

# Pre-training helps Bayesian optimization too

<b>Zi Wang</b>	WANGZI@GOOGLE.COM
<b>George E. Dahl</b>	GDAHL@GOOGLE.COM
<b>Kevin Swersky</b>	KSWERSKY@GOOGLE.COM
<b>Chansoo Lee</b>	CHANSOO@GOOGLE.COM
<b>Zelda Mariet</b>	ZMARIET@GOOGLE.COM
<b>Zachary Nado</b>	ZNADO@GOOGLE.COM
<b>Justin Gilmer</b>	GILMER@GOOGLE.COM
<b>Jasper Snoek</b>	JSNOEK@GOOGLE.COM
<b>Zoubin Ghahramani</b>	ZOUBIN@GOOGLE.COM

## Abstract

Bayesian optimization (BO) has become a popular strategy for global optimization of many expensive real-world functions. Contrary to a common belief that BO is suited to optimizing black-box functions, it actually requires domain knowledge on characteristics of those functions to deploy BO successfully. Such domain knowledge often manifests in Gaussian process priors that specify initial beliefs on functions. However, even with expert knowledge, it is not an easy task to select a prior. This is especially true for hyperparameter tuning problems on complex machine learning models, where landscapes of tuning objectives are often difficult to comprehend. We seek an alternative practice for setting these functional priors. In particular, we consider the scenario where we have data from similar functions that allow us to pre-train a tighter distribution a priori. Theoretically, we show a bounded regret of BO with pre-trained priors. To verify our approach in realistic model training setups, we collected a large multi-task hyperparameter tuning dataset by training tens of thousands of configurations of near-state-of-the-art models on popular image and text datasets, as well as a protein sequence dataset. Our results show that on average, our method is able to locate good hyperparameters at least 3 times more efficiently than the best competing methods.

## 1. Introduction

Bayesian optimization (BO) has been successfully applied in numerous real-world global optimization problems, ranging broadly from hyperparameter tuning (Snoek et al., 2012; Kotthoff et al., 2019) to chemical synthesis (Shields et al., 2021; Griffiths and Hernández-Lobato, 2020), drug discovery (Pyzer-Knapp, 2018), aerospace engineering (Lam et al., 2018), robotics (Dri   et al., 2017; Wang et al., 2017a) and the list goes on. However, in some scenarios, Bayesian optimization has been reported to under-perform naive strategies including random search (Li et al., 2017). While recent collective efforts have shown that "Bayesian optimization is superior to random search" (Turner et al., 2021), we seek more understandings on why BO works in some hands but not others.

Many successful BO applications benefit from expert knowledge on characteristics of the function to be optimized and hands-on experience with BO on similar tasks in the past. Such knowledge or experience can give intuitions about a functional form of the problem and thus specifications of

a functional prior, e.g. a Gaussian process (GP) with squared exponential kernels for smoothness. Sometimes people may be uncertain about their own understanding, and as a result they might choose to use a hierarchical model (Cowen-Rivers et al., 2020) or Bayesian neural nets (Springenberg et al., 2016), such that observed data can play a more important role in modeling. Despite having almost no information about a function, we can guess a generic prior from past experience with BO on other functions (Turner et al., 2021). But, what if we have neither domain knowledge nor hands-on experience to set an informative prior? In this case, we often have no alternative but a misspecified prior. However, we found very little theoretical or empirical evidence on a positive correlation between good performance of BO and priors that mismatch the distribution of the optimization target function. Our experiments in §5 further verified that model misspecification can be detrimental to off-the-shelf BO with type-II maximum likelihood.

For BO to gain more impact on complex real-world problems, it is important to continue developing more convenient and accessible BO methods. Barriers of understanding on priors from a target domain and enough experience with BO can often turn away potential practitioners even within the machine learning (ML) community (Bouthillier and Varoquaux, 2020). Similar to how we as engineers and researchers learn how to set good priors from past experience, we seek to automate the prior determination process by *pre-training priors* on data that are available on different but related tasks. Our prior pre-training approach is also known as prior learning or a version of meta learning.

Given the benefit of pre-trained priors on synthetic functions, simple tuning tasks and complex robotics tasks (Wang et al., 2018b; Kim et al., 2019; Perrone et al., 2018), can we take it to the level of real-world hyperparameter tuning problems for modern deep learning models (e.g. ResNet50) and large-scale datasets (e.g. ImageNet)? To the best of our knowledge, there is no such multi-task tuning benchmark available for modern large models and datasets, but this kind of tuning task is most prevalent in hyperparameter tuning in recent years of ML-related publications given the success of large models (He et al., 2016; Raffel et al., 2019; Brown et al., 2020). For these problems, it is difficult to understand the landscapes of tuning objectives, hindering the use of Bayesian optimization with almost unobtainable expert interventions on priors.

To fill the vacancy of a dataset for hyperparameter tuning in modern ML, we collected a large multi-task hyperparameter tuning dataset by training tens of thousands of configurations of near-state-of-the-art models on popular image and text datasets, as well as on a protein sequence dataset. Our open-sourced dataset can save roughly 12,000 machine-days of computation for anyone who makes use of it.

On the modeling side, most existing meta BO methods either scale cubically in the number of evaluations and tasks (Swersky et al., 2013; Bardenet et al., 2013), impose a restrictive set of assumptions on the available data (Wang et al., 2018b; Swersky et al., 2013) for efficient solutions, or make assumptions on the availability of GP parameters (Volpp et al., 2020) or descriptive task-level features (Brazdil et al., 1994; Bardenet et al., 2013; Yogatama and Mann, 2014). To address these issues, we introduce HyperBO: a meta BO method that builds upon Wang et al. (2018b) with a simple assumption: all the related functions being optimized are samples from the same GP prior distribution. Concretely, HyperBO assumes the functions are conditionally independent given the hyperparameters, mean and covariance function of the GP. Compared to Wang et al. (2018b), HyperBO does not impose any strict conditions on data or model structures, and a special case of HyperBO retains strong regret bounds. From a computational perspective, HyperBO scales linearly in the number of tasks during training, and does not depend on the number of tasks when deployed.

By not imposing assumptions about the data collection conditions, it can be used with large offline datasets or a few related optimization trajectories.

Our empirical results show that HyperBO is at least 3 times more efficient in function evaluations than recent baseline methods to locate the best hyperparameters. Our main contributions are two-fold: (1) a new prior pre-training approach for BO that makes minimal assumptions while retaining theoretical guarantees; (2) a large multi-task hyperparameter tuning dataset that not only benefits our method but also serves as a realistic benchmark to test future methods. Both open-sourced code and dataset are available at <https://github.com/google-research/hyperbo>.

## 2. Related work

There is a rich literature of innovative methodologies to improve the efficiency of BO given related tasks or additional context. Here we discuss the most closely related work and explain why these don’t solve the specific scenario which we envision. Specifically, our goal is a methodology that is scalable enough to share information across thousands of tasks, each with potentially hundreds of observations, such as in the context of a large BO service or library.

Pre-training and prior learning is directly related to meta learning, learning to learn and learning multiple tasks (Baxter, 1996). We use the word pre-training to refer to supervised pre-training, which is a general approach in the deep learning community (Girshick et al., 2014) to transfer knowledge from prior tasks to a target task. The same as pre-training deep features on a variety of tasks, Wang et al. (2018b) proposed prior learning for GPs to learn the basis functions by treating the independent function outputs as individual heads of a neural network.

Several methods, including that which HyperBO extends, refer to their method as “meta-BO” (Wang et al., 2018b; Volpp et al., 2020). However, in this work we use the term *meta-BO* more generally to refer to the class of BO methods that use data from existing tasks to optimize a new task. Since standard BO is a learning process, it is consistent to call those methods meta BO methods given that they learn how to learn. Under this viewpoint, meta BO approaches also include multi-task BO (Swersky et al., 2013; Poloczek et al., 2017; Yogatama and Mann, 2014), transfer learning BO using contextual GPs (Krause and Ong, 2011; Bardenet et al., 2013; Poloczek et al., 2016) and transfer learning based on quantiles (Salinas et al., 2020). Some meta BO methods have also been studied for hyperparameter tuning tasks in machine learning (Feurer et al., 2015; Salinas et al., 2020).

HyperBO assumes all tasks are independent (after conditioning on the GP), whereas both multi-task and contextual BO rely heavily on the assumption that tasks are related. Thus the latter approaches typically scale cubically in both the number of tasks and observations in each task, meaning that they cannot gracefully scale across both without heavy approximations. When assuming that all inputs are equal across tasks, multi-task BO can be sped up using a Kronecker decomposition of the kernel to a task kernel and an input kernel which can be inverted separately; a similar assumption is made by Wang et al. (2018b). In comparison, HyperBO scales linearly in the number of tasks (see §4.3).

Another thread of meta BO literature was started in the robot learning area by Kim et al. (2017, 2019), which estimated a multivariate Gaussian to transfer knowledge on scoring functions for search strategies in robot manipulation tasks, and thus only considered finite discrete inputs. Wang et al. (2018b) provided regret bounds for Kim et al. (2017, 2019) and extended it to continuous search spaces by considering a GP as a Bayesian linear regressor with neural net basis functions.

Similar ideas were adopted by Perrone et al. (2018); Wistuba and Grabocka (2021) in the machine learning hyperparameter tuning literature. These ideas, on a high level, can be viewed as special cases of HyperBO or Wang et al. (2018b). Compared to Wang et al. (2018b), Perrone et al. (2018) and Wistuba and Grabocka (2021) only use zero means, which, as shown by Kim et al. (2017, 2019), is critical for learning the initial data points to acquire. As a remedy, Wistuba and Grabocka (2021) developed an evolutionary algorithm based data-driven strategy to warm start the initialization of data selection. Although different terms are used, Wang et al. (2018b) and Perrone et al. (2018) concurrently proposed the idea of learning parameters of GP priors from multi-task datasets, while Wang et al. (2018b) is the first to clarify the assumptions that those multi-task functions need to be conditionally independent so that regret bounds hold for BO with an unknown GP prior.

HyperBO builds upon Wang et al. (2018b) and Kim et al. (2017, 2019), yet principally resolves their limitations on search spaces and data availability. For both finite discrete search spaces and continuous ones, Wang et al. (2018b) requires observations on the same set of inputs across tasks, which is an assumption that is not required for HyperBO; HyperBO still inherits the same regret bound as Wang et al. (2018b) for the special case where the same-inputs assumption is satisfied. Both arbitrary data points from different tasks and observations on same inputs across tasks can be incorporated into HyperBO efficiently and effectively. Another critical advantage of HyperBO is accommodations of very flexible kernels and mean functions that are not limited by Bayesian linear regressors; this opens meta BO to a lot more GP architectures involving combinations of deep features and kernels with infinite basis functions.

### 3. Problem formulation

We consider the standard black-box function optimization scenario: given a real-valued function  $f$  defined over a compact, hyper-rectangular space  $\mathcal{X} \subset \mathbb{R}^d$  and given observations of similar functions  $f_1, \dots, f_N$ , we seek an  $x \in \mathcal{X}$  optimizing  $f$ . We inherit our problem formulation from Wang et al. (2018b), but we relax impractical assumptions on data availability (we do not require all observations to be made on the same inputs across tasks) and model restrictions.

**Assumptions and the goal.** Concretely, we assume that there exists a Gaussian process  $\mathcal{GP}(\mu, k)$  with unknown mean function  $\mu : \mathcal{X} \rightarrow \mathbb{R}$  and kernel  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ . Let  $N$  be the number of tasks and let  $M_i$  be the number of observations we have for the  $i$ th task. Conditioned on independent function samples  $f_i \sim \mathcal{GP}(\mu, k)$  and inputs  $x_j^{(i)} \in \mathcal{X}, i \in [N], j \in [M_i]$ , we observe evaluations  $y_j^{(i)} \sim \mathcal{N}(f_i(x_j^{(i)}), \sigma^2)$  perturbed by *i.i.d.* additive Gaussian noise  $\mathcal{N}(0, \sigma^2)$  with unknown variance  $\sigma^2$ . Taken together, the collection of sub-datasets  $D_{f_i} = \{(x_j^{(i)}, y_j^{(i)})\}_{j=1}^{M_i}$  define a dataset  $D_N = \{D_{f_i}\}_{i=1}^N$ . Finally, our goal is to maximize a new function independently sampled from the same GP,  $f \sim \mathcal{GP}(\mu, k)$ ; that is, solve  $\arg \max_{x \in \mathcal{X}} f(x)$  given dataset  $D_N$  but unknown functions  $\mu, k$  and unknown parameter  $\sigma^2$ .

**An example.** In our optimizer hyperparameter tuning application, a task corresponds to finding the best optimizer hyperparameters to train a given model on a particular dataset,<sup>1</sup> e.g. training a ResNet (He et al., 2016) on ImageNet (Russakovsky et al., 2015). Notice that we do not assume that the mean function  $\mu$ , kernel  $k$  and noise variance  $\sigma^2$  are given. This is consistent with the reality of

1. Technically, we also consider different batch sizes to be different tasks.

solving real-world black-box optimization problems including hyperparameter tuning. We must learn those unknown functions and parameters from data. However, in practice, searching in functional spaces to find the right mean  $\mu$  or kernel  $k$  is a daunting task. Hence for practical concerns, a well defined search space for functions is required. More details on this can be found at §4.2.

**Metrics.** For simplicity, throughout this paper, we focus on the setting where the target function  $f$  can only be optimized by iteratively choosing where to evaluate, and defer batch evaluation setups to Sec. 6. As we run BO on the target function  $f$  for  $T$  iterations, we accumulate a set of observations  $D_f = \{(x_t, y_t)\}_{t=1}^T$ ,  $y_t \sim \mathcal{N}(f(x_t), \sigma^2)$ . We evaluate the quality of the optimization using the *simple regret* metric:  $R_T = \max_{x \in \mathcal{X}} f(x) - f(\hat{x})$ , where  $\hat{x}$  is the final recommendation at the end of the optimization process. There are various ways of setting  $\hat{x}$  based on the observations  $D_f$ ; we use the input that achieved the best evaluation:  $\hat{x} = x_\tau$ ;  $\tau = \arg \max_{t \in [T]} y_t$ .

**Bayesian viewpoint.** As mentioned above, the observed functions  $f_1, \dots, f_N$  and the evaluation target  $f$  are assumed to be independent draws from the same GP. This assumption is consistent with a hierarchical Bayes interpretation (Fig. 1), where all observed functions are independent conditioned on the GP. Notice that for BO, each selected input  $x_j^{(i)}$  depends on all previous observations. But we only describe the generative model of a hierarchical GP for simplicity.

More specifically, we assume that the overall setting of the hyperparameter optimization task is defined by a parameter  $\theta \sim p(\theta; \alpha)$ ; mean and kernel functions  $\mu$  and  $k$  are drawn from  $p(\mu, k \mid \theta)$ . The independent function samples  $\{f_i\}_{i \in [N]}$  are themselves draws from  $\mathcal{GP}(\mu, k)$ . The generative story is as follows:

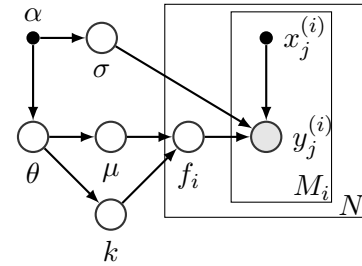
- Draw GP parameter  $\theta$  from  $p(\theta; \alpha)$  and observation noise parameter  $\sigma$  from  $p(\sigma; \alpha)$ .
- Draw mean function  $\mu$  and kernel function  $k$  from  $p(\mu, k \mid \theta)$ .
- For each task  $i$  from 1 to  $N$ ,
  - Draw a function  $f_i$  from  $\mathcal{GP}(\mu, k)$ .
  - For each data point  $j$  from 1 to  $M_i$ ,
    - \* Given input  $x_j^{(i)}$ , we draw the observation  $y_j^{(i)} \sim \mathcal{N}(f_i(x_j^{(i)}), \sigma^2)$ .

We simplify this hierarchical setting by defining  $p(\mu, k \mid \theta)$  to be a sum of Dirac delta functions: both mean function  $\mu$  and kernel  $k$  are deterministic functions parameterized by  $\theta$ . Thus, we can infer GP parameter  $\theta$  and noise  $\sigma$  from their posterior  $p(\theta, \sigma \mid D_N \cup D_f; \alpha)$  and obtain an informed prediction for the target function

$$\begin{aligned} p(f \mid D_N \cup D_f) &= \int_{\theta} p(f \mid \theta) p(\theta \mid D_N \cup D_f; \alpha) \\ &= \int_{\theta} p(f \mid \theta) \int_{\sigma} p(\theta, \sigma \mid D_N \cup D_f; \alpha) \end{aligned}$$

In other words, we learn function  $f$  from observations on all other conditionally *i.i.d.* function samples  $f_1, \dots, f_N$ . We forgo a fully Bayesian approach that samples from the posterior

Figure 1: Graphical model for a hierarchical Gaussian process.



over  $\theta$  at every BO iteration, although our method, HyperBO, can be viewed as a type-II maximum likelihood approximation of such a Bayesian solution.

**Notations.** Let  $[n]$  denote  $\{1, \dots, n\}$ ,  $\forall n \in \mathbb{Z}^+$ . For conciseness, we write the evaluation of a function  $f : \mathbb{R} \rightarrow \mathbb{R}$  on vector  $\mathbf{x} = [x_i]_{i=1}^n$  as  $\mu(\mathbf{x}) := [\mu(x_i)]_{i=1}^n$ . Similarly, for two vectors  $\mathbf{x}, \mathbf{x}'$ , we write the corresponding kernel matrix as  $k(\mathbf{x}, \mathbf{x}') := [k(x_i, x'_j)]_{i \in [n], j \in [n']}$ , and shorten  $k(\mathbf{x}) := k(\mathbf{x}, \mathbf{x})$ .

We denote a (multivariate) Gaussian distribution with mean  $u$  and variance  $\Sigma$  by  $\mathcal{N}(u, \Sigma)$ , and a Gaussian process (GP) with mean function  $\mu$  and covariance function  $k$  by  $\mathcal{GP}(\mu, k)$ . Let  $\sigma^2$  be the noise variance in observations. Given a set of observations  $D = \{(x_t, y_t)\}_{t=1}^T$ ,  $\mathbf{y}_T = [y_t]_{t=1}^T \sim \mathcal{N}(f(\mathbf{x}_T), \sigma^2 \mathbf{I})$ ,  $\mathbf{x}_T = [x_t]_{t=1}^T$  and  $f \sim \mathcal{GP}(\mu, k)$ , we denote the corresponding conditional GP distribution as  $\mathcal{GP}(\mu, k \mid D)$ . Recall that the conditional distribution  $\mathcal{GP}(\mu, k \mid D) = \mathcal{GP}(\mu_D, k_D)$ , is given for any  $x, x' \in \mathcal{X}$  as

$$\mu_D(x) = \mu(x) + \psi(x)(\mathbf{y}_T - \mu(\mathbf{x}_T)), \quad (1)$$

$$k_D(x, x') = k(x, x') - \psi(x)k(\mathbf{x}_T, x'), \quad (2)$$

where we set  $\psi(x) = k(x, \mathbf{x}_T)(k(\mathbf{x}_T) + \sigma^2 \mathbf{I})^{-1}$ .

## 4. Our method

As shown in Alg. 1, our approach pre-trains the GP hyperparameters on a representative set of datasets and fixes them for the duration of the optimization procedure; we refer to this approach as HyperBO. HyperBO runs in two steps. First, we learn a GP model  $\mathcal{GP}(\hat{\mu}, \hat{k})$  to approximate the ground-truth (unknown) GP that generated the dataset  $D_N$ . Then, we do standard BO to optimize a new function  $f$  with the learned GP  $\mathcal{GP}(\hat{\mu}, \hat{k})$ . The initial pre-training process (Alg. 1, line 2) is the critical difference between HyperBO and standard BO algorithms, as well as the key contribution of this paper.

Based on the Bayesian graphical model interpretation (Fig. 1), our goal is to obtain a point estimate  $\hat{\theta}$  for the parameter  $\theta$ . Given this estimate, we can then estimate the mean function  $\hat{\mu}$  and the kernel  $\hat{k}$ , which defines our learned model  $\mathcal{GP}(\hat{\mu}, \hat{k})$ . During the BO iterations (Alg. 1, lines 4-8), we update the conditional GP, but do not re-estimate the GP mean and kernel. By separating the data for conditional GP update and GP parameter pre-training, we minimize the computational cost while still maintaining good performance both theoretically and empirically. Moreover, we avoid the BO chicken-and-egg dilemma (Wang et al., 2018b) where the search strategy is trained on data collected in the BO process and the data points are selected by the search strategy simultaneously.

Next, we introduce our GP pre-training strategy based on two types of objectives: KL divergence between estimates and model predictions (§ 4.1) and negative log likelihood (§ 4.2). §4.6 details

---

### Algorithm 1 HyperBO with acquisition function $\alpha(\cdot)$ .

---

```

1: function HYPERBO( $f, D_N$ )
2:    $\mathcal{GP}(\hat{\mu}, \hat{k}) \leftarrow \text{PRE-TRAIN}(D_N)$ 
3:    $D_f \leftarrow \emptyset$ 
4:   for  $t = 1, \dots, T$  do
5:      $x_t \leftarrow \arg \max_{x \in \mathcal{X}} \alpha(x; \mathcal{GP}(\hat{\mu}, \hat{k} \mid D_f))$ 
6:      $y_t \leftarrow \text{OBSERVE}(f(x_t))$ 
7:      $D_f \leftarrow D_f \cup \{(x_t, y_t)\}$ 
8:   end for
9:   return  $D_f$ 
10: end function

```

---



the difference between these two objectives and our recommendations on when to use which. §4.3 reveals the complexity of HyperBO that is linear in the number of tasks and § 4.5 shows that a special case of HyperBO retains strong regret bounds (Wang et al., 2018b).

#### 4.1 Pre-training with empirical KL divergence

We first investigate the case where observations on the same set of inputs are available. This is the main scenario considered by Wang et al. (2018b), but their method only works for Bayesian linear regression. We now present HyperBO based on an empirical KL divergence to allow highest flexibility on mean functions and kernels, e.g. a Matérn kernel on deep features shared with a mean function. The objective is called *empirical KL divergence* because it is the KL divergence between an empirically estimated multivariate Gaussian and model predictions from a GP.

Here we consider a special case of dataset  $D_N$  which contains matching inputs across some tasks. More formally, suppose we have a *matching dataset*  $D'_N = \{(x_j, \mathbf{y}_j)\}_{j=1}^M$  where  $M$  is the number of shared inputs across  $N$  tasks. For each input index  $j \in [M]$  and input  $x_j \in \mathcal{X}$ , we have  $N$  observed values  $\mathbf{y}_j = [y_j^{(i)}]_{i=1}^N \in \mathbb{R}^N$  and each observation  $y_j^{(i)} \sim \mathcal{N}(f(x_j), \sigma^2)$  corresponds to evaluating input  $x_j$  on a different task. In practice, dataset  $D'_N$  can be constructed by querying a set of functions  $f_1, \dots, f_N$  at the same set of input locations  $\mathbf{x} = [x_j]_{j=1}^M \in \mathbb{R}^{M \times d}$  to obtain an observation matrix  $\mathbf{y} = [\mathbf{y}_j]_{j=1}^M \in \mathbb{R}^{M \times N}$ .

By definition of a GP, the vector of all function queries  $f(\mathbf{x})$  is distributed according to a multivariate Gaussian distribution  $\mathcal{N}(\mu(\mathbf{x}), k(\mathbf{x}))$ . With our observation model, we get the distribution for observations  $\mathbf{y} \sim \mathcal{N}(\mu(\mathbf{x}), k(\mathbf{x}) + \mathbf{I}\sigma^2)$  for some unknown mean function  $\mu$  and kernel  $k$ .

However, given that we have access to all observations  $\mathbf{y}$ , we can estimate the mean on inputs  $\mathbf{x}$  as  $\tilde{\mu} = \frac{1}{N}\mathbf{y}\mathbf{1}_N \in \mathbb{R}^M$  and estimated covariance as  $\tilde{K} = \frac{1}{N}(\mathbf{y} - \tilde{\mu}\mathbf{1}_N^\top)(\mathbf{y} - \tilde{\mu}\mathbf{1}_N^\top)^\top \in \mathbb{R}^{M \times M}$ ; here  $\mathbf{1}_N$  is a column vector of size  $N$  filled with 1s. We use a biased estimate of covariance to be consistent with the maximum likelihood estimator in §4.2. But one may choose to re-scale learned kernel by  $\frac{N}{N-1}$  to be unbiased. Notice that the estimated covariance includes in diagonal terms the variance of the observation noise.

For any divergence function between the estimate  $\mathcal{N}(\tilde{\mu}, \tilde{K})$  and model prediction  $\mathcal{N}(\mu(\mathbf{x}), k(\mathbf{x}) + \mathbf{I}\sigma^2)$ , we obtain an objective to minimize,  $\mathcal{D}(\mathcal{N}(\tilde{\mu}, \tilde{K}), \mathcal{N}(\mu(\mathbf{x}), k(\mathbf{x}) + \mathbf{I}\sigma^2))$ . While there are different measures of distributional discrepancy, we adopt the KL divergence. Let  $\mu = \mu(\mathbf{x})$  and  $K = k(\mathbf{x}) + \mathbf{I}\sigma^2$ . The empirical KL divergence is defined as

$$\mathcal{D}_{\text{KL}}(\mathcal{N}(\tilde{\mu}, \tilde{K}), \mathcal{N}(\mu, K)) = \frac{1}{2} \left( \text{tr}(K^{-1}\tilde{K}) + (\mu - \tilde{\mu})^\top K^{-1}(\mu - \tilde{\mu}) + \ln \frac{|K|}{|\tilde{K}|} - M \right), \quad (3)$$

and we can estimate the mean, kernel and noise variance by minimizing  $\mathcal{D}_{\text{KL}}$ .

#### 4.2 Pre-training with negative log likelihood

If we have arbitrary data points from each task, a straightforward way to pre-train a GP is by optimizing a negative log likelihood (NLL) over parameters of the GP. The regression-based NLL objective corresponds to how supervised pre-training is done with a cross-entropy loss on deep learning models for classification tasks.

In our regression problem, we use the NLL on the given observations from multiple functions that are assumed to be independently sampled from the GP. The independence naturally results in a summation over NLLs for all observed functions, which is a key difference to the widely used type II

maximum likelihood approximation for GP inference on a single function in BO setups. The NLL loss function for our method is

$$\begin{aligned}
L(\mu, k, \sigma^2) &= -\log p(D_N \mid \mu, k, \sigma^2) \\
&= -\sum_{i=1}^N \log p(D_{f_i} \mid \mu, k, \sigma^2) \\
&= \sum_{i=1}^N \left( \frac{1}{2} \left( \mathbf{y}^{(i)} - \mu(\mathbf{x}^{(i)}) \right)^\top K^{-1} \left( \mathbf{y}^{(i)} - \mu(\mathbf{x}^{(i)}) \right) + \frac{1}{2} \log |K| + \frac{M_i}{2} \log 2\pi \right),
\end{aligned} \tag{4}$$

where  $K = k(\mathbf{x}^{(i)}) + \sigma^2 \mathbf{I}$ ,  $\mathbf{x}^{(i)} = [x_j^{(i)}]_{j=1}^{M_i}$  and  $\mathbf{y}^{(i)} = [y_j^{(i)}]_{j=1}^{M_i}$ . We then obtain a solution to the choice of mean function, kernel function and noise variance by minimizing the NLL loss function.

### 4.3 Computational complexity

To better understand the computational efficiency of HyperBO, we analyze the computational complexity of our approach. Let  $N$  be the number of sub-datasets and  $M = \max_{i=1}^N M_i$  be the maximum number of data points for these sub-datasets. For optimization, we limit our search space to a single architecture of GP mean and kernel. And we assume the constant number of variables to optimize is much smaller than  $\min(M, N)$ . Tab. 1 summarizes the following analyses.

The empirical KL divergence (Eq. 3) introduced in §4.1 requires an overhead of estimating mean and covariance, which takes  $\mathcal{O}(M^2 N)$  for matrix multiplication. The KL divergence in Eq. 3 has a time complexity of  $\mathcal{O}(M^3)$  to compute and a memory complexity of  $\mathcal{O}(M^2)$ .

The NLL objective in Eq. 4 naturally decomposes into a sum of GP data likelihood terms on each sub-dataset  $D_{f_i}$ . The time complexity to compute Eq. 4 is  $\mathcal{O}(M^3 N)$  and the memory complexity is  $\mathcal{O}(M^2)$ . If we have  $N$  processes computing each additive component of Eq. 4 in parallel, the time complexity can be reduced to  $\mathcal{O}(M^3)$  while the memory complexity becomes  $\mathcal{O}(M^2 N)$ .

Notice that our method scales (at most) linearly in the number of tasks,  $N$ , in contrast to the cubic  $\mathcal{O}(M^3 N^3)$  scaling of multi-task or contextual BO methods (Swersky et al., 2013; Bardenet et al., 2013; Poloczek et al., 2016; Yogatama and Mann, 2014). The only cubic cost of HyperBO is on the number of data points in sub-datasets.

If there is any better probabilistic model than a GP to fit the data with less compute time, we can easily swap the  $\mathcal{O}(M^3)$  cost in Eq. 3 or Eq. 4 for a more efficient one. For example, if we approximate a GP with a linear model on  $V$  random features (Rahimi et al., 2007), the time complexity becomes  $\mathcal{O}(V^3)$  for Eq. 3 and  $\mathcal{O}(V^3 N)$  for Eq. 4.

If we pre-train the model via NLL (Eq. 4) with gradient descent (GD) or stochastic gradient descent (SGD) on partitions of tasks and data points per task, the time complexity of the NLL (Eq. 4) on the full dataset can be reduced to  $\mathcal{O}(B^2 M N)$ , where  $B$  is the mini-batch size of data points per task. Here,  $B \equiv M$  for GD or SGD with partitions only on tasks. Running stochastic optimization will then take  $\mathcal{O}(B^2 M N K)$ , where  $K$  is the number of optimization epochs. The memory complexity reduces from  $\mathcal{O}(M^2)$  for the original case to  $\mathcal{O}(B^2)$  for the loss on mini-batches. Note that the larger the mini-batch size  $B$ , the better the approximation to Eq. 4. Setting  $B = 1$  is equivalent to assuming each data point is from a different task and hence it is unlikely to obtain a good pre-trained GP. In terms of complexity, SGD with partitions only on tasks has no difference to GD by computing the additive components of the NLL (Eq. 4) sequentially; however, the optimization landscape changes.



Table 1: Time and memory complexity of HyperBO.  $K$  is the number of optimization steps (or epochs in stochastic optimization).  $B$  is the mini-batch size of SGD over data points per task.

		Time	Memory
KL (Eq. 3)	Overhead	$\mathcal{O}(M^2 N)$	$\mathcal{O}(M^2)$
	Loss function	$\mathcal{O}(M^3)$	$\mathcal{O}(M^2)$
	GD	$\mathcal{O}(M^3 K)$	$\mathcal{O}(M^2)$
	SGD	$\mathcal{O}(B^2 M K)$	$\mathcal{O}(B^2)$
NLL (Eq. 4)	Loss function	$\mathcal{O}(M^3 N)$	$\mathcal{O}(M^2)$
	Parallel	$\mathcal{O}(M^3)$	$\mathcal{O}(M^2 N)$
	GD	$\mathcal{O}(M^3 N K)$	$\mathcal{O}(M^2)$
	SGD	$\mathcal{O}(B^2 M N K)$	$\mathcal{O}(B^2)$

#### 4.4 How to optimize objectives

Optimizing the KL (Eq. 3) or NLL (Eq. 4) objectives over mean function  $\mu$ , kernel  $k$  and noise variance  $\sigma^2$  requires an optimization procedure done in functional spaces for mean and kernel. While methods exist to search for functional structures (Kemp and Tenenbaum, 2008; Malkomes and Garnett, 2018), one may opt for a simple search strategy within a group of functional structures (e.g. mean  $\mu \in \{\text{linear, constant}\}$  and kernel  $k \in \{\text{exponentiated quadratic, Matérn, dot-product}\}$ ). For all combinations of mean/kernel structures or functional classes, we then optimize the parameterization of them and noise variance  $\sigma^2$  to eventually optimize the objective directly or perform cross-validation on a held-out validation dataset.

In this paper, we simplify this procedure by limiting the search space to one specific architecture while ensuring that the architecture allows flexibility. In particular, we used a neural network mean functions and Matérn kernels on the same deep features. Details of how we defined the search space can be found in §5.

Now our problem becomes how to optimize the KL (Eq. 3) or NLL (Eq. 4) objectives over fixed-dimension real variables. There are many optimization methods that can be used here. We found L-BFGS (Liu and Nocedal, 1989) to perform consistently well. Alternatively, it’s also possible to use other popular gradient based optimization methods, e.g. from Hessel et al. (2020).

#### 4.5 Theoretical analyses

Though it is nontrivial to prove regret bounds for general scenarios without strict assumptions using the NLL objective (Eq. 4), it is relatively straightforward to show a regret bound for our method with objective  $\mathcal{D}_{\text{KL}}$  of Eq. 3 in the matching dataset case where BO is running on a finite set of same inputs across all tasks.

**Proposition 1.** *For any  $M, d, N \in \mathbb{Z}^+$ ,  $\mathbf{x} \in \mathbb{R}^{M \times d}$ ,  $\boldsymbol{\mu} \in \mathbb{R}^M$ ,  $V \in \mathbb{R}^{N \times M}$  and  $K = V^\top V$ , there exists a GP  $\mathcal{GP}(\hat{\mu}, \hat{k})$  such that  $\mathcal{D}_{\text{KL}}(\mathcal{N}(\boldsymbol{\mu}, K), \mathcal{N}(\hat{\mu}(\mathbf{x}), \hat{k}(\mathbf{x}))) \equiv 0$ .*

Proposition 1 is easy to show. We can train a simple memory based model for mean function  $\hat{\mu}$  and kernel  $\hat{k}$ . The model stores each element of vector  $\boldsymbol{\mu}$  and matrix  $K$  at the corresponding locations in input  $\mathbf{x}$ . When making a prediction on a new input  $\mathbf{x}' \in \mathbb{R}^d$ , the model simply retrieves the values

of the closest element in  $\mathbf{x}$ . Given Proposition 1, the proof for regret bounds follows Wang et al. (2018b).

**Theorem 2.** *Let  $N \geq 4 \log \frac{6}{\delta} + T + 2$ . With probability at least  $1 - \delta$ , simple regret in  $T$  iterations of Alg. 1 with special cases of either GP-UCB or PI satisfies*

$$R_T < O \left( \sqrt{\frac{1}{N-T}} + \sqrt{\log \frac{1}{\delta}} \right) O(\rho_T/T + \sigma), \quad (5)$$

where  $\rho_T = \max_{A \subset \mathcal{X}, |A|=T} \frac{1}{2} \log |\mathbf{I} + \sigma^{-2} k(A)|$ .

We describe details of the proof and the special cases of GP-UCB and PI in Appendix B. Theorem 2 shows that the regret bound always has a linear dependency on the observation noise  $\sigma$ . This is expected because in practice, we select the best observation rather than best function value (before observing a noisy version of it) to compute the simple regret. Another reason is that we learn the noise parameter  $\sigma$  jointly with the kernel, which is clear in Eq. 3. Hence when computing acquisition functions, the noise  $\sigma$  is always included in the predicted variance.

Intuitively, the more sub-datasets we have in the dataset, the larger  $N$  is, the better we are able to estimate the GP model, and the closer the regret bound is to the case where the GP model is assumed known. Interestingly, the number of BO iterations  $T$  makes the regret smaller in the second term but larger in the first term in Eq. 5. Usually as we get more observations, we get more information about the maximizer, and we are able to optimize the function better. However, as we get more observations on the new function, GP conditional predictions have more freedom to deviate from the ground truth (see Lemma 1 of Wang et al. (2018b)). As a result, we get less and less confident about our predictions, which is eventually reflected in a looser regret upper bound.

It is tempting to prove similar bounds for more general settings where inputs are not the same across all sub-datasets and BO happens in a continuous space. Though the only prerequisite is to show that the difference between the learned mean/kernel and the ground truth mean/kernel is small, this prerequisite is as difficult as showing we can find a model that has bounded generalization error across the entire continuous input space of an arbitrary function. Instead of making impractical assumptions just to satisfy such prerequisites, we leave the regret bound for general settings as an open question.

## 4.6 Remarks on objective functions

Like most of the optimization literature, it is difficult to claim one approach is always better than the other. This phenomenon is widely recognized: “no free lunch in optimization”. In this work, however, we narrow down the assumptions in Bayesian optimization, and show that if we pre-train the GP prior from multi-tasks, we have more guarantees to escape “no free lunch”. As shown in §5, BO with NLL or KL based pre-trained GPs both improved upon other alternatives.

Yet, NLL and KL objectives are developed for different scenarios: NLL works unanimously with any types of multi-task data sizes or collection strategies. KL only works with the same data sizes and same inputs across tasks. If the dataset only has data arbitrarily collected for each task, the only option is to use NLLs as objectives. If most of the dataset is composed of same inputs across tasks, we’d then recommend the empirical KL as the objective for the reasons below.

Although NLLs are straightforward objectives to optimize, it can be difficult to interpret how high of a likelihood is high enough for us to stop our search for a decent model. The empirical KL

divergence in Eq. 3, on the other hand, is a “distance” that goes to 0 as the difference between two distributions reduces. One may choose to do early stopping or model selection based on how close Eq. 3 is to 0. From information theory, we know that the KL objective in Eq. 3 describes the number of extra bits (or nats) to encode the multivariate Gaussian  $\mathcal{N}(\tilde{\mu}, \tilde{K})$ . Overall we found the empirical divergence in Eq. 3 relatively more interpretable than the data likelihood in Eq. 4.

The KL objective in Eq. 3 also introduces a different optimization landscape than the NLL one in Eq. 4. The KL objective makes use of the matching dataset  $D'_N$  in a way that the NLL objective cannot. In fact, all matching inputs in the NLL (Eq. 4) are implicit: all inputs are passed in to mean/kernel functions, and so there is no way that the NLL can be informed that some inputs are the same across tasks. As shown in §5, the empirical KL objective in Eq. 3 interestingly led to better results in our experiments.

Alternatively, we can additively combine the NLL and KL objectives that treats KL as a regularizer (see §A more discussions), though, their performances can be similar (Fig. 9) if the matching dataset is already large enough. Based on our experiments, we’d suggest using KL when the data allows and using NLL when there are no data points with inputs shared across tasks. If the subset of data with the same inputs across tasks is very small, using NLL with KL as a regularizer could be an alternative to get the best of both worlds.

## 5. Experiments

Our goal in this paper is to provide a practical approach for hyperparameter optimization when we are given data on a range of tasks over the same search space. To analyze the effectiveness of our proposal, we take the optimizer hyperparameter tuning problem in deep learning as a case study. Our JAX-based (Bradbury et al., 2018) implementation of HyperBO can be found at <https://github.com/google-research/hyperbo>, which was used for all of our experiments. To accommodate needs for more modular use cases, we also provide a Flax (Heek et al., 2020) and TensorFlow-Probability (Dillon et al., 2017) based implementation for GP pre-training at <https://github.com/google-research/gpax>.

For empirical validation, we first collected a dataset composed of hyperparameter evaluations on various deep neural network training tasks. The tasks included optimizing deep models on image, text, and other datasets (see more details in Sec. 5.1). We then compared our method to several competitive baselines in realistic hyperparameter tuning scenarios in deep neural net optimizers to understand HyperBO’s properties better.

To reduce ambiguity, we distinguish between datasets that individual neural networks are trained on and the dataset we collected that includes optimizer hyperparameter points with their validation errors (and other metrics). We will call the former (e.g. MNIST, CIFAR10) task datasets and call the latter the tuning dataset. The tuning dataset is what we described as dataset  $D_N$  in §3.

### 5.1 Hyperparameter tuning dataset

In order to collect our hyperparameter tuning dataset, the PD1 Neural Net Tuning Dataset, we defined a set of 24 neural network tuning tasks<sup>2</sup> and a single, broad search space for Nesterov momentum. Each task is defined by a task dataset (e.g. ImageNet), a specific neural network model

2. The batch size 1024 ResNet50 ImageNet task only has 100 hyperparameter points because we abandoned it when scaling up data collection in order to save compute resources. It is used in training, but not evaluation.

(e.g. ResNet50), and a batch size. Tab. 2 shows all the tasks that we consider in the tuning dataset. We used an existing code base (Gilmer et al., 2021) for neural network model training. The dataset used roughly 12,000 machine-days of computation for approximately 50,000 hyperparameter evaluations.

For each task, we trained the model on the task dataset repeatedly using Nesterov momentum (Nesterov, 1983; Sutskever et al., 2013), with the task’s minibatch size, with different hyperparameter settings drawn from the 4-dimensional search space detailed in Tab. 3. We tuned the base learning rate,  $\eta$ , on a log scale, the momentum,  $\beta$ , with  $1 - \beta$  on a log scale, and the polynomial learning rate decay schedule power  $p$  and decay steps fraction  $\lambda$ . We used a polynomial decay schedule with the following form:

$$\eta_\tau = \frac{\eta}{1000} + \left( \eta - \frac{\eta}{1000} \right) \left( 1 - \frac{\min(\tau, \lambda\mathcal{T})}{\lambda\mathcal{T}} \right)^p, \quad (6)$$

where  $\tau$  is the training step and  $\mathcal{T}$  is the total number of training steps for the task.

We collected two types of data: matched and unmatched data. Matched data used the same set of uniformly-sampled hyperparameter points across all tasks and unmatched data sampled new points for each task. All other training pipeline hyperparameters were fixed to hand-selected, task-specific default values. All of our tasks are classification problems, so they all used the same training loss, although occasionally task-specific regularization terms were added. For each trial (training run for a single hyperparameter point), we recorded validation error (both cross entropy error and misclassification rate). In many cases, poor optimizer hyperparameter choices can cause training to diverge. We detected divergent training when the training cost became NaN and then marked the trial but did not discard it. Please download the dataset (<http://storage.googleapis.com/gresearch/pint/pd1.tar.gz>) and see its descriptions for additional details about the tasks and training procedure. The different tuning tasks vary in difficulty and numbers of data points, but generally there are roughly 500 matched datapoints and 1500 unmatched datapoints per tuning task. For unmatched data only, we attempted to generate roughly similar numbers of non-divergent points across tasks, so tasks with a higher probability of sampling a hyperparameter point that causes training to diverge will tend to have more trials.

Table 2: Tasks

Task Dataset	Model	Batch Sizes
CIFAR10	Wide ResNet	{256, 2048}
CIFAR100	Wide ResNet	{256, 2048}
Fashion MNIST	Max pool CNN ReLU	{256, 2048}
Fashion MNIST	Max pool CNN tanh	{256, 2048}
Fashion MNIST	Simple CNN	{256, 2048}
ImageNet	ResNet50	{512, 1024, 2048}
LM1B	Transformer	{2048}
MNIST	Max pool CNN relu	{256, 2048}
MNIST	Max pool CNN tanh	{256, 2048}
MNIST	Simple CNN	{256, 2048}
SVHN (no extra)	Wide ResNet	{256, 1024}
WMT15 German-English	xformer	{64}
uniref50	Transformer	{128}

Table 3: 4-dimensional input search space (see Eq.6)

Hyperparameter	Range	Scaling
$\eta$	$[10^{-5}, 10]$	Log
$p$	$[0.1, 2.0]$	Linear
$1 - \beta$	$[10^{-3}, 1.0]$	Log
$\lambda$	$[0.01, 0.99]$	Linear

## 5.2 Description of all compared methods

Our method HyperBO has several variants including using different acquisition functions and different objectives. In §5, unless otherwise mentioned, we used a thresholded probability of im-

provement (PI) as the acquisition function<sup>3</sup>. We set PI in line 5 of Alg. 1 as  $\alpha \left( x; \mathcal{GP}(\hat{\mu}, \hat{k} \mid D_f) \right) = \frac{\hat{\mu}_{D_f}(x) - \max_t(y_t + 0.1)}{\hat{\sigma}_{D_f}(x)}$ .

- H\* NLL: HyperBO with PI as the acquisition function and pre-trained with the NLL objective. The GP has a one-hidden layer neural network of size 8 as mean function and a Matérn32 covariance on the feature layer of the mean function as kernel.
- H\* KL: HyperBO with PI as the acquisition function and pre-trained with the empirical KL divergence on matching datapoints as objective. The same GP model is used as H\* NLL.

These two settings of HyperBO are representative of variants of HyperBO. We provide more comparisons over acquisition functions and other objective functions in Appendix C.

Our first set of baselines include those that *do not* use information from training tasks:

- Rand: Random search in the corresponding scaled search space in §5.1 (see Tab. 3).
- STBO: Single task BO with a constant mean function, Matérn32 kernel and PI acquisition function (same as above). Every BO iteration, STBO optimizes the GP hyperparameters via type-II maximum likelihood on the test task. This implementation corresponds to the basic off-the-shelf BO setups.
- STBOH: Single task GP-UCB (coefficient=1.8) with constant mean, Matérn52 kernel and *hand-tuned* prior on hyper-parameters including UCB coefficient (Srinivas et al., 2010; Golovin et al., 2017). Specifically, log amplitude follows Normal(-1, 1), log length scale (one per input parameter) follows Normal(0,1), and log observation noise variance follows Normal(-6, 3). The hyperparameters are post-processed by tensorflow-probability’s `SoftClip` bijector to constrain the values between 1-st and 99-th quantiles. These prior distributions were manually tuned to obtain reasonable convergence rates on 24 analytic functions in COCO (Hansen et al., 2021). The GP parameters are then optimized via maximum marginal likelihood every BO iteration.

For multi-task BO baselines, we included scalable methods that replace the GP with a regression model that can be trained using SGD and thus scales linearly in the number of observations. Following the multi-task setup of Springenberg et al. (2016), we jointly trained a 5-dimensional embedding of each task, which was then added to the input of the following two models.

- MIMO: Multi-task BO with GP bases as an ensemble of feedforward neural networks with shared subnetworks (Havasi et al., 2020; Kim et al., 2021). We trained an ensemble of feedforward neural networks with shared subnetworks (Havasi et al., 2020). We used 1 shared dense layer of size 10 and 2 unshared layers of size 10. We used tanh activation based on Figure 2 from Snoek et al. (2015). The network has one output unit with linear activation and another with  $\text{softmax}(10^{-4}, 1)$  activation, corresponding respectively to the mean and

3. We empirically evaluated a variety of acquisition functions, but found PI thresholded at 0.1 to be surprisingly effective. Because we model the observations as log error rate, this actually trades off exploration and exploitation - i.e. with larger error rates this seeks relatively more substantial improvements than with small error rates. The list of 5 different acquisition functions we tested is as follows: PI with 0.1 threshold, expected improvement and UCB with 2, 3, 4 coefficients. More results can be found in Appendix C.2.

standard deviation parameters of a normal distribution. We trained for 1000 epochs using the Adam optimizer with learning rate  $10^{-4}$  and batch size 64.

- RFGP: Multi-task BO with GP bases as random features (Snoek et al., 2015; Krause and Ong, 2011). We used the open-source implementation of approximate GP by Liu et al. (2020). We trained for 1000 epochs using the Adam optimizer with learning rate  $10^{-3}$  and batch size 64.
- MAF: We refer to the meta BO method from Volpp et al. (2020) as MAF (Meta Acquisition Function) to avoid confusion. MAF used reinforcement learning to learn an acquisition function modeled by a neural network over a set of transfer learning tasks. All MAF results were generated using the code from Volpp et al. (2020). See App. C.3 for experimental details. As MAF takes significantly longer to run than HyperBO and other methods, we only include its results for §5.3.4 and §5.3.1.

All methods share the same input and output warping. The input warping is done according to the scaling function in Tab. 3:  $\eta \leftarrow \log \eta$ ,  $1 - \beta \leftarrow \log(1 - \beta)$ . The output warping is done for the validation error rate  $r \leftarrow -\log(r + 10^{-10})$ .

### 5.3 Results on offline optimizer hyperparameter tuning tasks

Many tasks in §5.1 can use up a lot of compute resources and time, which makes it infeasible to perform a wide variety of experiments to analyze the characteristics of BO methods. Hence we adopt an offline approximation, which runs BO only on the finite set of points that each tuning sub-dataset contains. In §5.5, we show some BO comparisons in the online setting.

In all the experiments in this section, we ran offline BO on the data from the test task starting from zero initial data from this task. Each method was repeated 5 times with different random seeds to initialize its model. We ran all methods without de-duplication to best simulate online BO settings. We evaluate methods on regret on error rate which denotes the simple regret on the finite set of data points in each tuning sub-dataset.

#### 5.3.1 HOLDING OUT RELEVANT TASKS

We first conducted experiments in a setting where a new task dataset is presented, and a BO method is trying to tune the optimizer hyperparameters for a selected model on that task dataset. A training dataset for meta BO is composed of at most 18 tuning sub-datasets on training tasks that do not involve the same task dataset as the test task. All methods then proceed to solve the test task on the new task dataset. Fig. 2 shows *performance profiles* of the BO methods described in §5.2. The performance profiles show the fraction of all test tasks that each method is able to outperform a baseline criterion at each BO iteration.<sup>4</sup> We chose the criterion to be the median of best error rates achieved by all methods at 3 different BO iterations: 25th, 50th or 100th. The larger the fraction of tasks at each BO iteration, the better the method is. From all 3 criteria, we can see that MIMO is able to outperform other methods in the beginning 10 to 20 BO iterations, but its leading position soon gets surpassed by HyperBO (H\* NLL and H\* KL) and STBOH. HyperBO methods are gaining a similar if not larger fraction than the best alternative, STBOH, throughout BO iterations. Fig. 2 (c) has the most stringent performance criterion, and it shows that HyperBO with the KL objective

4. We show performance relative to a baseline because of varying scales across tasks.



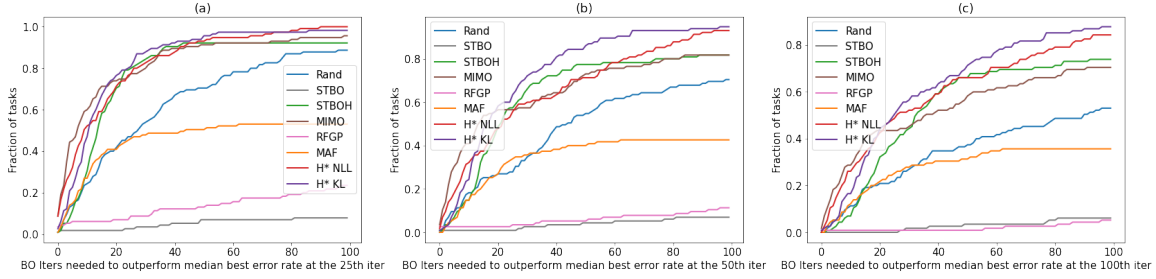


Figure 2: Performance profiles for outperforming the median of best error rates at the (a) 25th BO iteration, (b) 50th BO iteration and (c) 100th BO iteration.

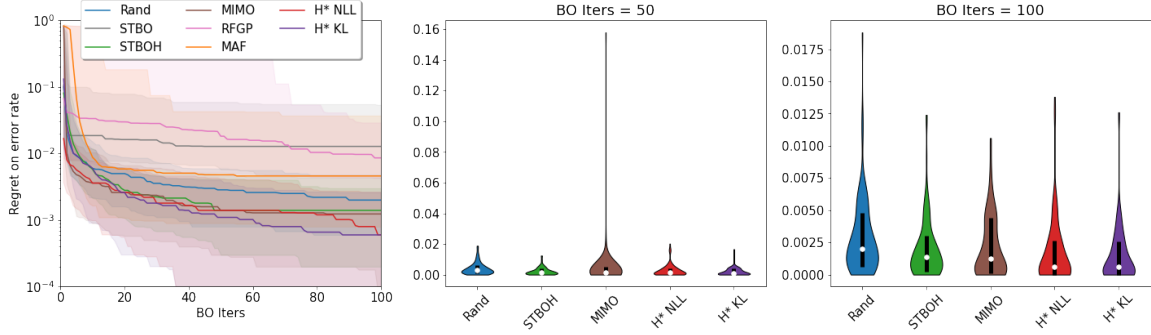


Figure 3: The left most is a summary of the BO convergence of all methods: the median and 20/80 percentiles of the regrets on error rates over 115 BO runs: 23 tasks and each with 5 repeats of different random seeds. We also show violin plots on its two vertical slices at 50th and 100th iteration, where the white dot is the median and the black line is the 20/80 percentile. Overall, HyperBO methods H\* NLL and H\* KL are able to achieve the lowest regret on error rate on the majority of tasks.

outperforms HyperBO with the NLL objective in this set of experiments with a small margin. And both methods in HyperBO are doing considerably better than others.

Fig. 3 illustrates the BO convergence curves of all competing methods, together with the vertical slices at the 50th and 100th iterations. RFGP and STBO are both falling much behind Random search. STBO trains the GP on the data that the GP suggests to query, which creates a loop that could be harmful for data acquisition. Optimizing the marginal data likelihood on at most 100 datapoints in fact may not lead to a better model than random initialization (see Tab. 6 in §6). Surprisingly, RFGP, though equipped with the tuning dataset and initially reached some good values, performed similarly to STBO in the end. Clearly, the contextual information learned by RFGP did not generalize to a new task. On the other hand, MIMO is able to obtain a slightly better error rate than STBOH.

Fig. 2 and Fig. 3 both show that learning the GP prior through data as with HyperBO outperforms other meta BO methods, and is a more principled and effective way to obtain the GP prior when compared with hand-tuning. As a reference, we include Tab. 4 which shows the task-wise best validation error rates obtained by the top 5 methods in 100 BO iterations.

To more precisely quantify HyperBO’s advantage, we also computed how much faster HyperBO can get a better error rate than best alternatives, which can be different from task to task. We found that on average, on over 50% tasks, H\* NLL is at least 2.86 times faster than best non-HyperBO alternatives; while on over 57% tasks, H\* KL is at least 3.26 times faster than best non-HyperBO

Table 4: The mean and standard error of best validation error rates (%) for each test task in the offline optimizer hyperparameter tuning experiments. Meta BO methods including MIMO and HyperBO variants (H\* NLL and H\* KL) have access to training tasks that do not share the same task dataset as the test task. We show results of random search and the top 5 methods, and we highlight the lowest error rates in bold.

	Rand	STBOH	MIMO	MAF	H* NLL	H* KL
WMT XFormer 64	34.27 $\pm$ 0.16	34.15 $\pm$ 0.15	34.29 $\pm$ 0.16	47.32 $\pm$ 8.77	<b>33.94 <math>\pm</math> 0.01</b>	33.99 $\pm$ 0.03
Uniref50 Transformer 128	79.06 $\pm$ 0.04	78.92 $\pm$ 0.12	78.93 $\pm$ 0.11	79.88 $\pm$ 0.39	<b>78.64 <math>\pm</math> 0.00</b>	78.74 $\pm$ 0.09
LM1B Transformer 2048	61.96 $\pm$ 0.03	61.95 $\pm$ 0.04	61.95 $\pm$ 0.01	62.26 $\pm$ 0.29	<b>61.83 <math>\pm</math> 0.01</b>	<b>61.82 <math>\pm</math> 0.01</b>
SVHN WRN 1024	3.99 $\pm$ 0.04	4.05 $\pm$ 0.10	<b>3.82 <math>\pm</math> 0.04</b>	4.24 $\pm$ 0.08	4.11 $\pm$ 0.04	4.06 $\pm$ 0.02
SVHN WRN 256	3.71 $\pm$ 0.01	3.72 $\pm$ 0.02	<b>3.62 <math>\pm</math> 0.02</b>	3.69 $\pm$ 0.03	3.79 $\pm$ 0.01	3.77 $\pm$ 0.02
ImageNet ResNet50 256	23.03 $\pm$ 0.07	22.66 $\pm$ 0.07	22.69 $\pm$ 0.06	23.75 $\pm$ 0.25	22.57 $\pm$ 0.02	<b>22.55 <math>\pm</math> 0.02</b>
ImageNet ResNet50 512	23.02 $\pm$ 0.11	22.74 $\pm$ 0.05	22.99 $\pm$ 0.05	24.85 $\pm$ 0.86	<b>22.65 <math>\pm</math> 0.02</b>	22.75 $\pm$ 0.03
MNIST CNNPoolTanh 2048	0.55 $\pm$ 0.01	<b>0.53 <math>\pm</math> 0.01</b>	<b>0.52 <math>\pm</math> 0.01</b>	0.56 $\pm$ 0.02	<b>0.53 <math>\pm</math> 0.01</b>	0.54 $\pm$ 0.01
MNIST CNNPoolTanh 256	0.51 $\pm$ 0.01	0.48 $\pm$ 0.01	<b>0.46 <math>\pm</math> 0.00</b>	<b>0.46 <math>\pm</math> 0.01</b>	<b>0.46 <math>\pm</math> 0.01</b>	<b>0.46 <math>\pm</math> 0.01</b>
MNIST CNNPoolReLU 2048	0.69 $\pm$ 0.01	0.73 $\pm$ 0.02	0.66 $\pm$ 0.01	0.73 $\pm$ 0.03	<b>0.64 <math>\pm</math> 0.01</b>	0.65 $\pm$ 0.01
MNIST CNNPoolReLU 256	0.51 $\pm$ 0.01	0.55 $\pm$ 0.03	0.52 $\pm$ 0.01	0.55 $\pm$ 0.01	<b>0.48 <math>\pm</math> 0.00</b>	0.49 $\pm$ 0.00
MNIST CNNReLU 2048	1.14 $\pm$ 0.03	1.20 $\pm$ 0.09	1.11 $\pm$ 0.02	71.21 $\pm$ 15.65	<b>1.06 <math>\pm</math> 0.00</b>	1.08 $\pm$ 0.01
MNIST CNNReLU 256	1.09 $\pm$ 0.02	1.06 $\pm$ 0.01	1.08 $\pm$ 0.02	18.61 $\pm$ 15.67	<b>1.03 <math>\pm</math> 0.00</b>	<b>1.03 <math>\pm</math> 0.00</b>
Fashion CNNPoolTanh 2048	7.14 $\pm$ 0.06	7.10 $\pm$ 0.05	<b>7.05 <math>\pm</math> 0.06</b>	39.91 $\pm$ 17.98	<b>7.03 <math>\pm</math> 0.03</b>	7.16 $\pm$ 0.02
Fashion CNNPoolTanh 256	6.51 $\pm$ 0.03	6.67 $\pm$ 0.18	6.41 $\pm$ 0.07	45.06 $\pm$ 16.72	6.38 $\pm$ 0.02	<b>6.28 <math>\pm</math> 0.01</b>
Fashion CNNPoolReLU 2048	7.47 $\pm$ 0.02	7.48 $\pm$ 0.04	7.52 $\pm$ 0.06	7.73 $\pm$ 0.14	<b>7.42 <math>\pm</math> 0.03</b>	7.53 $\pm$ 0.04
Fashion CNNPoolReLU 256	6.78 $\pm$ 0.04	6.74 $\pm$ 0.01	6.89 $\pm$ 0.06	7.03 $\pm$ 0.10	6.79 $\pm$ 0.04	<b>6.70 <math>\pm</math> 0.01</b>
Fashion CNNReLU 2048	7.70 $\pm$ 0.03	<b>7.47 <math>\pm</math> 0.09</b>	7.64 $\pm$ 0.06	56.60 $\pm$ 17.83	7.57 $\pm$ 0.01	7.56 $\pm$ 0.02
Fashion CNNReLU 256	7.70 $\pm$ 0.04	7.46 $\pm$ 0.11	7.83 $\pm$ 0.05	40.46 $\pm$ 17.78	7.44 $\pm$ 0.12	<b>7.25 <math>\pm</math> 0.05</b>
CIFAR100 WRN 2048	21.28 $\pm$ 0.27	<b>20.78 <math>\pm</math> 0.19</b>	<b>20.94 <math>\pm</math> 0.13</b>	22.63 $\pm$ 0.68	21.26 $\pm$ 0.23	20.98 $\pm$ 0.24
CIFAR100 WRN 256	19.17 $\pm$ 0.19	19.02 $\pm$ 0.03	19.15 $\pm$ 0.06	27.28 $\pm$ 7.19	19.07 $\pm$ 0.04	<b>18.98 <math>\pm</math> 0.01</b>
CIFAR10 WRN 2048	3.73 $\pm$ 0.05	<b>3.43 <math>\pm</math> 0.07</b>	<b>3.40 <math>\pm</math> 0.04</b>	37.16 $\pm$ 18.26	3.66 $\pm$ 0.11	<b>3.43 <math>\pm</math> 0.07</b>
CIFAR10 WRN 256	2.84 $\pm$ 0.04	2.88 $\pm$ 0.06	2.84 $\pm$ 0.05	20.00 $\pm$ 15.00	2.83 $\pm$ 0.03	<b>2.80 <math>\pm</math> 0.02</b>

alternatives. Moreover, on over 73% tasks, H\* NLL is at least 7.74 times faster than random search; and on over 75% tasks, H\* KL is at least 6.07 times faster than random search.

### 5.3.2 EFFECT OF NUMBER OF TRAINING TASKS

We now investigate the impact of number of training tasks on the performance of meta BO methods. In Fig 4 we show the BO simple regrets on tasks from Table 2 (except ImageNet ResNet50 2048) that use meta BO models trained on different number of training tasks. To analyze the performance of all methods on less-related tasks, we first remove training tasks that have the same task dataset as our current tuning task for testing, and then remove randomly selected training datasets from the rest.

HyperBO variants were able to reduce the simple regret as more training tasks are given. Interestingly, both H\* NLL and H\* KL are already slightly better than Rand and STBOH when they started off with only 3 training tasks. There are reasonable fluctuations in the results but overall the trend of regret is going down as the number of training tasks increases. MIMO also reduced regret when the number of tasks increased from 8 to 18. RFGP, however, fails to learn from training tasks possibly because it did not learn good task embeddings for GP regression models.

### 5.3.3 EFFECT OF NUMBER OF DATA POINTS IN TRAINING TASKS

One remaining question is, how does  $M_i$  in §3, the number of data points in each training tasks, affect the performance of meta BO methods. We analyze the impact of  $M_i$  by removing a portion of all data that we have access to for each task. In particular, we set the percentage of remaining data to be 0.2%, 0.5%, 1%, 3%, 5%, 10%, 30%, 50%, 70%, 90%. Remaining datapoints are selected uniformly randomly, which breaks the structure of matching data. Hence we do not include H\* KL in this comparison, as H\* KL only makes use of matching data.

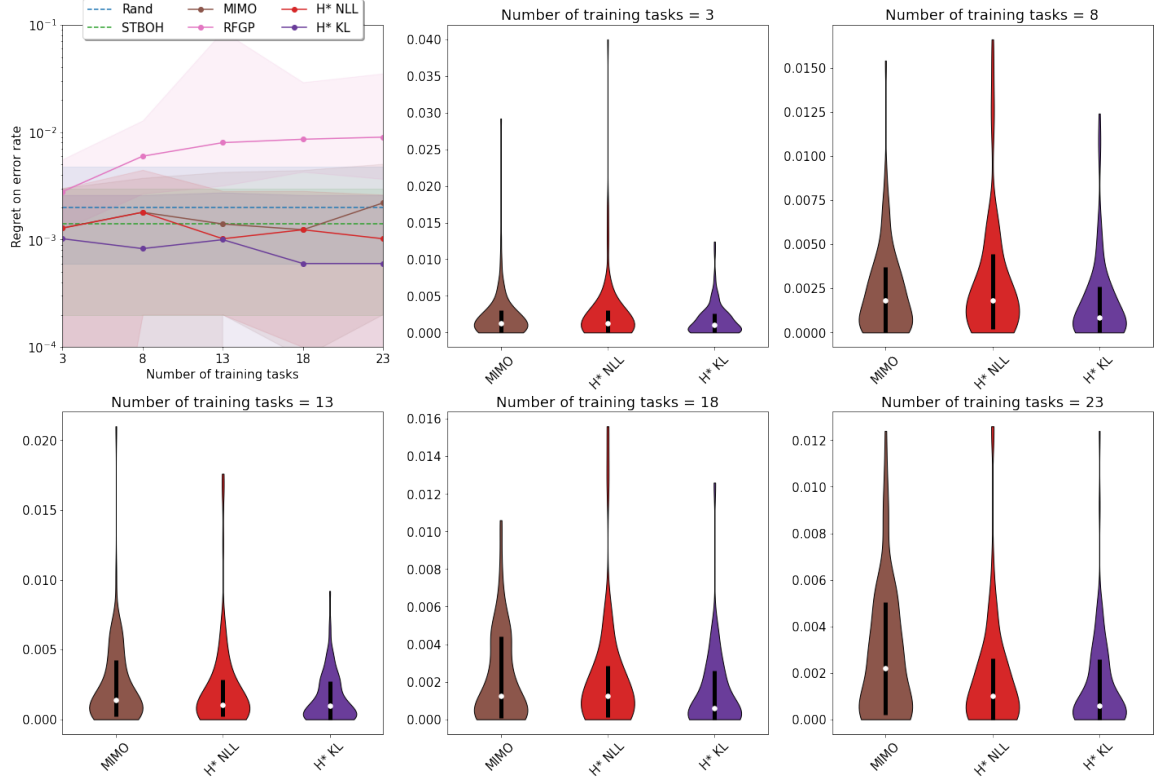


Figure 4: Aggregated BO results on 23 tasks (all in Table 2 except ImageNet ResNet50 2048 because of insufficient data) that uses models trained on 3 to 23 training tasks. Note that the models are never trained on the data from the test task that we run BO on. If the number of training tasks is less than 23, we first remove the tasks that involve the same task dataset as the test task and then remove others randomly until we reach the designated number of training tasks. The top left shows the median and 20/80 percentiles of regret on best validation error rate for each method. The rest are violin plots showing the regret for MIMO, H\* NLL and H\* KL, where white dots indicate the median and black lines the 20/80 percentiles.

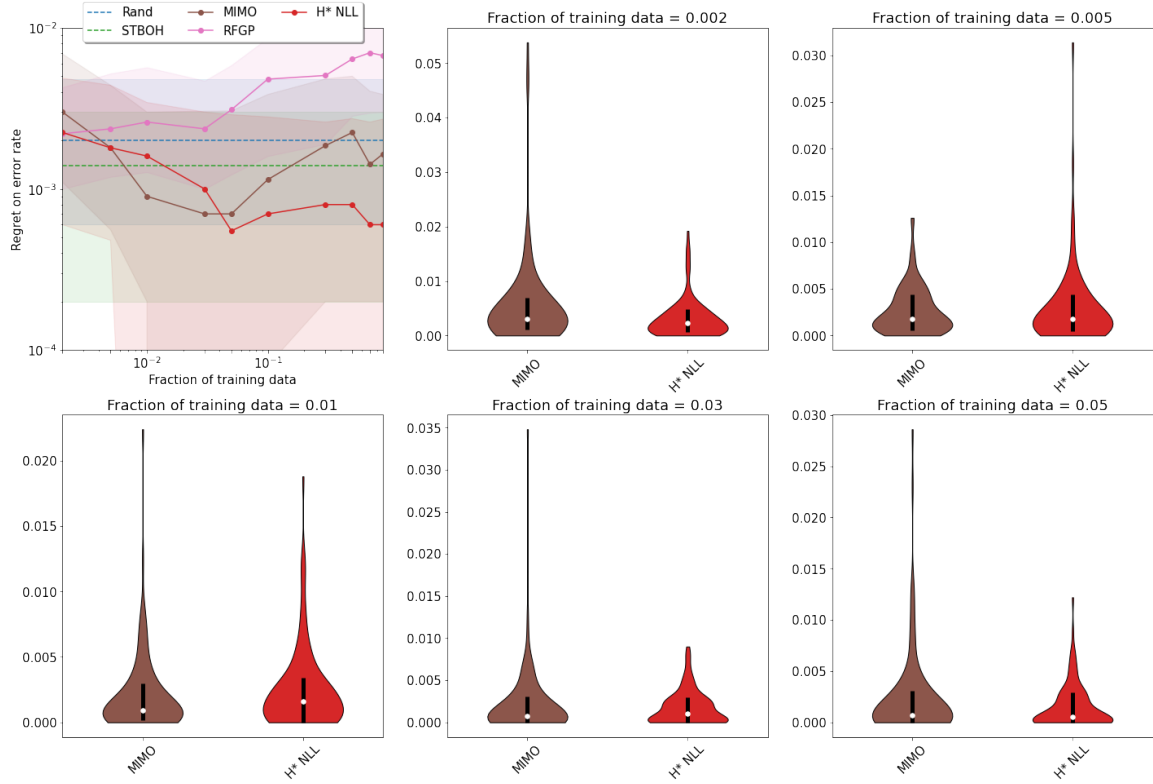


Figure 5: Aggregated BO results on 23 tasks (all in Table 2 except ImageNet ResNet50 2048 because of insufficient data) that uses models trained on 0.2% to 90% of data in each task. Note that the models are never trained on the data from the test task that we run BO on. The top left is the median and 20/80 percentiles of simple regret in log scale. The rest of the figures are simple regret violin plots for MIMO and H\* NLL

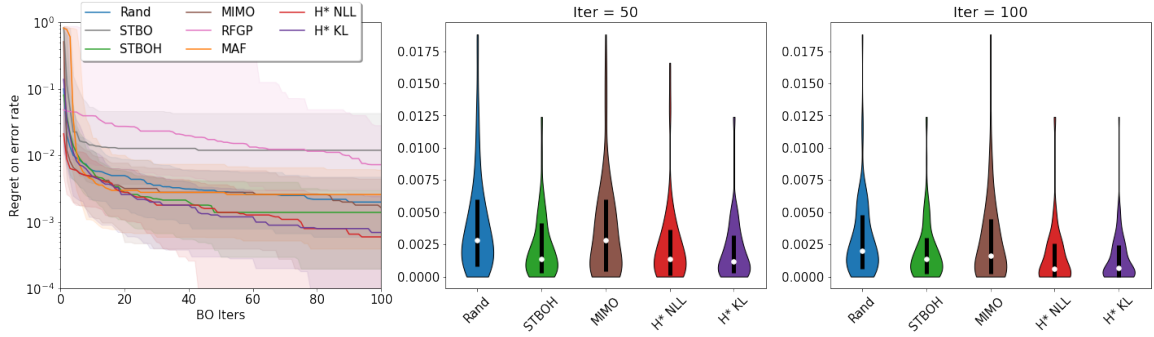


Figure 6: Aggregated leave-one-out BO convergence results on 23 tasks, each with 5 repeats using different random seeds. The left most is the median and 20/80 percentiles of the regrets on error rates. We also show violin plots on its two vertical slices at 50th and 100th iteration, where the white dot is the median and the black line is the 20/80 percentile.

Fig. 5 shows how the simple regret changes as the fraction of training data grows. Below 10% training data, we observe clear trend that more data lead to lower regret for both H\* NLL and MIMO, and relatively no change for RFGP. We also found that the performance of HyperBO (H\* NLL) does not change much as the fraction of training data increases from 5% to 90%. However, MIMO and RFGP suffers significantly from more data as the fraction of training data increases from 5% to 50%. It is not entirely clear why MIMO and RFGP have such behaviors. One conjecture is that neural network based Bayesian linear regression models may get too confident once the amount of data reaches a certain threshold. This means much less exploration if those models are used for BO.

#### 5.3.4 TRAINING ON ALL BUT ONE TASK

We also studied the case where meta BO approaches have access to both training tasks that do not use the same task dataset and training tasks that use the same task dataset but different model configurations. This is especially common when we do architecture search: we aim to find the best model and we are tuning the optimizer hyperparameters for a new machine learning model given tuning data of the same task dataset on some other models.

For this section only, we added a new baseline, MAF: We refer to the meta BO method from Volpp et al. (2020) as MAF (Meta Acquisition Function) to avoid confusion. MAF used reinforcement learning to learn an acquisition function modeled by a neural network over a set of transfer learning tasks. All MAF results were generated using the code from Volpp et al. (2020). See App. C.3 for experimental details. As MAF takes significantly longer to run than HyperBO and other methods, we only include its results for this section.

We carried out a series of leave-one-out experiments, where we picked one task as the BO test function and let meta BO methods train on the remaining tasks. In Fig. 6, we aggregated results from all 23 tasks to show the trend of how each method performs.

The conclusions are similar to those from §5.3.1. As expected, STBO without any tricks to avoid pitfalls of vanilla BO did not show very good results. We inspected its learned GP which mimicked a Dirac function that is flat almost everywhere except some locations, and hence it got very confident that it landed at a good spot and lost its ability to explore.

Table 5: The mean and standard error of best validation error rates (%) for each test task in the offline leave-one-out experiments. We show results of the top 6 methods, and we highlight the lowest error rates in bold.

	Rand	STBOH	MIMO	MAF	H* NLL	H* KL
WMT XFormer 64	34.27 $\pm$ 0.16	34.15 $\pm$ 0.15	34.40 $\pm$ 0.13	34.09 $\pm$ 0.09	<b>33.91 <math>\pm</math> 0.01</b>	33.97 $\pm$ 0.02
Uniref50 Transformer 128	79.06 $\pm$ 0.04	78.92 $\pm$ 0.12	79.17 $\pm$ 0.13	79.34 $\pm$ 0.27	78.71 $\pm$ 0.06	<b>78.64 <math>\pm</math> 0.00</b>
LM1B Transformer 2048	61.96 $\pm$ 0.03	61.95 $\pm$ 0.04	61.96 $\pm$ 0.05	62.02 $\pm$ 0.10	<b>61.81 <math>\pm</math> 0.01</b>	<b>61.81 <math>\pm</math> 0.01</b>
SVHN WRN 1024	3.99 $\pm$ 0.04	4.05 $\pm$ 0.10	<b>3.83 <math>\pm</math> 0.04</b>	4.10 $\pm$ 0.09	4.10 $\pm$ 0.02	4.08 $\pm$ 0.01
SVHN WRN 256	3.71 $\pm$ 0.01	3.72 $\pm$ 0.02	<b>3.65 <math>\pm</math> 0.01</b>	3.69 $\pm$ 0.03	3.78 $\pm$ 0.01	3.72 $\pm$ 0.03
ImageNet ResNet50 256	23.03 $\pm$ 0.07	22.66 $\pm$ 0.07	22.73 $\pm$ 0.07	26.44 $\pm$ 1.98	<b>22.53 <math>\pm</math> 0.00</b>	22.58 $\pm$ 0.04
ImageNet ResNet50 512	23.02 $\pm$ 0.11	22.74 $\pm$ 0.05	23.01 $\pm$ 0.05	25.46 $\pm$ 1.41	<b>22.65 <math>\pm</math> 0.02</b>	22.79 $\pm$ 0.03
MNIST CNNPoolTanh 2048	0.55 $\pm$ 0.01	<b>0.53 <math>\pm</math> 0.01</b>	<b>0.53 <math>\pm</math> 0.01</b>	<b>0.52 <math>\pm</math> 0.01</b>	0.59 $\pm$ 0.02	0.54 $\pm$ 0.00
MNIST CNNPoolTanh 256	0.51 $\pm$ 0.01	0.48 $\pm$ 0.01	<b>0.47 <math>\pm</math> 0.00</b>	<b>0.47 <math>\pm</math> 0.01</b>	<b>0.46 <math>\pm</math> 0.01</b>	<b>0.47 <math>\pm</math> 0.01</b>
MNIST CNNPoolReLU 2048	0.69 $\pm$ 0.01	0.73 $\pm$ 0.02	0.67 $\pm$ 0.02	0.68 $\pm$ 0.01	<b>0.64 <math>\pm</math> 0.00</b>	0.70 $\pm$ 0.03
MNIST CNNPoolReLU 256	0.51 $\pm$ 0.01	0.55 $\pm$ 0.03	0.50 $\pm$ 0.01	0.51 $\pm$ 0.01	<b>0.49 <math>\pm</math> 0.00</b>	<b>0.49 <math>\pm</math> 0.00</b>
MNIST CNNReLU 2048	1.14 $\pm$ 0.03	1.20 $\pm$ 0.09	1.10 $\pm$ 0.01	1.17 $\pm$ 0.02	<b>1.06 <math>\pm</math> 0.00</b>	1.11 $\pm$ 0.02
MNIST CNNReLU 256	1.09 $\pm$ 0.02	1.06 $\pm$ 0.01	1.08 $\pm$ 0.02	1.07 $\pm$ 0.02	<b>1.03 <math>\pm</math> 0.00</b>	1.04 $\pm$ 0.01
Fashion CNNPoolTanh 2048	7.14 $\pm$ 0.06	7.10 $\pm$ 0.05	<b>7.01 <math>\pm</math> 0.04</b>	7.12 $\pm$ 0.04	<b>7.00 <math>\pm</math> 0.04</b>	<b>7.02 <math>\pm</math> 0.07</b>
Fashion CNNPoolTanh 256	6.51 $\pm$ 0.03	6.67 $\pm$ 0.18	6.40 $\pm$ 0.05	6.47 $\pm$ 0.03	6.40 $\pm$ 0.04	<b>6.34 <math>\pm</math> 0.04</b>
Fashion CNNPoolReLU 2048	<b>7.47 <math>\pm</math> 0.02</b>	<b>7.48 <math>\pm</math> 0.04</b>	7.54 $\pm$ 0.06	7.63 $\pm$ 0.04	<b>7.47 <math>\pm</math> 0.03</b>	<b>7.47 <math>\pm</math> 0.02</b>
Fashion CNNPoolReLU 256	6.78 $\pm$ 0.04	<b>6.74 <math>\pm</math> 0.01</b>	7.03 $\pm$ 0.07	6.84 $\pm$ 0.05	<b>6.74 <math>\pm</math> 0.03</b>	6.81 $\pm$ 0.05
Fashion CNNReLU 2048	7.70 $\pm$ 0.03	<b>7.47 <math>\pm</math> 0.09</b>	7.60 $\pm$ 0.04	40.40 $\pm$ 17.80	<b>7.54 <math>\pm</math> 0.01</b>	7.57 $\pm$ 0.02
Fashion CNNReLU 256	7.70 $\pm$ 0.04	7.46 $\pm$ 0.11	7.84 $\pm$ 0.06	24.13 $\pm$ 14.54	<b>7.29 <math>\pm</math> 0.05</b>	<b>7.25 <math>\pm</math> 0.05</b>
CIFAR100 WRN 2048	21.28 $\pm$ 0.27	<b>20.78 <math>\pm</math> 0.19</b>	21.75 $\pm$ 0.15	50.70 $\pm$ 15.44	21.22 $\pm$ 0.23	<b>20.82 <math>\pm</math> 0.19</b>
CIFAR100 WRN 256	19.17 $\pm$ 0.19	19.02 $\pm$ 0.03	19.12 $\pm$ 0.04	19.84 $\pm$ 0.13	<b>19.00 <math>\pm</math> 0.00</b>	19.04 $\pm$ 0.05
CIFAR10 WRN 2048	3.73 $\pm$ 0.05	<b>3.43 <math>\pm</math> 0.07</b>	<b>3.46 <math>\pm</math> 0.05</b>	<b>3.40 <math>\pm</math> 0.06</b>	3.55 $\pm$ 0.10	<b>3.43 <math>\pm</math> 0.05</b>
CIFAR10 WRN 256	2.84 $\pm$ 0.04	2.88 $\pm$ 0.06	2.89 $\pm$ 0.06	3.04 $\pm$ 0.05	2.82 $\pm$ 0.03	<b>2.74 <math>\pm</math> 0.01</b>

STBOH, on the other hand, achieved very competitive results. This is because it used hand-tuned priors on all of its GP parameters. STBOH hence represents meta BO where meta learning is performed by experts with years of experience. All of our meta BO methods here, however, trains for at most a few hours. As part of the goals of meta learning, we would like to show that it is possible for meta BO methods to exceed or at least match STBOH.

Both HyperBO variants obtained better results than the hand-tuned STBOH. Especially in the beginning few BO iterations, it is able to locate much better hyperparameters than all other methods.

Tab. 5 presents mean and standard error of the best validation error rates achieved in 100 BO iterations on the 23 tasks. HyperBO and its variants were able to achieve the best performance on 20 out of 23 tasks. In Fig. 7, we show the optimization curves of 4 individual tasks that are considered most difficult because few similar task datasets are present in their training data. On all of these 4 difficult tasks, HyperBO identified good hyperparameters much sooner than its competitors.

#### 5.4 Empirical analyses on data likelihoods of HyperBO v.s. misspecified GP priors

To get a better idea on how our “i.i.d functions sampled from the same GP” assumption help on training the GP, we compared NLLs associated with 23 tasks in §5.1 with models obtained via 3 scenarios:

- (a) No training: a randomly initialized GP model with no training;
- (b) Single task: a GP model initialized the same way as “No training” and trained on 100 randomly selected data points of the test task via type-II maximum likelihood (STBO with 100 initial observations on the test task);
- (c) H\* NLL: a GP model initialized the same way as “No training” and pre-trained on 18 irrelevant tasks selected in §5.3.2 via §4.2 (HyperBO with held-out training tasks).



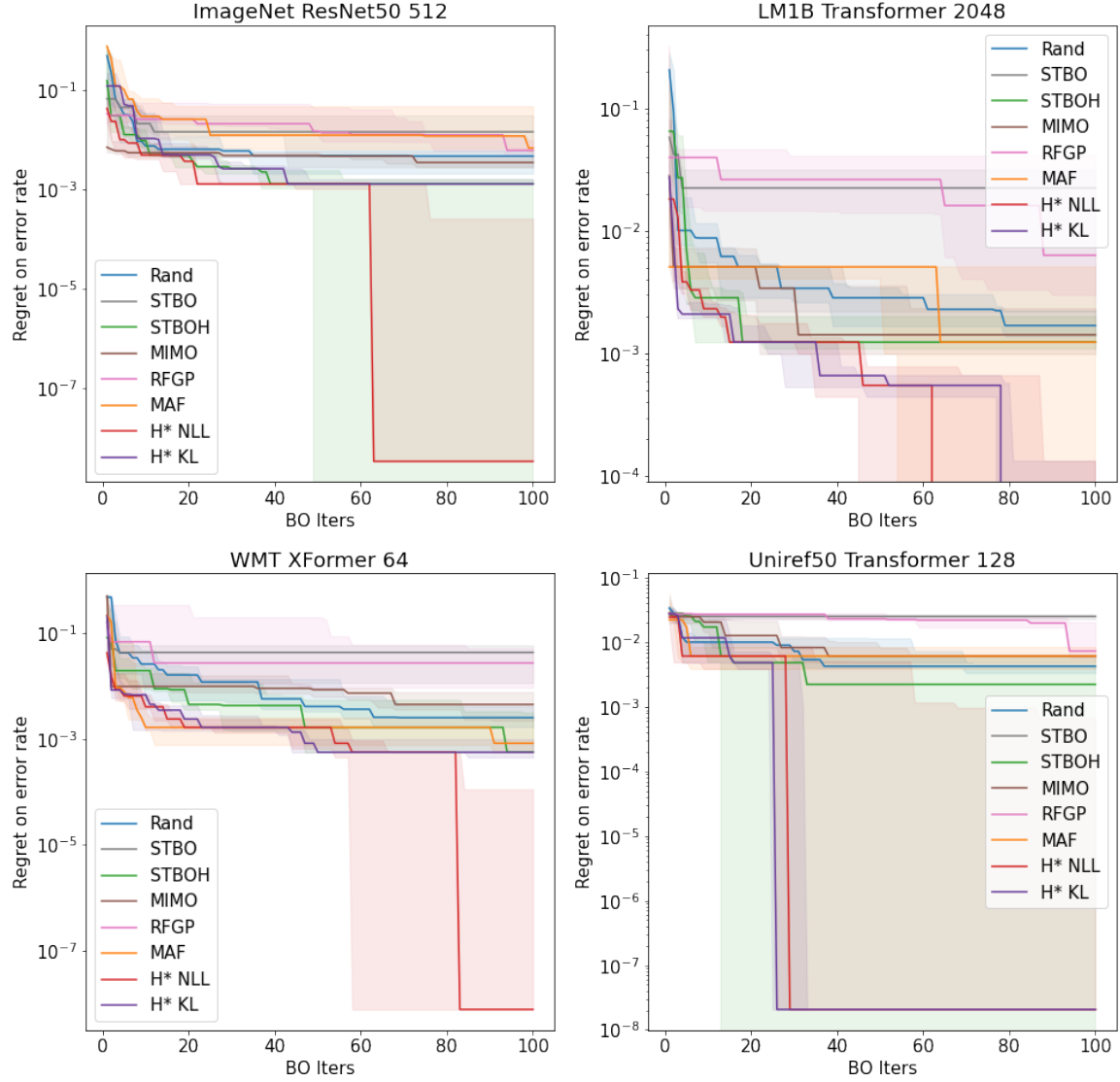


Figure 7: Leave-one-out log regret mean and standard deviation results on ImageNet ResNet50 512, LM1B Transformer 2048, WMT XFormer 64 and Uniref50 Transformer 128. All methods were repeated 5 times with different random seeds to initialize their models. In LM1B Transformer 2048, H\* NLL and H\* KL disappeared around 60 to 80 BO iterations because they reached 0 regret.

In Tab. 6, for each of these 3 methods, we show (1) the NLL of all available data from the test task only; (2) the NLL (Eq. 4) of all available data from all tasks<sup>5</sup>; and (3) the empirical (pseudo) KL divergence (Eq. 3) across all matching datasets. Note that the held-out tasks for some test tasks are the same because of the held-out rules in §5.3.1.

Table 6: NLLs on 23 tasks and (pseudo) KL divergences on matching datasets with trained and randomly initialized GP models. The NLL of randomly initialized model (No training) on all tasks is 148211.2. The KL value of randomly initialized model (No training) is 2177.2. “Single task” training often leads to much worse marginal likelihood on the entire sub-dataset. Training on irrelevant tasks (H\* NLL) achieves much lower (pseudo) KLs on matching datasets and lower NLLs for both the test task only and all tasks.

Test task	NLL of the test task only			NLL of all tasks		(Pseudo) KL	
	No training	Single task	H* NLL	Single task	H* NLL	Single task	H* NLL
WMT XFormer 64	-301.1	159.1	<b>-1735.0</b>	1147900.5	<b>2264.5</b>	9651.9	<b>-40.2</b>
Uniref50 Transformer 128	-651.7	<b>-6829.4</b>	-1850.0	106348128.0	<b>867.9</b>	316672.2	<b>-25.1</b>
LM1B Transformer 2048	-540.6	<b>-2009.7</b>	-1692.7	18840458.0	<b>3565.7</b>	57744.1	<b>-23.5</b>
SVHN WRN 1024	9703.1	72407.5	<b>4267.1</b>	3399330.0	<b>9346.5</b>	4677.9	<b>-0.9</b>
SVHN WRN 256	10770.0	53245.5	<b>3794.8</b>	1164804.5	<b>9346.5</b>	3092.7	<b>-0.9</b>
ImageNet ResNet50 256	1196.7	7483.0	<b>-746.3</b>	7925583.5	<b>-74.2</b>	15028.1	<b>-30.6</b>
ImageNet ResNet50 512	1300.2	6930.3	<b>-673.1</b>	1778823.5	<b>-74.2</b>	9462.1	<b>-30.6</b>
MNIST CNNPoolTanh 2048	10079.7	38871.9	<b>794.8</b>	1375930.1	<b>97.0</b>	3165.5	<b>-32.4</b>
MNIST CNNPoolTanh 256	12147.7	25607.9	<b>550.0</b>	556254.6	<b>-606.0</b>	1255.1	<b>-41.9</b>
MNIST CNNPoolReLU 2048	26870.5	7149.3	<b>5506.6</b>	46538.2	<b>1542.2</b>	113.8	<b>-59.4</b>
MNIST CNNPoolReLU 256	15601.6	6734.6	<b>51.0</b>	88687.7	<b>-782.2</b>	361.2	<b>-41.5</b>
MNIST CNNReLU 2048	13939.2	40619.2	<b>3153.2</b>	743233.1	<b>-231.4</b>	877.6	<b>-61.7</b>
MNIST CNNReLU 256	10111.0	34412.4	<b>1365.3</b>	977295.0	<b>-779.8</b>	1373.3	<b>-46.2</b>
Fashion CNNPoolTanh 2048	2072.8	11433.0	<b>-381.0</b>	1139702.4	<b>-1051.7</b>	1910.5	<b>-37.8</b>
Fashion CNNPoolTanh 256	2800.7	4115.6	<b>-251.4</b>	1278018.0	<b>-1051.7</b>	4208.3	<b>-37.8</b>
Fashion CNNPoolReLU 2048	4677.4	725.2	<b>-405.2</b>	69173.3	<b>-1051.7</b>	205.1	<b>-37.8</b>
Fashion CNNPoolReLU 256	3925.7	4254.4	<b>-755.7</b>	296739.1	<b>-1051.7</b>	1027.1	<b>-37.8</b>
Fashion CNNReLU 2048	4667.3	6778.1	<b>251.9</b>	193488.4	<b>-1051.7</b>	597.0	<b>-37.8</b>
Fashion CNNReLU 256	3295.1	29348.6	<b>-235.1</b>	1526829.2	<b>-1051.7</b>	3341.4	<b>-37.8</b>
CIFAR100 WRN 2048	1271.5	15813.7	<b>-467.4</b>	3306556.5	<b>312.3</b>	25593.7	<b>-19.2</b>
CIFAR100 WRN 256	1957.6	5950.8	<b>-510.9</b>	3468309.0	<b>11.7</b>	9288.4	<b>-25.9</b>
CIFAR10 WRN 2048	5220.6	4917.6	<b>832.9</b>	334488.8	<b>1127.1</b>	1040.4	<b>-14.8</b>
CIFAR10 WRN 256	7819.1	32995.8	<b>463.4</b>	895691.2	<b>847.4</b>	1946.0	<b>-19.6</b>

Comparing NLLs of each test task using models without training and trained via type-II maximum likelihood, it is surprising to see that training on a subset of data points of the sub-dataset of the test task not only did not contribute to lowering NLL on the entire sub-dataset, but it even made it worse in 20 out of 23 test tasks. The training process by optimizing the NLL on a part of a sub-dataset leads to severe over-fitting. We can observe the same results of NLLs on all tasks. Without any training, our NLL is 148211.2. Yet single-task training leads to higher NLLs for all models trained on different sub-datasets.

Our method H\* NLL, on the other end, consistently achieves lower NLLs on both the test task and all tasks. Although it is not entirely clear what the relation is between a better NLL of the GP and better BO results, achieving lower NLLs typically means that the model has a better fit to the dataset. Hence, by the assumption of typical BO methods, the test function should look like a sample from our model, and so lower NLLs of model will contribute to matching the assumption of typical

5. All tasks include ImageNet ResNet50 2048. But it is excluded in the test tasks in Tab. 6 because it has much fewer data points than the others.

BO methods. By enhancing the assumption with ours on *i.i.d.* GP samples, Tab. 6 shows we then will be able to obtain models with a much better fit to the data.

We also computed the (pseudo) KL divergence across all matching datasets in the last columns of Tab. 6. See Appendix A for a comprehensive analysis on pseudo KL divergence for degenerate multivariate Gaussians. Note that pseudo KL divergence can be negative. Here we use pseudo KL divergence if required by the matching dataset. Again, single-task training leads to unstable (pseudo) KL values, sometimes even higher than without training (2177.2). On the contrary, training with H\* NLL leads to much more stable and lower values for KL. This indicates that the model learned to predict similarly to the sample mean/covariance estimate, which is known to help better selection of BO query points by Theorem 2.

### 5.5 Results on online optimizer hyperparameter tuning tasks

Finally, we look into the online BO setting where we optimize over the full hypercube. In the online setting, some combinations of hyperparameters may be *infeasible* to evaluate. For example, an overly big learning rate may lead to divergence in gradients, in which case we do not obtain a valid model. To address this, we pre-process the function values to  $[-2, 2)$  such that *infeasible* evaluations map to  $-2$ , while bad evaluations approach asymptotically to  $-2$ . More precisely, for each subdataset  $D_{f_i}$ , we applied for each successful  $y \in \{y_{j(i)}\}_{j=1}^{M_i}$  the following mapping:

$$y \leftarrow \frac{\text{softplus}(y - \bar{y})}{\text{softplus}(y_{\max} - \bar{y})} * 4 - 2$$

where  $\bar{y}$  is the median of  $\{y_{j(i)}\}_{j=1}^{M_i}$ .

In this section, we set HyperBO variants and STBO to share exactly the same GP-UCB acquisition function as STBOH, MIMO and RFGP. The UCB coefficient for all methods is 1.8. The variants of HyperBO are as follows:

- H\* NLL: HyperBO with UCB as the acquisition function and negative log marginal likelihood (NLL) as the objective function.
- H\* NLLKL: HyperBO with UCB as the acquisition function and NLL plus 10 times KL divergence on matching datapoints as the objective function. See §A for more details.

In Fig. 8, we include the online tuning results for selected tasks due to limited compute resources. We noticed that for some methods, e.g. STBO and MIMO, it is very difficult for them to recover from a “bad” datapoint. This is partly because predictions from these models are significantly tied to the initial observations. For example, STBO may overfit to the initial bad value and believe there are bad values in the entire search space. Nevertheless, in 7 out of 9 tasks, HyperBO methods performed the best among all methods being compared.

## 6. Discussion

In this work, we focused on the question of how to efficiently and effectively make use of multi-task data to enable better BO with pre-trained priors. We simplified other aspects of BO that are orthogonal to our focuses, such as parallel queries or different search spaces. Here we discuss extensions to our work that would enable even more flexible uses.

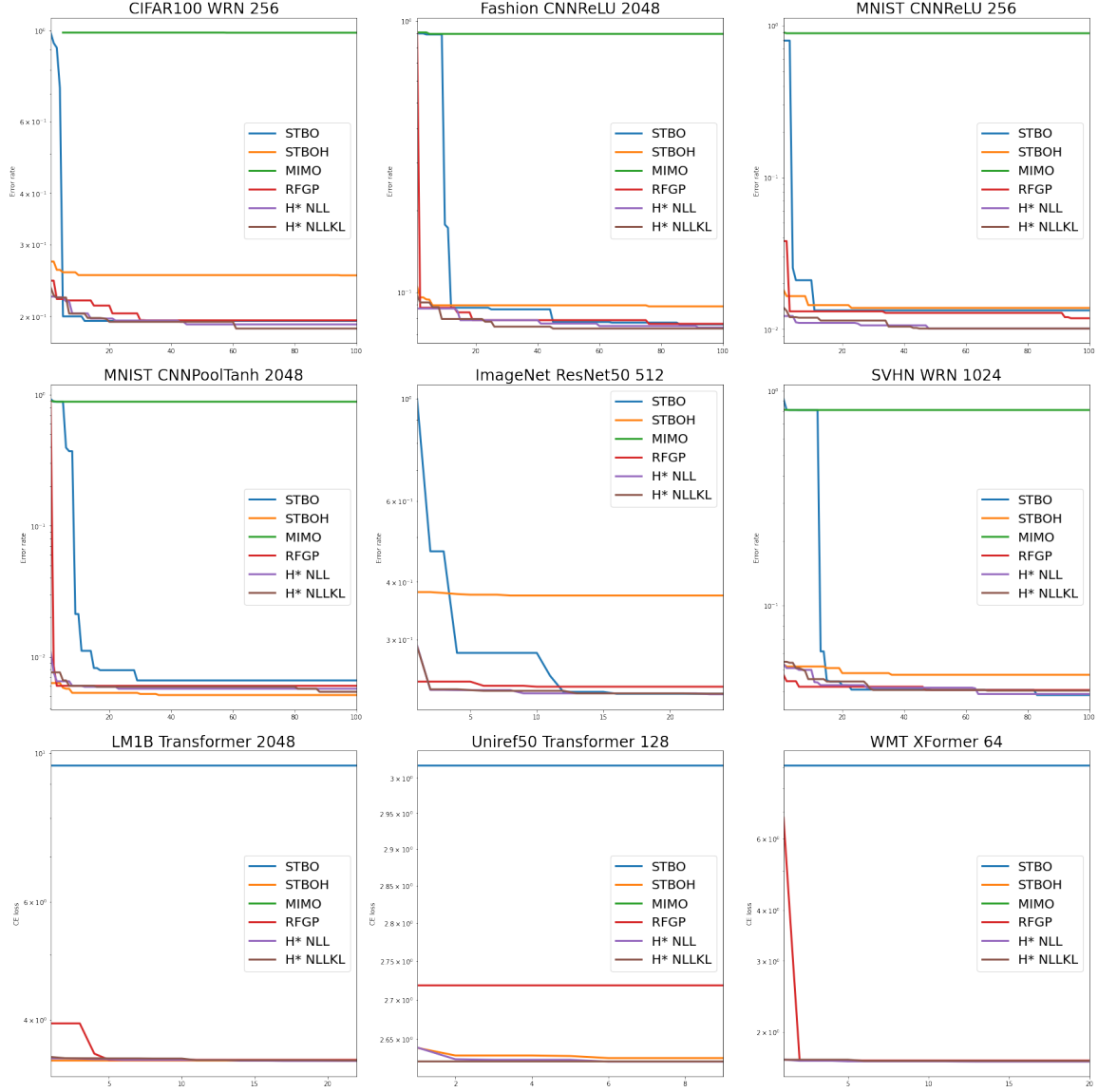


Figure 8: Results of running BO methods in the online setting on 9 different tasks. The image based tasks all use best validation error rate as objective while the text based tasks (LM1B, Uniref50 and WMT) use best validation ce loss. HyperBO methods achieved better results in 7 out of 9 tasks.

**Batch evaluation.** For simplicity of this paper, we did not consider batch evaluation but rather only focused on the prior selection dimension of the challenges in BO. However, it is straightforward to adopt any batch BO methods in conjunction with HyperBO to support obtaining observations in parallel. For example, we can directly use batch methods from Snoek et al. (2012); Kathuria et al. (2016); Wang et al. (2017b) etc. to replace line 5 of Alg. 1.

**High-dimensional and large scale data.** Similar to batch BO, our method can also be naturally combined with most high-dimensional and large scale BO methods to offer more capabilities. For these cases, typically a probabilistic model different from vanilla GPs may be adopted. In line 2 of Alg. 1, it is straightforward to adapt our method to optimize the cumulative marginal likelihood in Eq. 4 instead for the new model. Our meta-learning idea in this paper in fact also brings benefit to high-dimensional and large scale BO methods so that they can better identify their critical special structures, e.g. low-dimensional embedding (Wang et al., 2016), cylindrical kernels (Oh et al., 2018) or additive Mondrian kernels (Wang et al., 2018a).

**Different search spaces.** Roughly speaking, there could be two circumstances for difference search spaces. Case I is that tasks share the same search variables, but the search ranges for some variables are different. For example, we may have each function  $f_i : \mathcal{X}_i \rightarrow \mathbb{R}, i \in [N]$  and  $\mathcal{X}_i = \prod_{j=1}^d [l_{ij}, h_{ij}] \subset \mathbb{R}^d$ . In this case, our solution still applies by simply setting a union search space as  $\mathcal{X} = \bigcup_{i=1}^N \mathcal{X}_i$  for learning and use the designated search space of new tasks for optimization.

Case II is more complicated: the search space for each function  $f_i$  is  $\mathcal{X}_i \subset \mathbb{R}^{d_i}$  and each dimension of  $\mathcal{X}_i$  may have a different meaning than another search space  $\mathcal{X}_j$  ( $i \neq j$ ). This paper does not have a solution for this scenario. Further research will be needed to reduce Case II to Case I which can be then immediately combined with HyperBO.

## 7. Conclusion

We proposed HyperBO: a novel meta BO approach that supports practical applications that involve continuous inputs queried at possibly non-aligned locations across tasks. HyperBO uses a simple yet effective idea that is easy to implement and efficient to run. We evaluated HyperBO on real-world big model optimizer tuning tasks, and the results demonstrated its superior performance over state-of-the-art competing methods.

## References

- Rémi Bardenet, Mátyás Brendel, Balázs Kégl, and Michele Sebag. Collaborative hyperparameter tuning. In *ICML*, 2013.
- J Baxter. A Bayesian/information theoretic model of bias learning. In *COLT*, New York, New York, USA, 1996.
- Xavier Bouthillier and Gaël Varoquaux. Survey of machine-learning experimental methods at NeurIPS2019 and ICLR2020. Research report, Inria Saclay Ile de France, January 2020. URL <https://hal.archives-ouvertes.fr/hal-02447823>.
- James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and

- Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>.
- Pavel Brazdil, João Gama, and Bob Henery. Characterizing the applicability of classification algorithms using meta-level learning. In *ECML*, 1994.
- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- Alexander I Cowen-Rivers, Wenlong Lyu, Rasul Tutunov, Zhi Wang, Antoine Grosnit, Ryan Rhys Griffiths, Alexandre Max Maraval, Hao Jianye, Jun Wang, Jan Peters, et al. An empirical study of assumptions in bayesian optimisation. *arXiv preprint arXiv:2012.03826*, 2020.
- Joshua V Dillon, Ian Langmore, Dustin Tran, Eugene Brevdo, Srinivas Vasudevan, Dave Moore, Brian Patton, Alex Alemi, Matt Hoffman, and Rif A Saurous. Tensorflow distributions. *arXiv preprint arXiv:1711.10604*, 2017.
- Danny Drieß, Peter Englert, and Marc Toussaint. Constrained bayesian optimization of combined interaction force/task space controllers for manipulations. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 902–907. IEEE, 2017.
- Matthias Feurer, Aaron Klein, Katharina Eggensperger, Jost Springenberg, Manuel Blum, and Frank Hutter. Efficient and robust automated machine learning. In *NeurIPS*, 2015.
- Justin M. Gilmer, George E. Dahl, and Zachary Nado. init2winit: a jax codebase for initialization, optimization, and tuning research, 2021. URL <http://github.com/google/init2winit>.
- Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 580–587, 2014.
- Daniel Golovin, Benjamin Solnik, Subhodeep Moitra, Greg Kochanski, John Elliot Karro, and D. Sculley. Google vizier: A service for black-box optimization. In *KDD*, 2017.
- Ryan-Rhys Griffiths and José Miguel Hernández-Lobato. Constrained bayesian optimization for automatic chemical design using variational autoencoders. *Chemical Science*, 11(2):577–586, 2020.
- Nikolaus Hansen, Anne Auger, Raymond Ros, Olaf Mersmann, Tea Tušar, and Dimo Brockhoff. Coco: A platform for comparing continuous optimizers in a black-box setting. *Optimization Methods and Software*, 36(1):114–144, 2021. URL <https://arxiv.org/pdf/1603.08785.pdf>.
- Marton Havasi, Rodolphe Jenatton, Stanislav Fort, Jeremiah Zhe Liu, Jasper Snoek, Balaji Lakshminarayanan, Andrew M Dai, and Dustin Tran. Training independent subnetworks for robust prediction. *arXiv preprint arXiv:2010.06610*, 2020.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.



- Jonathan Heek, Anselm Levskaya, Avital Oliver, Marvin Ritter, Bertrand Rondepierre, Andreas Steiner, and Marc van Zee. Flax: A neural network library and ecosystem for JAX, 2020. URL <http://github.com/google/flax>.
- Matteo Hessel, David Budden, Fabio Viola, Mihaela Rosca, Eren Sezener, and Tom Hennigan. Optax: composable gradient transformation and optimisation, in *jax!*, 2020. URL <http://github.com/deepmind/optax>.
- Tarun Kathuria, Amit Deshpande, and Pushmeet Kohli. Batched gaussian process bandit optimization via determinantal point processes. *NeurIPS*, 2016.
- Charles Kemp and Joshua B Tenenbaum. The discovery of structural form. *Proceedings of the National Academy of Sciences*, 105(31):10687–10692, 2008.
- Beomjoon Kim, Leslie Pack Kaelbling, and Tomás Lozano-Pérez. Learning to guide task and motion planning using score-space representation. In *ICRA*, 2017.
- Beomjoon Kim, Zi Wang, Leslie Pack Kaelbling, and Tomás Lozano-Pérez. Learning to guide task and motion planning using score-space representation. *The International Journal of Robotics Research*, 38(7):793–812, 2019.
- Samuel Kim, Peter Y Lu, Charlotte Loh, Jamie Smith, Jasper Snoek, and Marin Soljačić. Scalable and flexible deep bayesian optimization with auxiliary information for scientific problems. *arXiv preprint arXiv:2104.11667*, 2021.
- Lars Kotthoff, Chris Thornton, Holger H Hoos, Frank Hutter, and Kevin Leyton-Brown. Auto-weka: Automatic model selection and hyperparameter optimization in weka. In *Automated Machine Learning*, pages 81–95. Springer, Cham, 2019.
- Andreas Krause and Cheng S Ong. Contextual Gaussian process bandit optimization. In *NeurIPS*, 2011.
- Rémi Lam, Matthias Poloczek, Peter Frazier, and Karen E Willcox. Advances in bayesian optimization with applications in aerospace engineering. In *AIAA Non-Deterministic Approaches Conference*, page 1656, 2018.
- Lisha Li, Kevin Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. Hyperband: A novel bandit-based approach to hyperparameter optimization. *The Journal of Machine Learning Research*, 18(1):6765–6816, 2017.
- Dong C Liu and Jorge Nocedal. On the limited memory bfgs method for large scale optimization. *Mathematical Programming*, 45(1):503–528, 1989.
- Jeremiah Zhe Liu, Zi Lin, Shreyas Padhy, Dustin Tran, Tania Bedrax-Weiss, and Balaji Lakshminarayanan. Simple and principled uncertainty estimation with deterministic deep learning via distance awareness. *arXiv preprint arXiv:2006.10108*, 2020.
- Gustavo Malkomes and Roman Garnett. Automating Bayesian optimization with Bayesian optimization. *Advances in Neural Information Processing Systems*, 31:5984–5994, 2018.

- Yurii E Nesterov. A method for solving the convex programming problem with convergence rate  $O(1/k^2)$ . In *Dokl. akad. nauk Sssr*, volume 269, pages 543–547, 1983.
- ChangYong Oh, Efstratios Gavves, and Max Welling. Bock: Bayesian optimization with cylindrical kernels. In *ICML*, 2018.
- Valerio Perrone, Rodolphe Jenatton, Matthias Seeger, and Cédric Archambeau. Scalable hyperparameter transfer learning. In *NeurIPS*, pages 6846–6856, 2018.
- Matthias Poloczek, Jialei Wang, and Peter I Frazier. Warm starting Bayesian optimization. In *Winter Simulation Conference (WSC)*. IEEE, 2016.
- Matthias Poloczek, Jialei Wang, and Peter Frazier. Multi-information source optimization. In *NeurIPS*, 2017.
- Edward O Pyzer-Knapp. Bayesian optimization for accelerated drug discovery. *IBM Journal of Research and Development*, 62(6):2–1, 2018.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683*, 2019.
- Ali Rahimi, Benjamin Recht, et al. Random features for large-scale kernel machines. In *NeurIPS*, 2007.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015.
- David Salinas, Huibin Shen, and Valerio Perrone. A quantile-based approach for hyperparameter transfer learning. In *ICML*, 2020.
- Benjamin J Shields, Jason Stevens, Jun Li, Marvin Parasram, Farhan Damani, Jesus I Martinez Alvarado, Jacob M Janey, Ryan P Adams, and Abigail G Doyle. Bayesian reaction optimization as a tool for chemical synthesis. *Nature*, 590(7844):89–96, 2021.
- Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical Bayesian optimization of machine learning algorithms. In *NeurIPS*, 2012.
- Jasper Snoek, Oren Rippel, Kevin Swersky, Ryan Kiros, Nadathur Satish, Narayanan Sundaram, Mostofa Patwary, Mr Prabhat, and Ryan Adams. Scalable bayesian optimization using deep neural networks. In *International conference on machine learning*, pages 2171–2180. PMLR, 2015.
- Jost Tobias Springenberg, Aaron Klein, Stefan Falkner, and Frank Hutter. Bayesian optimization with robust bayesian neural networks. *Advances in Neural Information Processing Systems*, 29: 4134–4142, 2016.
- Niranjan Srinivas, Andreas Krause, Sham M Kakade, and Matthias Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. In *ICML*, 2010.

- Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *ICML*, 2013.
- Kevin Swersky, Jasper Snoek, and Ryan P Adams. Multi-task Bayesian optimization. In *NeurIPS*, 2013.
- 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. *arXiv preprint arXiv:2104.10201*, 2021.
- 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 *International Conference on Learning Representations (ICLR)*, 2020.
- Zi Wang, Stefanie Jegelka, Leslie Pack Kaelbling, and Tomás Lozano-Pérez. Focused model-learning and planning for non-Gaussian continuous state-action systems. In *ICRA*, 2017a.
- Zi Wang, Chengtao Li, Stefanie Jegelka, and Pushmeet Kohli. Batched high-dimensional Bayesian optimization via structural kernel learning. In *ICML*, 2017b.
- Zi Wang, Clement Gehring, Pushmeet Kohli, and Stefanie Jegelka. Batched large-scale bayesian optimization in high-dimensional spaces. In *AISTATS*, 2018a.
- Zi Wang, Beomjoon Kim, and Leslie Pack Kaelbling. Regret bounds for meta Bayesian optimization with an unknown gaussian process prior. In *NeurIPS*, 2018b.
- Ziyu Wang, Frank Hutter, Masrour Zoghi, David Matheson, and Nando de Freitas. Bayesian optimization in a billion dimensions via random embeddings. *Journal of Artificial Intelligence Research*, 55:361–387, 2016.
- Martin Wistuba and Josif Grabocka. Few-shot bayesian optimization with deep kernel surrogates. In *International Conference on Learning Representations (ICLR)*, 2021.
- Dani Yogatama and Gideon Mann. Efficient transfer learning method for automatic hyperparameter tuning. In *AISTATS*, 2014.

## Appendix A. Objective functions

In §4, we presented NLL and KL divergence as objectives. Below we derive the KL divergence between a regular multivariate Gaussian and a degenerate multivariate Gaussian, which is the case for most of our matching data settings in §5.1: the number of matching data points is greater than the number of training tasks. In the end of this section, we introduce a new objective function, combining NLL and KL, interpreted as MAP with a data-dependent prior.

**KL divergence for a degenerate multivariate Gaussian** Eq. 3 of §4.1 gives the KL divergence between two Gaussians in the non-degenerate case. In practice, when we minimize Eq. 3, we can simply remove the constants and do the following

$$\begin{aligned}\hat{\mu}, \hat{k}, \hat{\sigma}^2 &= \arg \min_{\mu, k, \sigma^2} \mathcal{D}_{\text{KL}} \left( \mathcal{N}(\tilde{\mu}, \tilde{K}), \mathcal{N}(\mu, K) \right) \\ &= \arg \min_{\mu, k, \sigma^2} \frac{1}{2} \left( \text{tr}(K^{-1} \tilde{K}) + (\mu - \tilde{\mu})^\top K^{-1} (\mu - \tilde{\mu}) + \ln \frac{|K|}{|\tilde{K}|} - M \right) \\ &= \arg \min_{\mu, k, \sigma^2} \text{tr}(K^{-1} \tilde{K}) + (\mu - \tilde{\mu})^\top K^{-1} (\mu - \tilde{\mu}) + \ln |K|. \end{aligned} \quad (7)$$

Here the variables we care about,  $\mu, k, \sigma^2$ , only appear in mean vector  $\mu$  and covariance matrix  $K$  over the matching data. Even if the sample mean and covariance estimate  $\mathcal{N}(\tilde{\mu}, \tilde{K})$  is degenerate, the optimization objective stays the same as reflected by the derivation below.

If  $\mathcal{N}(\tilde{\mu}, \tilde{K})$  is degenerate, its base measure is at most  $N$ -dimensional rather than  $M$ -dimensional, given that there exists a full rank matrix  $A \in \mathbb{R}^{M \times R}$  such that  $\tilde{K} = AA^\top$  ( $R \leq N$ ). Note that  $M$  is the number of matching data points,  $N$  the number of training tasks and  $R$  is the rank of matrix  $A$  and  $\tilde{K}$ . The KL divergence  $\mathcal{D}_{\text{KL}} \left( \mathcal{N}(\tilde{\mu}, \tilde{K}), \mathcal{N}(\mu, K) \right)$  is not well-defined because the base measure of  $\mathcal{N}(\tilde{\mu}, \tilde{K})$  is different from the base measure of  $\mathcal{N}(\mu, K)$ , given  $K$  is full-rank. However, it is still possible to derive a pseudo KL divergence  $\mathcal{D}_{\text{KL}}^* \left( \mathcal{N}(\tilde{\mu}, \tilde{K}), \mathcal{N}(\mu, K) \right)$  as below.

Let the degenerate Gaussian be  $p(x) = \mathcal{N}(\tilde{\mu}, \tilde{K}) = |2\pi \tilde{K}|_*^{-\frac{1}{2}} \exp \left( -\frac{1}{2} (x - \tilde{\mu}) \tilde{K}^+ (x - \tilde{\mu})^\top \right)$  and the non-degenerate one be  $q(x) = \mathcal{N}(\mu, K)$ , where  $|\cdot|_*$  is the pseudo-determinant and  $\tilde{K}^+$  the pseudo-inverse of  $\tilde{K}$ . We define the support of distribution  $p$  as  $S(p) = \{\tilde{\mu} + \tilde{K}^{\frac{1}{2}} v \mid v \in \mathbb{R}^M\}$ . The pseudo KL divergence between  $p(x)$  and  $q(x)$  now becomes

$$\begin{aligned}\mathcal{D}_{\text{KL}}^* \left( \mathcal{N}(\tilde{\mu}, \tilde{K}), \mathcal{N}(\mu, K) \right) &= \int_{S(p)} p(x) (\ln p(x) - \ln q(x)) \\ &= -\frac{1}{2} \int_{S(p)} p(x) \left( \ln |2\pi \tilde{K}|_* - \ln |2\pi K| + (x - \tilde{\mu})^\top \tilde{K}^+ (x - \tilde{\mu}) - (x - \mu)^\top K^{-1} (x - \mu) \right) \\ &= \frac{1}{2} \left( (M - R) \ln 2\pi + \ln \frac{|K|}{|A^\top A|} - \mathbb{E}_p[\text{tr}(\tilde{K}^+ (x - \tilde{\mu})(x - \tilde{\mu})^\top) + \text{tr}(K^{-1} (x - \mu)(x - \mu)^\top)] \right) \\ &= \frac{1}{2} \left( (M - R) \ln 2\pi + \ln \frac{|K|}{|A^\top A|} - \text{tr}(\tilde{K}^+ \tilde{K}) + \mathbb{E}_p[\text{tr}(K^{-1} (x - \mu)(x - \mu)^\top)] \right) \\ &= \frac{1}{2} \left( (M - R) \ln 2\pi + \ln \frac{|K|}{|A^\top A|} - \text{tr}(AA^+) + \mathbb{E}_p[\text{tr}(K^{-1} (x - \tilde{\mu})(x - \tilde{\mu})^\top + K^{-1} (2x\tilde{\mu}^\top - 2x\mu^\top - \tilde{\mu}\tilde{\mu}^\top + \mu\mu^\top))] \right) \\ &= \frac{1}{2} \left( (M - R) \ln 2\pi + \ln \frac{|K|}{|A^\top A|} - R + \text{tr}(K^{-1} \tilde{K} + K^{-1} (\tilde{\mu} - \mu)(\tilde{\mu} - \mu)^\top) \right) \\ &= \frac{1}{2} \left( \text{tr}(K^{-1} \tilde{K}) + (\mu - \tilde{\mu})^\top K^{-1} (\mu - \tilde{\mu}) + \ln |K| - \ln |A^\top A| - R + (M - R) \ln 2\pi \right). \end{aligned} \quad (8)$$

If the covariance matrix  $\tilde{K}$  is in fact full rank, i.e.  $R = M$ , pseudo KL in Eq. 8 then recovers the KL divergence  $\mathcal{D}_{\text{KL}} \left( \mathcal{N}(\tilde{\mu}, \tilde{K}), \mathcal{N}(\mu, K) \right)$  in Eq. 3 for non-degenerate Gaussians. If we were to minimize this pseudo KL divergence  $\mathcal{D}_{\text{KL}}^* \left( \mathcal{N}(\tilde{\mu}, \tilde{K}), \mathcal{N}(\mu, K) \right)$ , we still would get the minimization task in Eq. 7. However, pseudo KL divergence  $\mathcal{D}_{\text{KL}}^*$  does not satisfy properties including non-negativity. In practice, we may also choose to add small epsilon values to the diagonal terms of both  $K$  and  $\tilde{K}$  to make  $\tilde{K}$  “less degenerate” and enable the existence of  $\mathcal{D}_{\text{KL}} \left( \mathcal{N}(\tilde{\mu}, \tilde{K}), \mathcal{N}(\mu, K) \right)$ .

**Combining marginal likelihood and empirical divergence** It is also possible to additively combine the KL divergence with the negative log marginal likelihood objective, and treat the empirical divergence as a regularizer. In the case of KL divergence, it is equivalent to adding a data-dependent prior on the GP itself:  $\mu, K \sim \exp(-\lambda \mathcal{D}_{\text{KL}}(\mathcal{N}(\tilde{\mu}, \tilde{K}), \mathcal{N}(\mu, K)))/Z$  for some normalization constant  $Z$ , and the posterior is

$$p(\mu, k, \sigma^2 \mid D_N; \tilde{\mu}, \tilde{K}) \propto p(D_N \mid \mu, k, \sigma^2) \exp(-\lambda \mathcal{D}_{\text{KL}}(\mathcal{N}(\tilde{\mu}, \tilde{K}), \mathcal{N}(\mu, K))). \quad (9)$$

We can then obtain an MAP estimation for the unknown functions and variables  $\mu, k, \sigma^2$ .

## Appendix B. Details of regret bounds

Theorem 2 is a direct result of Theorem 16 in Wang et al. (2018b). The only subtle difference is that we adopted a biased estimate of the covariance matrix, a factor of  $\frac{N}{N-1}$  different from the unbiased estimate. But this can be corrected in the acquisition functions. We provide more details below.

By Proposition 1, we are able to obtain a  $\mathcal{GP}(\mu, k)$  where  $\mu(x)$  and  $k(x)$  are equal to the sample mean and covariance on a matching dataset  $(x, y)$ , following the notations in §4.1. Hence, our problem setup is consistent with the case with discrete input space in Wang et al. (2018b). The following theorem is a rewrite of Theorem 16 in Wang et al. (2018b), taking into account our biased estimators.

**Theorem 3.** Assume there exist constant  $c \geq \max_{x \in \mathcal{X}} k(x)$  and a training dataset is available whose size is  $N \geq 4 \log \frac{6}{\delta} + T + 2$ . Define

$$\iota_{t-1} = \sqrt{\frac{6 \left( N - 3 + t + 2\sqrt{t \log \frac{6}{\delta}} + 2 \log \frac{6}{\delta} \right)}{\delta N(N - t - 1)}}, \quad b_{t-1} = \frac{1}{N - t} \log \frac{6}{\delta}, \quad \text{for any } t \in [T],$$

and  $\rho_T = \max_{A \in \mathcal{X}, |A|=T} \frac{1}{2} \log |\mathbf{I} + \sigma^{-2} k(A)|$ . Then, with probability at least  $1 - \delta$ , the best-sample simple regret in  $T$  iterations of meta BO with GP-UCB that uses

$$\zeta_t = \frac{\left( 6N(N - 3 + t + 2\sqrt{t \log \frac{6}{\delta}} + 2 \log \frac{6}{\delta}) / (\delta N(N - t - 1)) \right)^{\frac{1}{2}} + (2N \log(\frac{3}{\delta}))^{\frac{1}{2}}}{\left( (N - 1)(1 - 2(\frac{1}{N-t} \log \frac{6}{\delta})^{\frac{1}{2}}) \right)^{\frac{1}{2}}} \quad (10)$$

as its hyperparameter in  $\alpha_{t-1}^{\text{GP-UCB}}(x) = \hat{\mu}_{t-1}(x) + \zeta_t \hat{k}_{t-1}(x)^{\frac{1}{2}}$  satisfies

$$r_T^{\text{GP-UCB}} \leq \eta^{\text{GP-UCB}} \sqrt{\frac{2c\rho_T}{T \log(1 + c\sigma^{-2})} + \sigma^2} - \frac{(2 \log(\frac{3}{\delta}))^{\frac{1}{2}} \sigma^2}{\sqrt{c + \sigma^2}},$$

where  $\eta^{\text{GP-UCB}} = \left( \frac{\iota_{T-1} + (2 \log(\frac{3}{\delta}))^{\frac{1}{2}}}{\sqrt{1 - 2\sqrt{b_{T-1}}}} \sqrt{1 + 2\sqrt{b_{T-1}} + 2b_{T-1}} + \iota_{T-1} + (2 \log(\frac{3}{\delta}))^{\frac{1}{2}} \right)$ .

With probability at least  $1 - \delta$ , the best-sample simple regret in  $T$  iterations of meta BO with PI that uses  $\hat{f}^* \geq \max_{x \in \mathcal{X}} f(x)$  as its target value satisfies

$$r_T^{\text{PI}} < \eta^{\text{PI}} \sqrt{\frac{2c\rho_T}{T \log(1 + c\sigma^{-2})} + \sigma^2} - \frac{(2 \log(\frac{3}{2\delta}))^{\frac{1}{2}} \sigma^2}{2\sqrt{c + \sigma^2}},$$

where  $\eta^{PI} = (\frac{\hat{f}^* - \mu_{\tau-1}(x_*)}{\sqrt{k_{\tau-1}(x_*) + \sigma^2}} + \iota_{\tau-1}) \sqrt{\frac{1 + 2b_{\tau-1}^{\frac{1}{2}} + 2b_{\tau-1}}{1 - 2b_{\tau-1}^{\frac{1}{2}}}} + \iota_{\tau-1} + (2 \log(\frac{3}{2\delta}))^{\frac{1}{2}}$ ,  $\tau = \arg \min_{t \in [T]} k_{t-1}(x_t)$ .

The proof can be found in Wang et al. (2018b). Theorem 2 is a condensed version of Theorem 3.

While Theorem 3 provides us with some understanding of HyperBO in a specific setting, in practice, we need to query in a continuous input space that goes beyond the finite set of points present in the training dataset. It may or may not be possible to obtain data on a wide range of tasks to ensure  $N \geq 4 \log \frac{6}{\delta} + T + 2$ . In fact, in all of our experiment, this criterion on number of tasks is not satisfied. However, we still obtained good performance.

## Appendix C. Experiment details and more results

In this section, we provide more empirical results on the impact of objective functions and acquisition functions in HyperBO. All experiment setups are the same as §5.3.1: offline and holding out related tasks.

### C.1 Impact of objective functions

Here we investigate how different objective functions in HyperBO can impact its performance. Besides NLL and KL, which are already described in details in §4, we also include NLL+KL, which corresponds to Eq. 9 with  $\lambda = 10$ .  $\lambda = 10$  is an arbitrary choice, and one may find other better options to set  $\lambda$  in Eq. 9.

Figure 9 shows the performance profiles and BO simple regret curves of NLL, KL and NLL+KL when HyperBO uses different acquisition functions. As comparisons, we also included the better performing baselines: Rand, STBOH and MIMO. While the ranks of performance by different objectives do vary depending on which acquisition function we used, it is clear that all HyperBO variants outperform the baselines. For EI and PI, the KL objective gives better performance, but NLL or NLL+KL might be preferred for UCB with coefficient 2, 3 or 4.

### C.2 Impact of acquisition functions

As explained briefly in §5, we used 5 acquisition functions in our experiments: the vanilla EI method, PI with coefficient 0.1,  $\alpha^{PI}(x; \mathcal{GP}(\hat{\mu}, \hat{k} \mid D_f)) = \frac{\hat{\mu}_{D_f}(x) - \max_t(y_t + 0.1)}{\hat{\sigma}_{D_f}(x)}$ , and UCB with coefficient  $\zeta = 2, 3$  and 4 in  $\alpha^{UCB}(x; \mathcal{GP}(\hat{\mu}, \hat{k} \mid D_f)) = \hat{\mu}_{D_f}(x) + \zeta \hat{\sigma}_{D_f}(x)$ .

Our goal is to verify that HyperBO maintains good performance across different choices of acquisition functions. To do so, we investigated in the performance of HyperBO variants under different objectives. We avoid over cluttering the figures by only including STBOH as baseline, since it is roughly the best baseline according to the main results in Fig. 2.

As shown in Figure 10, HyperBO with EI and PI as acquisition functions perform relatively better than HyperBO with UCB variants. However, HyperBO with UCB3 can still be very competitive when it is coupled with NLL objective. Overall, HyperBO with all of the 5 acquisition function options outperforms the best performing baselines.



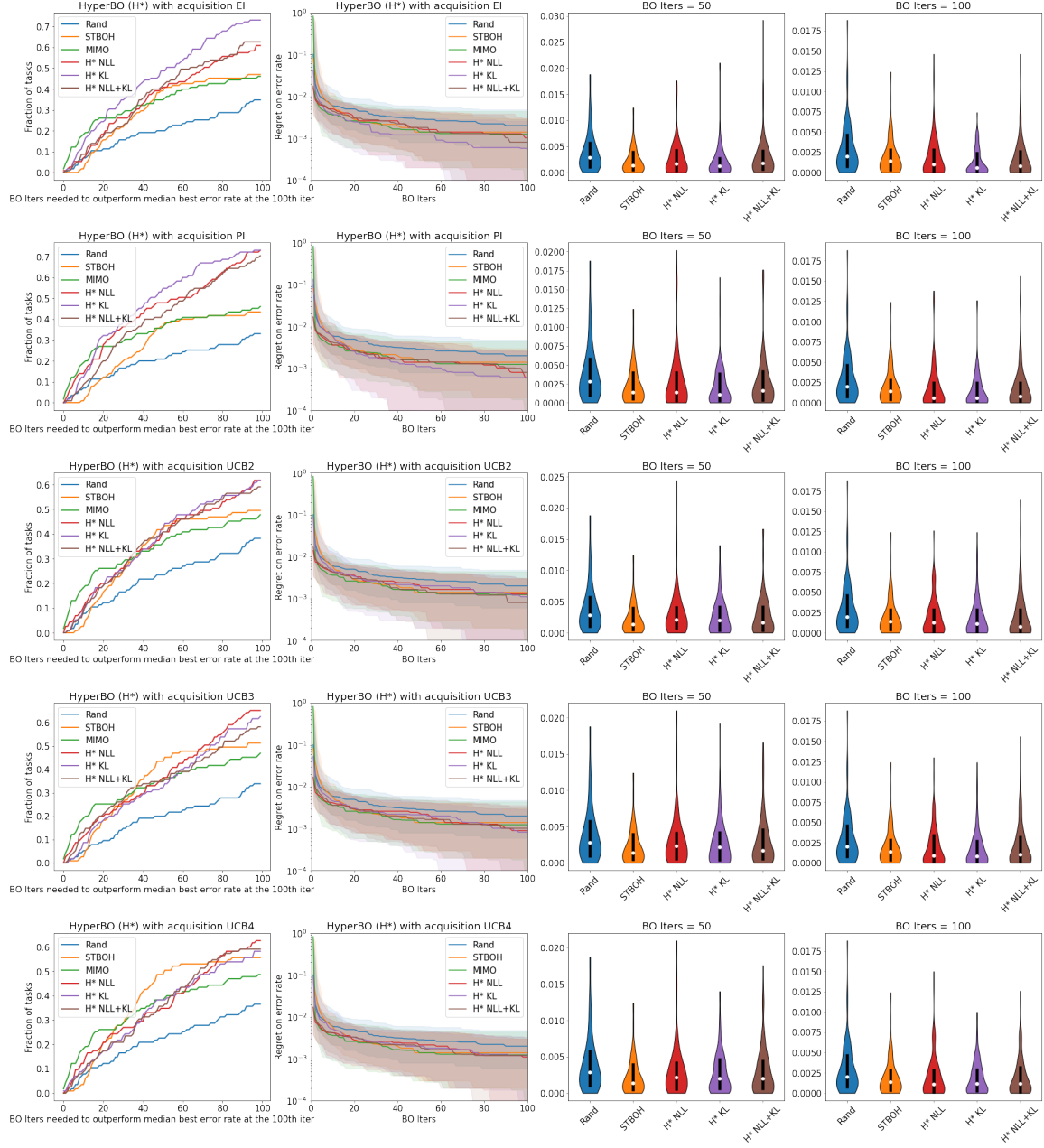


Figure 9: We compare the performance of 3 different objective functions in HyperBO under 5 settings of acquisition functions. For EI and PI, using KL as the objective for HyperBO is slightly better than NLL or NLL+KL. However, different conclusions can be drawn for UCB2, UCB3 and UCB4. Nevertheless, all HyperBO variants still outperform the best alternatives.

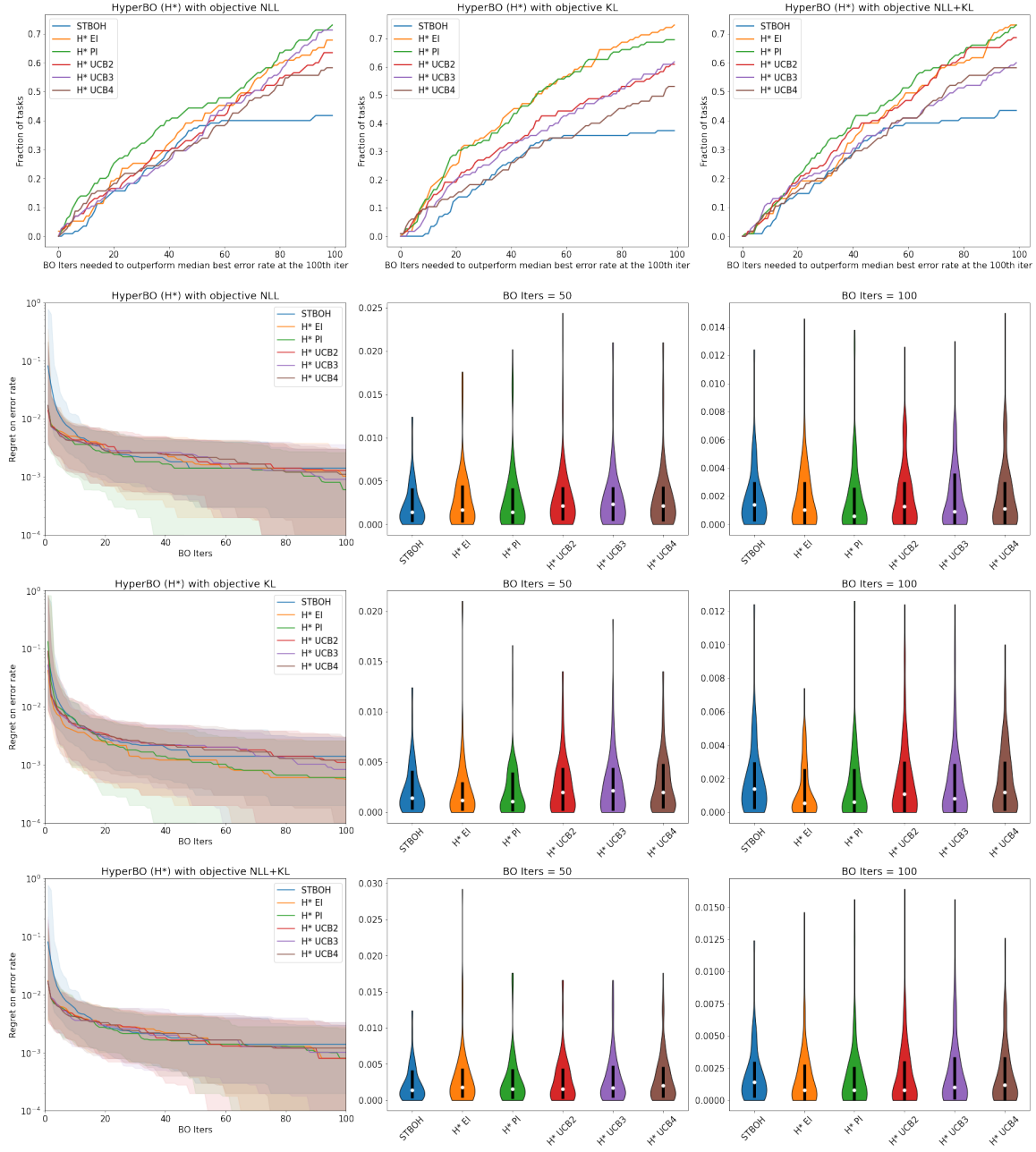


Figure 10: We compare the performance of 5 different acquisition functions under 3 settings of objectives in HyperBO. Overall, PI and EI outperform UCB with different coefficient values. But HyperBO with UCB variants still outperforms STBOH, which is roughly the best baseline according to the main results in Fig. 2.

### C.3 MAF implementation details

We compared to (Volpp et al., 2020) using the code and default hyperparameters provided by the authors.<sup>6</sup> This method assumes the availability of the optimal set of GP hyperparameters for each task (including the task used for evaluation). Following Volpp et al. (2020), these GP hyperparameters for the MAF algorithm are learned by optimizing the marginal likelihood on each training and evaluation task using the GPY library. Given that MAF takes significantly longer to run than HyperBO and other baselines, subdataset in each task was evaluated using limited random seeds.

Each neural acquisition function was trained for a total of 1000 iterations. As was done in (Volpp et al., 2020), we selected the optimal training iteration for the neural acquisition function by cross-validation on the transfer learning tasks; in this case, we randomly sampled 3 transfer learning task, and chose the training iteration with the lowest average simple regret.

Finally, to make use of the MAF code, we also had to ensure that (a) each subtask had the same number of evaluation points, and (b) that there were no duplicated tuning parameters. For this reason, we first removed all duplicate hyperparameters within each subdataset, then capped each subdataset to the first 1559 points (the size of the smallest sub-dataset) while retaining the best possible data point.

---

6. <https://github.com/boschresearch/MetaBO>