Nguyen, Tea Tusar, Dimo Brockhoff, Nikolaus Hansen, and Anne Auger. COCO: The large scale black-box optimization benchmarking (BBOB-largescale) test suite. arXiv preprint arXiv:1903.06396, 2019.
- David Eriksson, Michael Pearce, Jacob R. Gardner, Ryan Turner, and Matthias Poloczek. Scalable
 global optimization via local Bayesian optimization. In *Advances in Neural Information Process- Systems 32 (NeurIPS'19)*, pp. 5497–5508, Vancouver, Canada, 2019.
- Matthias Feurer, Jost Springenberg, and Frank Hutter. Initializing Bayesian hyperparameter optimization via meta-learning. In *Proceedings of the 29th AAAI Conference on Artificial Intelligence (AAAI'15)*, pp. 1128–1135, Austin, TX, 2015.
- Matthias Feurer, Jan N Van Rijn, Arlind Kadra, Pieter Gijsbers, Neeratyoy Mallik, Sahithya Ravi,
 Andreas Müller, Joaquin Vanschoren, and Frank Hutter. OpenML-Python: An extensible Python
 API for OpenML. *The Journal of Machine Learning Research*, 22(1):4573–4577, 2021.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation
 of deep networks. In *Proceedings of the 34th International Conference on Machine Learning* (*ICML'17*), pp. 1126–1135, Sydney, Australia, 2017.
- ⁶¹⁷ Peter I Frazier. A tutorial on Bayesian optimization. *arXiv preprint arXiv:1807.02811*, 2018.
- 619 Peter I Frazier and Jialei Wang. *Bayesian Optimization for Materials Design*. Springer, 2016.
 - Marta Garnelo, Dan Rosenbaum, Christopher Maddison, Tiago Ramalho, David Saxton, Murray Shanahan, Yee Whye Teh, Danilo Rezende, and SM Ali Eslami. Conditional neural processes. In *Proceedings of the 35th International Conference on Machine Learning (ICML'18)*, pp. 1704–1713, Stockholm, Sweden, 2018.
- R. Gómez-Bombarelli, D. K. Duvenaud, J. M. Hernández-Lobato, J. Aguilera-Iparraguirre, T.D.
 Hirzel, R. P. Adams, and A. Aspuru-Guzik. Automatic chemical design using a data-driven continuous representation of molecules. *ACS Central Science*, 4(2):268 – 276, 2018.
 - Stewart Greenhill, Santu Rana, Sunil Gupta, Pratibha Vellanki, and Svetha Venkatesh. Bayesian optimization for adaptive experimental design: A review. *IEEE Access*, 8:13937–13948, 2020.
- Tianyu Guo, Wei Hu, Song Mei, Huan Wang, Caiming Xiong, Silvio Savarese, and Yu Bai. How do
 transformers learn in-context beyond simple functions? A case study on learning with representations. *arXiv preprint arXiv:2310.10616*, 2023.
- Nikolaus Hansen. The CMA evolution strategy: A tutorial. *arXiv preprint arXiv:1604.00772*, 2016.
- ⁶³⁶ Noah Hollmann, Samuel Müller, Katharina Eggensperger, and Frank Hutter. TabPFN: A trans ⁶³⁷ former that solves small tabular classification problems in a second. In *Proceedings of the 11th International Conference on Learning Representations (ICLR'23)*, Kigali, Rwanda, 2023.
- Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. Meta-learning in neural networks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9): 5149–5169, 2021.
- Bing-Jing Hsieh, Ping-Chun Hsieh, and Xi Liu. Reinforced few-shot acquisition function learning for Bayesian optimization. In Advances in Neural Information Processing Systems 34 (NeurIPS'21), pp. 7718–7731, Virtual, 2021.
- 647 Donald R Jones, Matthias Schonlau, and William J Welch. Efficient global optimization of expensive black-box functions. *Journal of Global Optimization*, 13(4):455–492, 1998.

658

659

660

661

- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. arXiv preprint arXiv:2001.08361, 2020.
- Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and
 Mubarak Shah. Transformers in vision: A survey. *ACM Computing Surveys*, 54(10s):1–41, 2022.
- Siddarth Krishnamoorthy, Satvik Mashkaria, and Aditya Grover. Generative pretraining for blackbox optimization. In *Proceedings of the 40th International Conference on Machine Learning* (*ICML'23*), pp. 24173–24197, Honolulu, HI, 2023.
 - Robert Tjarko Lange, Tom Schaul, Yutian Chen, Chris Lu, Tom Zahavy, Valentin Dalibard, and Sebastian Flennerhag. Discovering attention-based genetic algorithms via meta-black-box optimization. In *Proceedings of the 25th Conference on Genetic and Evolutionary Computation* (GECCO'23), pp. 929–937, Lisbon, Portugal, 2023a.
- Robert Tjarko Lange, Tom Schaul, Yutian Chen, Tom Zahavy, Valentin Dalibard, Chris Lu, Satin der Singh, and Sebastian Flennerhag. Discovering evolution strategies via meta-black-box op timization. In *Proceedings of the 11th International Conference on Learning Representations* (*ICLR'23*), Kigali, Rwanda, 2023b.
- Michael Laskin, Luyu Wang, Junhyuk Oh, Emilio Parisotto, Stephen Spencer, Richie Steigerwald, DJ Strouse, Steven Stenberg Hansen, Angelos Filos, Ethan Brooks, maxime gazeau, Himanshu Sahni, Satinder Singh, and Volodymyr Mnih. In-context reinforcement learning with algorithm distillation. In *Proceedings of the 11th International Conference on Learning Representations* (*ICLR'23*), Kigali, Rwanda, 2023.
- Yang Li, Yu Shen, Huaijun Jiang, Tianyi Bai, Wentao Zhang, Ce Zhang, and Bin Cui. Transfer learning based search space design for hyperparameter tuning. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD'22)*, pp. 967–977, Washington, DC, 2022.
- Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers
 as algorithms: Generalization and stability in in-context learning. In *Proceedings of the 40th In- ternational Conference on Machine Learning (ICML'23)*, pp. 19565–19594, Honolulu, HI, 2023.
- Minfang Lu, Shuai Ning, Shuangrong Liu, Fengyang Sun, Bo Zhang, Bo Yang, and Lin Wang.
 OPT-GAN: A broad-spectrum global optimizer for black-box problems by learning distribution. In *Proceedings of the 37th AAAI Conference on Artificial Intelligence (AAAI'23)*, pp. 12462–12472, Washington, DC, 2023.
- Alexandre Max Maraval, Matthieu Zimmer, Antoine Grosnit, and Haitham Bou Ammar. End-to-end meta-Bayesian optimisation with transformer neural processes. In *Advances in Neural Information Processing Systems 36 (NeurIPS'23)*, New Orleans, LA, 2023.
- Samuel Müller, Noah Hollmann, Sebastian Pineda Arango, Josif Grabocka, and Frank Hutter. Trans formers can do Bayesian inference. In *Proceedings of the 10th International Conference on Learning Representations (ICLR'22)*, Virtual, 2022.
- Samuel Müller, Matthias Feurer, Noah Hollmann, and Frank Hutter. PFNs4BO: In-context learning for Bayesian optimization. In *Proceedings of the 40th International Conference on Machine Learning (ICML'23)*, pp. 25444–25470, Honolulu, HI, 2023.
- Diana M. Negoescu, Peter I. Frazier, and Warren B. Powell. The knowledge-gradient algorithm for
 sequencing experiments in drug discovery. *INFORMS Journal on Computing*, 23(3):346–363,
 2011.
- Tung Nguyen and Aditya Grover. Transformer neural processes: Uncertainty-aware meta learning via sequence modeling. In *Proceedings of the 39th International Conference on Machine Learning (ICML'22)*, pp. 16569–16594, Baltimore, MD, 2022.
- 701 Tung Nguyen, Sudhanshu Agrawal, and Aditya Grover. ExPT: Synthetic pretraining for few-shot experimental design. *arXiv preprint arXiv:2310.19961*, 2023.

702 703 704	Valerio Perrone and Huibin Shen. Learning search spaces for Bayesian optimization: Another view of hyperparameter transfer learning. In <i>Advances in Neural Information Processing Systems 32 (NeurIPS'19)</i> , pp. 12751–12761, Vancouver, Canada, 2019.
705 706 707 708	Valerio Perrone, Rodolphe Jenatton, Matthias W. Seeger, and Cédric Archambeau. Scalable hyperparameter transfer learning. In <i>Advances in Neural Information Processing Systems 31 (NeurIPS'18)</i> , pp. 6846–6856, Montreal, Canada, 2018.
709 710 711	Matthias Poloczek, Jialei Wang, and Peter I Frazier. Warm starting Bayesian optimization. In <i>Proceedings of the 46th Winter Simulation Conference (WSC'16)</i> , pp. 770–781, Washington, DC, 2016.
712 713 714 715	Ofir Press, Noah A. Smith, and Mike Lewis. Train short, test long: Attention with linear biases en- ables input length extrapolation. In <i>Proceedings of the 10th International Conference on Learning</i> <i>Representations (ICLR'22)</i> , Virtual, 2022.
716 717 718	Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language under- standing by generative pre-training. <i>OpenAI Blog</i> , 2018.
719 720	C. E. Rasmussen and C. K. I. Williams. <i>Gaussian Processes for Machine Learning</i> . The MIT Press, 2006.
721 722 723 724	Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V Le. Regularized evolution for image classifier architecture search. In <i>Proceedings of the 33rd AAAI Conference on Artificial Intelligence (AAAI'19)</i> , pp. 4780–4789, Honolulu, HI, 2019.
725 726	Jürgen Schmidhuber. Reinforcement learning upside down: Don't predict rewards - Just map them to actions. <i>arXiv preprint arXiv:1912.02875</i> , 2019.
727 728 729 730	Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P. Adams, and Nando de Freitas. Taking the human out of the loop: A review of Bayesian optimization. <i>Proceedings of the IEEE</i> , 104(1): 148–175, 2016.
731 732 733 734	Xingyou Song, Sagi Perel, Chansoo Lee, Greg Kochanski, and Daniel Golovin. Open Source Vizier: Distributed infrastructure and API for reliable and flexible black-box optimization. In <i>Proceedings</i> <i>of the 1st International Conference on Automated Machine Learning (AutoML Conference'22)</i> , pp. 1–17, Baltimore, MD, 2022.
735 736 737	Kei Terayama, Masato Sumita, Ryo Tamura, and Koji Tsuda. Black-box optimization for automated discovery. <i>Accounts of Chemical Research</i> , 54(6):1334–1346, 2021.
738 739 740 741 742	Ryan Turner, David Eriksson, Michael McCourt, Juha Kiili, Eero Laaksonen, Zhen Xu, and Isabelle Guyon. Bayesian optimization is superior to random search for machine learning hyperparameter tuning: Analysis of the black-box optimization challenge 2020. In Hugo Jair Escalante and Katja Hofmann (eds.), <i>Advances in Neural Information Processing Systems 33 (NeurIPS'20) Competition and Demonstration Track</i> , pp. 3–26, Virtual, 2020.
743 744 745 746	Vishnu TV, Pankaj Malhotra, Jyoti Narwariya, Lovekesh Vig, and Gautam Shroff. Meta-learning for black-box optimization. In <i>Proceedings of European Conference on Machine Learning and Knowledge Discovery in Databases (ECML PKDD'19)</i> , pp. 366–381, Würzburg, Germany, 2019.
747 748 749	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Advances in Neural Infor- mation Processing Systems 30 (NIPS'17), pp. 5998–6008, Long Beach, Canada, 2017.
750 751 752	Ricardo Vilalta and Youssef Drissi. A perspective view and survey of meta-learning. <i>Artificial Intelligence Review</i> , 18(2):77–95, 2002.
753 754 755	Michael Volpp, Lukas P. Fröhlich, Kirsten Fischer, Andreas Doerr, Stefan Falkner, Frank Hutter, and Christian Daniel. Meta-learning acquisition functions for transfer learning in Bayesian optimization. In <i>Proceedings of the 8th International Conference on Learning Representations (ICLR'20)</i> , Addis Ababa, Ethiopia, 2020.

- 756 Zi Wang, Clement Gehring, Pushmeet Kohli, and Stefanie Jegelka. Batched large-scale Bayesian optimization in high-dimensional spaces. In Proceedings of the 21st International Conference on 758 Artificial Intelligence and Statistics (AISTATS'18), pp. 745–754, Playa Blanca, Spain, 2018. 759 Zi Wang, George E Dahl, Kevin Swersky, Chansoo Lee, Zachary Nado, Justin Gilmer, Jasper 760 Snoek, and Zoubin Ghahramani. Pre-trained Gaussian processes for Bayesian optimization. arXiv 761 preprint arXiv:2109.08215, 2021. 762 763 Qingsong Wen, Tian Zhou, Chaoli Zhang, Weiqi Chen, Ziqing Ma, Junchi Yan, and Liang Sun. 764 Transformers in time series: A survey. arXiv preprint arXiv:2202.07125, 2022. 765 James T. Wilson, Frank Hutter, and Marc Peter Deisenroth. Maximizing acquisition functions for 766 Bayesian optimization. In Advances in Neural Information Processing Systems 31 (NeurIPS'18), 767 pp. 9906–9917, Montréal, Canada, 2018. 768 769 Martin Wistuba and Josif Grabocka. Few-shot Bayesian optimization with deep kernel surrogates. In Proceedings of the 9th International Conference on Learning Representations (ICLR'21), Virtual, 770 2021. 771 772 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, 773 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick 774 von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, 775 Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural 776 language processing. In Proceedings of the 25th Conference on Empirical Methods in Natural Language Processing: System Demonstrations (EMNLP'20), pp. 38-45, Virtual, 2020. 777 778 Xin-She Yang and Suash Deb. Eagle strategy using Lévy walk and firefly algorithms for stochastic 779 optimization. In Proceedings of the 4th Nature Inspired Cooperative Strategies for Optimization 780 (*NICSO'10*), pp. 101–111, Granada, Spain, 2010. 781 Qi Zhao, Tengfei Liu, Bai Yan, Qiqi Duan, Jian Yang, and Yuhui Shi. Automated metaheuristic 782 algorithm design with autoregressive learning. arXiv preprint arXiv:2405.03419, 2024. 783 784 Jianan Zhou, Yaoxin Wu, Wen Song, Zhiguang Cao, and Jie Zhang. Towards omni-generalizable 785 neural methods for vehicle routing problems. In Proceedings of the 40th International Conference 786 on Machine Learning (ICML'23), pp. 42769–42789, Honolulu, HI, 2023. 787 Zhi-Hua Zhou, Yang Yu, and Chao Qian. Evolutionary Learning: Advances in Theories and Algo-788 rithms. Springer, 2019. 789 790 791 А MODEL DETAILS 792 793 We employ the commonly used GPT architecture Radford et al. (2018) and the hyper-parameters are 794 maintained consistently across the problems. Details can be found in Table 1. For BC, BC Filter and OptFormer, the hyper-parameters are set as same as those of our model. The training takes about 20 796 hours on 1 GPU (Nvidia RTX 4090). 797 798 В DETAILS OF EXPERIMENTAL SETUP 799 800 **B**.1 BENCHMARKS 801 802 • **BBOB** Elhara et al. (2019) is a widely used synthetic BBO benchmark, consisting of 24 synthetic functions in the continuous domain. This benchmark makes a series of trans-804
- formations in the search space, such as linear transformations (e.g., translation, rotation, scaling) and non-linear transformations (e.g., Tosz, Tasy), to create a distribution of functions while retaining similar properties. According to the properties of functions, these
 synthetic functions can be divided into 5 categories, i.e., (1) separable functions, (2) moderately conditioned functions, (3) ill-conditioned and unimodal functions, (4) multi-modal functions with adequate global structure, and (5) multi-modal functions with weak global structures. We select one function from each category to evaluate our algorithm, Rastrigin

827

828

829

830

831

832

833

834

835

836

837

838 839

840

841

842

843

844

845

846

847

848

849

850 851

852

853

855

858

859

861

862

863

811	Table 1: List of hyper-parameter settings in RIBBO.		
812			
813	KIBBO		
814	Embedding dimension	256	
815	Number of self-attention layers	12	
816	Number of self-attention heads	8	
817	Point-wise feed-forward dimension	1024	
017	Dropout rate	0.1	
010	Batch size	64	
819	Learning rate	0.0002	
820	Learning rate decay	0.01	
821	Optimizer	Adam	
822	Optimizer scheduler	Linear warm up and cosine annealing	
823	Number of training steps	500,000	
824	Length of subsequence τ	50	
005			

Table 1. I ist af 1

from (1), Rosenbrock Rotated from (2), Sharp Ridge from (3), Greiwank Rosenbrock from (4), and Lunacek from (5). We use the BBOB benchmark implementation in Open Source Vizier¹, and the dimension is set to 10 for all functions.

- SVM and XGBoost. HPO-B Arango et al. (2021) is the most commonly used HPO benchmark and is grouped by search space id. Each search space id corresponds to a machine learning model, e.g., SVM or XGBoost. Each such search space has multiple associated dataset id, which is a particular HPO problem, i.e., optimizing the performance of the corresponding model on a dataset. For the continuous domain, it fits an XGBoost model as the objective function for each HPO problem. These datasets for each search space id are divided into training and test datasets. We examine our method on two selected search space id, i.e., 5527 and 6767, which are two representative HPO problems tuning SVM and XGBoost, respectively. SVM has 3 parameters, while XGBoost has 18 parameters, which is the most in HPO-B. We use the official implementations².
- Rover Trajectory Planning Eriksson et al. (2019); Wang et al. (2018) is a trajectory optimization task designed to emulate a rover navigation task. The trajectory is determined by fitting a B-spline to 30 points in a 2D plane, resulting in a total of 60 parameters to optimize. Given $x \in [0,1]^{60}$, the objective function is $f(x) = c(x) + \lambda (||x_{0,1} - s||_1 + ||x_{0,1} - s||_1 + ||x_{0,1} - s||_1)$ $\|\boldsymbol{x}_{58,59} - \boldsymbol{g}\|_1 + b$, where $c(\boldsymbol{x})$ is the cost of the given trajectory, $\boldsymbol{x}_{0,1}$ denotes the first and second dimensions of x (similarly for $x_{58,59}$), s and g are 2D points that specify the starting and goal positions in the plane, λ and b are parameters to define the problem. This problem is non-smooth, discontinuous, and concave over the first two and last two dimensions. To construct the distribution of functions, we applied translations in $[-0.1, 0.1]^d$ and scalings in [0.9, 1.1], similar to previous works Elhara et al. (2019); Volpp et al. (2020). The training and test tasks are randomly sampled from the distribution. We use the standard implementation for the rover problem³.

B.2 BAHAVIOR ALGORITHMS

The datasets are generated with several representatives of heuristic search, EA, and BO as behavior algorithms. We use the implementation in Open Source Vizier Song et al. (2022) for Random Search, 854 Shuffled Grid Search, Eagle Strategy, CMA-ES. For Hill Climbing and Regularized Evolution, we provide a simple re-implementation. We use the implementation in BoTorch Balandat et al. (2020) 856 for GP-EI. The details of these behavior algorithms are summarized as follows.

- Random Search selects a point uniformly at random from the domain at each iteration.
- Shuffled Grid Search discretizes the ranges of real parameters into 100 equidistant points and selects a random point from the grid without replacement at each iteration.

¹https://github.com/google/vizier

²https://github.com/releaunifreiburg/HPO-B

³https://github.com/uber-research/TuRBO

• Hill Climbing. At each iteration t, the current best solution x_{best} is mutated (using the same 865 mutation operation as Regularized Evolution) to generate x_{next} , which is then evaluated. 866 If $f(x_{next}) > f(x_{best})$, x_{best} is updated to x_{next} . 867 • Regularized Evolution Real et al. (2019) is an evolutionary algorithm with tournament 868 selection and age-based replacement. We use a population size of 25 and a tournament size of 5. At each iteration, a tournament subset is randomly selected from the current 870 population, and the solution with the maximum value is mutated. The mutation operation 871 uniformly selects one of the parameters and mutates it to a random value within the domain. 872 • Eagle Strategy Yang & Deb (2010) without the Levy random walk, aka Firefly Algorithm, 873 maintains a population of fireflies. Each firefly emits light whose intensity corresponds to 874 the objective value. At each iteration, for each firefly, a weight is calculated to chase after 875 a brighter firefly and actively move away from darker ones in its vicinity. The position is 876 updated based on the calculated weight. 877 • CMA-ES Hansen (2016) is a popular evolutionary algorithm. At each iteration, candidate 878 solutions are sampled from a multivariate normal distribution and evaluated, then the mean 879 and covariance matrix are updated. • GP-EI employs GP as the surrogate model and EI Jones et al. (1998) as the acquisition 881 function for BO. 882 883 RIBBO is trained exclusively from offline datasets derived from various source tasks, which are sampled from the task distribution. Each dataset consists of several optimization histories, generated 885 by running a behavior algorithm on a task. No additional assumptions are imposed to the data collection process or the behavior algorithms, thereby permitting the use of any behavior algorithm. 887

Nevertheless, we can choose appropriate behavior algorithms to help the training. The model is 888 trained on the datasets collected by the behavior algorithms, and the characteristics of these datasets 889 can influence the model's efficacy. It is advisable to employ well-performing and diverse behavior 890 algorithms to create these datasets. If the datasets are from suboptimal behavior algorithms pre-891 dominantly, there might be a decline in performance, even with the incorporation of RTG tokens. 892 If the datasets are predominantly produced by behavior algorithms with homogeneous properties, 893 the model may only be adept at addressing specific problems characterized by those properties. 894 Conversely, if the behavior algorithms are diverse, the model can learn the strengths from various 895 algorithms.

B.3 DATA

For BBOB and rover problem, a set of tasks is sampled from the task distribution and the above 899 behavior algorithms are used to collect data. Specifically, for the BBOB suite, a total of 200 func-900 tions are sampled, and each behavior algorithm is executed 500 times, except for GP-EI, which is 901 limited to 100 function samples due to its high time cost. The total number of optimization histories 902 is about 100,000 and the length of each history is 150. For the rover problem, 300 functions are 903 sampled with each being run 500 times, whereas for GP-EI, a smaller number of 50 functions is 904 selected, each being run 500 times. The total number of optimization histories is about 150,000 and 905 the length of each history is 150. For testing, we randomly sample 10 functions for each problem 906 and run each algorithm multiple times to report the average performance and standard deviation. 907

For HPO-B, we use the "v3" meta-training/test splits provided by the authors, which consist of 51 training and 6 test tasks for SVM, and similarly, 52 training and 6 test tasks for XGBoost. All behavior algorithms employ 500 random seeds to collect the training datasets. The total number of optimization histories is about 25,000 and the length of each history is 100.

911 912

896 897

864

C RE-IMPLEMENTATION OF OPTFORMER

913 914

915 OptFormer Chen et al. (2022) is a general optimization framework based on transformer Vaswani 916 et al. (2017) and provides an interface to learn policy or function prediction. When provided with 917 textual representations of historical data, OptFormer can determine the next query point x_t , acting as a policy. Additionally, if the context incorporate a possible query point x_t , the framework is capable of predicting the corresponding y_t , thus serving as a prediction model. We focus on the aspect of policy learning, due to its greater relevance to our work.

The original implementation is built upon a text-based transformer and uses private datasets for training. In this paper, we have re-implemented a simplified version of OptFormer, where we omit the textual tokenization process and only retain the algorithm type within the metadata. Numerical inputs are fed into the model and we use different initial states x_0 and y_0 to distinguish among algorithms. The initial states are derived by indexing the algorithm type in an embedding layer, thereby enabling the identification of distinct behavior algorithms within OptFormer. The hyper-parameters are set as same as our method.

To examine the algorithm imitation capability, we compare our re-implementation with the corre-sponding behavior algorithms. The results are shown in Figure 4. For clarity in the visual rep-resentation, we only plot a subset of the behavior algorithms, including Shuffled Grid Search, Hill Climbing, Regularized Evolution and Eagle Strategy. These behavior algorithms are plotted by solid lines, while their OptFormer counterparts are shown by dashed lines with the same color. Note that the figure uses the immediate function values as the y-axis to facilitate a comprehensive observation of the optimization process.



Figure 4: Comparison of the behavior algorithms with the OptFormer re-implementation.

D GENERALIZATION

We conduct a series of experiments to examine the generalization of our method in this section. We train the model on 4 of 5 chosen synthetic functions and test on the remaining one. The results in Figure 5 have shown that RIBBO has a strong generalization ability to unseen function distributions.



Figure 5: Cross-distribution generalization by training on 4 of 5 chosen synthetic functions and testing on the remaining one.

We also train the model across all training datasets from BBOB synthetic functions simultaneously, and normalize all the test functions to aggregate results across functions with different output scaling. The results are shown in Figure 6(a), and RIBBO has the best performance.

The generalization across different problems is implemented in an in-context manner. The con-text data, collected from the new problems, provides insights for understanding the problems at hand. These contexts are integrated to construct the inputs, thereby influencing the resulting sam-pled points.

0.95

Normalized value 0.80 0.70 0.70 0.70 0.70 0.70

0.65

0.60

0

975 976 977

972

973

974

978 979

980 981

982

983 984

985 986

E FURTHER DISCUSSIONS ABOUT RTG TOKENS

Number of evaluations

std of x

50

(a) Aggregation on the entire BBOB

Number of evaluations

BBOB

We conduct more discussion about the setting of RTG tokens in this section. The immediate RTG 987 tokens are designed to represent the performance based on cumulative regret over the future trajec-988 tories (i.e., $R_i = \sum_{t'=i+1}^{t} (y^* - y_{t'})$), rather than focusing solely on instantaneous regret $y^* - y_{t+1}$. 989 Thus, setting the immediate RTG to 0 inherently accounts for future regret, thereby enabling the 990 algorithm to autonomously trade-off the exploration and exploitation. We conduct experiments to 991 examine the performance of RIBBO with different values of R_t . The results are shown in Fig-992 ure 6(b). Setting to a value larger than 0 will decrease the performance and converge to a worse 993 value, validating the effectiveness of setting the immediate RTG as 0. 994

150

100

Normalized value 86'0 e 0'32

0.96

0

GriewankRosenbrock

HRR $(R_t=0)$

HRR $(R_t=1)$

HRR $(R_t=5)$

150

100

Number of evaluations

50

(b) HRR with different R_t

90

0.11

120

0.03

150

0.01

Sampling the next query point x_t is influenced by all the model inputs, including the immediate 995 RTG token x_t and the previous history $(x_i, y_i, R_i)_{i=0}^{t-1}$. The immediate RTG token represents the 996 goal that is intended to be achieved, while the historical data contains information about the problem. 997 These elements are integrated to construct the inputs, influencing the resulting sampled point. Even 998 if the RTG token is set to 0, implying a desire to sample the optimum, the short length of the history 999 suggests limited knowledge about the problem, prompting the model to retain explorative behavior. 1000 As the history extends, suggesting a sufficient knowledge of the problem, the preference encoded 1001 by the RTG token shifts towards more greedy action. We conduct an experiment about the std of 1002 the output Gaussian head of x of RIBBO using GriewankRosenbrock function and the results are 1003 shown in Table 2. It is clear that the std is very large in the early stage and becomes small later, which demonstrates the ability to trade-off the exploration and exploitation. 1004

1

0.37

30

0.17

Table 2: Std of the output Gaussian head of x.

60

0.14

1005

1008

1009

1010 1011 1012

F ABLATION STUDIES

1013 1014 We provide further ablation studies to examine the influence of the key components and hyper-1015 parameters in RIBBO, including the method to aggregate (x_i, y_i, R_i) tokens, the normalization 1016 method for the function value y, the model size, and the sampled subsequence length τ during 1017 training.

1018 **Token Aggregation** aims to aggregate the information from (x_i, y_i, R_i) tokens and establish asso-1019 ciations between them. RIBBO employs the Concatenation method (Concat), i.e., concatenating to 1020 aggregate x_i, y_i and R_i to form a single token. This method is compared with two alternative token 1021 aggregation methods, i.e., Addition (Add) and Interleaving (Interleave). The addition method integrates the values of each token into one, while the interleaving method addresses tokens sequentially. 1023 The results are shown in Figure 6(c). The concatenation method surpasses both the addition and interleaving methods, with the addition method showing the least efficacy. The concatenation method 1024 employs a relatively straightforward approach to aggregate tokens, while the interleaving method 1025 adopts a more complex process to learn the interrelations among tokens. The inferior performance of the addition method can be due to the summation operator, which potentially eliminates critical details from the original tokens. These details are essential for generating the subsequent query x_t .

Normalization Method is to balance the scales of the function value y across tasks. We compare 1029 the employed random normalization with the dataset normalization and the absence of any normal-1030 ization. Dataset normalization adjusts the value of y by $(y - y_{\min}^i)/(y_{\max}^i - y_{\min}^i)$ for each task i, 1031 where y_{\min}^i and y_{\max}^i denote the observed minimum and maximum values of f_i . No normalization 1032 uses the original function values directly. The results in Figure 6(d) indicate that random normaliza-1033 tion and dataset normalization perform similarly, whereas the absence of normalization significantly 1034 hurts the performance. The choice of normalization method directly impacts the computation of 1035 RTG tokens and subsequently affects the generation of the desired regret. The lack of normaliza-1036 tion leads to significant variations in the scale of y, which complicates the training process. Both random and dataset normalization scale the value of y within a reasonable range, thus facilitating 1037 both training and inference. However, random normalization can bring additional benefits, such as 1038 invariance across various scales of y as mentioned in Wistuba & Grabocka (2021) and Chen et al. 1039 (2022). Therefore, we recommend using random normalization in practice. 1040

1041 Model Size has an effect on the in-context learning ability. We assess the effects of different model 1042 sizes by comparing the performance of the currently employed model size with both a smaller and 1043 a larger model. The smaller model consists of 8 layers, 4 attention heads, and 128-dimensional embedding space, while the larger model has 16 layers, 12 attention heads, and 384-dimensional 1044 embedding space. The results are shown in Figure 6(e), indicating that the model size has a minimal 1045 impact on overall performance. Specifically, the smaller model, due to its limited capacity, shows a 1046 reduced performance, while the larger one, potentially more powerful, requires more training data, 1047 which can lead to a slight decrease in performance due to overfitting or inefficiency in learning from 1048 a limited dataset. This analysis highlights the trade-offs involved in selecting the appropriate model 1049 size for optimal performance. 1050

Subsequence Length τ controls the context length during both the training and inference phases. 1051 The process of subsequence sampling acts as a form of data augmentation, enhancing training ef-1052 ficiency. As shown in Figure 6(f), sampling subsequences rather than using the entire history as 1053 context, particularly when $\tau = T = 150$, leads to improved performance. The computational com-1054 plexity for causal transformer training and inference scales quadratically with the context length. 1055 Therefore, utilizing a shorter τ can significantly reduce computational demands. However, it is im-1056 portant to note that a shorter τ might not capture sufficient historical data, potentially degrading the 1057 performance due to the insufficient contextual information. This highlights the trade-offs between 1058 computational complexity and performance.



1072

1073

1059

Figure 6: Ablation studies of token aggregation, normalization method, model size, and subsequence length τ .

G VISUALIZATION OF BRANIN FUNCTION

In several of our experiments, such as BBOB, rover problems, and the visualization analysis of the Branin function, we implement a series of transformations to the search space. This is designed to generate a distribution of functions with similar properties. A set of functions is sampled serving as the training and test tasks from the distribution.

1079 For the Branin function, random translations and scalings are applied to form the distribution. In this section, we present visualizations of the contour lines of the 2D Branin functions, sampled from the

1080 distribution to demonstrate the effects of the applied transformations. To this end, 4 distinct random 1081 seeds are used to draw samples from the distribution. The visualizations are shown in Figure 7. 1082 Note that the two parameters of the Branin function have been scaled to the range of [-1, 1] for the 1083 clarity of visual representations. Variations are observed in both the location of the optimum and the 1084 scales of the objective values. For instance, the optimum of the first subfigure locates more to the right compared to the second one, and the contour lines in the former are much more dispersed than those in the latter. The value scale for the first subfigure ranges from -8 to 1, while in the second, it 1086 ranges from -3.5 to 1. By applying these transformations, we can generate a function distribution 1087 that retains similar properties, enabling the sampling of training and test tasks for our model. 1088

For other functions, such as BBOB suite, besides simple random translations and scalings, more complex transformations such as non-linear Tosz and Tasy transformations are applied, leading to a more intricate landscape. However, due to the high dimensionality of these functions, direct visualization is impractical, so we present only the visualization of the 2D Branin function.



Figure 7: Contour line visualization for samples drawn from the 2D Branin function distribution using 4 distinct random seeds.

H COMPARISON TO BBO ALGORITHM SELECTION METHODS AND META-LEARNING INDIVIDUAL COMPONENTS



Figure 8: Comparison to BBO algorithm selection methods and meta-learning individual components.

1119 1120

1116

1117

1118

1101

1102

1107

1108

- 1121
- 1122 1123
- 1124
- 1125

Ι

J FURTHER DISCUSSIONS ABOUT DT

ADDITIONAL RESULTS ON BBOB FUNCTIONS

1126 DT represents a paradigm of upside down RL Schmidhuber (2019), that maps the desired rewards 1127 to corresponding actions to enhance the performance. The benefit of using regret-to-go is its inde-1128 pendence from a predetermined horizon length when using the HRR strategy, contrasting with the 1129 return-to-go method where the calculation is usually related to the product of the maximal achiev-1130 able reward and the horizon length. Once established, the return-to-go is expected to decrease as 1131 the optimization progresses. Additionally, from the perspective of optimization, the cumulative regret provides a more meaningful measure of the performance (measuring the gap to the optimum), 1132 compared to the sum of y. Therefore, regret-to-go has been adopted as the preferred metric in the 1133 algorithm's design.



1188 1189	thereby enhancing its generalization. A mathematical theoretical analysis about the cumulative re-
1100	gret analysis and generalization analysis based on RTG tokens is of interest and neipitul for relining
1191	the algorithmic designs.
1192	
1193	
1194	
1195	
1196	
1197	
1198	
1199	
1200	
1201	
1202	
1203	
1204	
1205	
1206	
1207	
1208	
1209	
1210	
1211	
1212	
1213	
1214	
1215	
1216	
1217	
1218	
1219	
1220	
1221	
1222	
1223	
1224	
1225	
1220	
1227	
1220	
1223	
1231	
1232	
1233	
1234	
1235	
1236	
1237	
1238	
1239	
1240	
1241	