Bayesian Kernelized Tensor Factorization as Surrogate for Bayesian Optimization

Anonymous Author(s) Affiliation Address email

Abstract

1	Bayesian optimization (BO) mainly uses Gaussian processes (GP) with a stationary
2	and separable kernel function (e.g., the squared-exponential kernel with automatic
3	relevance determination [SE-ARD]) as the surrogate model. However, such lo-
4	calized kernel specifications are deficient in learning complex functions that are
5	non-stationary, non-separable and multi-modal. In this paper, we propose using
6	Bayesian Kernelized Tensor Factorization (BKTF) as a new surrogate model for
7	Bayesian optimization (BO) in a D-dimensional grid with both continuous and
8	categorical variables. Our key idea is to approximate the underlying D-dimensional
9	solid with a fully Bayesian low-rank tensor CP decomposition, in which we place
10	GP priors on the latent basis functions for each dimension to encode local consis-
11	tency and smoothness. With this formulation, the information from each sample
12	can be shared not only with neighbors but also across dimensions, thus fostering a
13	more global search strategy. Although BKTF no longer has an analytical posterior,
14	we efficiently approximate the posterior distribution through Markov chain Monte
15	Carlo (MCMC). We conduct numerical experiments on several test functions with
16	continuous variables and two machine learning hyperparameter tuning problems
17	with mixed variables. The results show that BKTF offers a flexible and highly
18	effective approach to characterizing and optimizing complex functions, especially
19	in cases where the initial sample size and budget are severely limited.

20 **1** Introduction

For many applications in science and engineering, such as emulation-based studies, experiment 21 design, and automated machine learning, the goal is to optimize a complex black-box function f(x)22 in a D-dimensional space, for which we have limited prior knowledge. The main challenge in 23 such optimization problems is that we aim to efficiently find the global optima, while the objective 24 function f is often gradient-free, multimodal and computationally expensive to evaluate. Bayesian 25 optimization (BO) offers a powerful statistical approach to these problems, particularly when the 26 observation budgets are limited [1, 2, 3]. A typical BO framework consists of two components—a 27 surrogate model and an acquisition function (AF)-to balance exploitation and exploration. The 28 surrogate is a probabilistic model that allows us to estimate f(x) with uncertainty at a new location 29 \boldsymbol{x} , and the AF is used to determine which location to query next. 30

Gaussian process (GP) regression is the most widely used surrogate for BO [3, 4], thanks to its
appealing properties in providing analytical derivations and uncertainty quantification (UQ). The
choice of kernel/covariance function is a critical decision in GP models; for multidimensional
BO problems, perhaps the most popular kernel is the ARD (automatic relevance determination)—
Squared-Exponential (SE) or Matérn kernel [4]. Although this specification has certain numerical
advantages and can help automatically learn the importance of input variables, a key limitation is



Figure 1: BO for a 2D function: (a) True function surface, where the global maximum is marked; (b) Comparison between BO models using GP surrogates (with two AFs) and BKTF with 30 random initial observations, averaged over 20 replications; (c) Specific results of one run, including the final estimated mean surface for f, in which green dots denote the locations of the selected candidates, and the corresponding AF surface.

that it implies/assumes that the underlying stochastic process is stationary and separable, and the 37 value of the covariance function between two random points quickly goes to zero with the increase 38 of input dimensionality. These assumptions can be problematic for complex real-world processes 39 with long-range dependencies, because estimating the underlying function with a simple ARD kernel 40 would require a large number of observations. A potential solution to address this issue is to use 41 more flexible kernel structures. The additive kernel, for example, is designed to characterize a more 42 'global'' structure by restricting variable interactions [5]. However, in practice using additive kernels 43 requires strong prior knowledge to determine the proper interactions and involves a large number 44 of kernel hyperparameters to learn [6]. Another emerging solution is to use deep GP [7], such as in 45 46 [8, 9]; however, learning deep GP often becomes a more challenging task due to the inference of latent layers. In addition, these GP related surrogates can be more deficient to tune when taking into 47 account both continuous and categorical inputs. 48

In this paper, we propose using *Bayesian Kernelized Tensor Factorization* (BKTF) as a flexible 49 and adaptive surrogate model for BO in a D-dimensional Cartesian product space (i.e., grid) when 50 D is relatively small (say $D \leq 10$). BKTF is initially developed for modeling multidimensional 51 52 spatiotemporal data with UQ, for tasks such as spatiotemporal kriging/cokriging [10, 11]. This paper adapts BKTF to the BO setting, and our key idea is to characterize the multivariate objective 53 function $f(\mathbf{x}) = f(x_1, \dots, x_D)$ for a specific BO problem using the low-rank tensor CANDE-54 COMP/PARAFAC (CP) factorization with random basis functions. Unlike other basis-function 55 models that rely on known/deterministic basis functions [12], BKTF uses a hierarchical Bayesian 56 framework to achieve high-quality UQ in a more flexible way-GP priors are used to model the basis 57 functions, and hyperpriors are used to model kernel hyperparameters in particular for the lengthscale 58 that characterizes the scale of variation. In addition, BKTF also provides a natural solution for 59 60 categorical variables, for which we can simply introduce an inverse-Wishart prior on the covariance matrix of the basis functions. 61

Figure 1 shows the comparison between BKTF and GP surrogates when optimizing a bivariate 62 function (D = 2) that is nonstationary, nonseparable, and multimodal. The details of this function 63 and the BO experiments are provided in Appendix C. This 2D case clearly shows that GP surrogate 64 is limited by the local kernel and becomes ineffective in finding the global solution, while BKTF 65 offers superior flexibility and adaptability to characterize the multidimensional process from limited 66 data. Unlike GP-based surrogate models, BKTF no longer has an analytical posterior; however, 67 efficient inference and acquisition can be achieved through Markov chain Monte Carlo (MCMC) 68 in an element-wise learning way, in which we update basis functions and kernel hyperparameters 69 using Gibbs sampling and slice sampling, respectively [11]. For optimization, we first use MCMC 70 samples to approximate the posterior distribution of the entire tensor and then naturally define the 71 upper confidence bound (UCB) of the posterior as AF. This process is feasible for many real-world 72 applications that can be studied in a discretized tensor product space, such as experimental design. 73

74 We conduct extensive experiments on both standard optimization and ML hyperparameter tuning

tasks. Our results show that BKTF achieves a fast global search for optimizing complex objective

⁷⁶ functions with limited initial data and budget.

77 2 Bayesian optimization

⁷⁸ Let $f : \mathcal{X} = \mathcal{X}_1 \times \ldots \times \mathcal{X}_D \to \mathbb{R}$ be a black-box function that could be nonconvex, derivative-free, ⁷⁹ and expensive to evaluate. BO aims to address the global optimization problem:

$$\boldsymbol{x}^{\star} = \operatorname*{arg\,max}_{\boldsymbol{x} \in \mathcal{X}} f(\boldsymbol{x}). \tag{1}$$

BO solves this problem by first building a probabilistic model for f(x) (i.e., surrogate model) based on initial observations and then using the predictive distribution to decide where in \mathcal{X} to evaluate/query next. The overall goal of BO is to find the global optimum of the objective function using as few evaluations as possible. Most BO models rely on a GP prior for f(x) to achieve prediction and UQ:

$$f(\boldsymbol{x}) = f(x_1, \dots, x_D) \sim \mathcal{GP}(m(\boldsymbol{x}), k(\boldsymbol{x}, \boldsymbol{x}')), \text{ with } x_d \in \mathcal{X}_d, \ d = 1, \dots, D,$$
(2)

where $k(\cdot, \cdot)$ is a valid kernel/covariance function and $m(\cdot)$ is a mean function that can be generally assumed to be 0. Given a finite set of observation points $\{\boldsymbol{x}_i\}_{i=1}^n$ with $\boldsymbol{x}_i = (x_1^i, \dots, x_D^i)^{\top}$, the vector of function values $\boldsymbol{f} = (f(\boldsymbol{x}_1), \dots, f(\boldsymbol{x}_n))^{\top}$ has a multivariate Gaussian distribution $\boldsymbol{f} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{K})$, where \boldsymbol{K} is the $n \times n$ covariance matrix. For a set of observed data $\mathcal{D}_n = \{\boldsymbol{x}_i, y_i\}_{i=1}^n$ with i.i.d. Gaussian noise, i.e., $y_i = f(\boldsymbol{x}_i) + \epsilon_i$ where $\epsilon_i \sim \mathcal{N}(0, \tau^{-1})$, GP gives an analytical posterior distribution of $f(\boldsymbol{x})$ at an unobserved point \boldsymbol{x}^* :

$$f(\boldsymbol{x}^*) \mid \mathcal{D}_n \sim \mathcal{N}\left(\boldsymbol{k}_{\boldsymbol{x}^*\boldsymbol{X}} \left(\boldsymbol{K} + \tau^{-1} \boldsymbol{I}_n\right)^{-1} \boldsymbol{y}, \ \boldsymbol{k}(\boldsymbol{x}^*, \boldsymbol{x}^*) - \boldsymbol{k}_{\boldsymbol{x}^*\boldsymbol{X}} \left(\boldsymbol{K} + \tau^{-1} \boldsymbol{I}_n\right)^{-1} \boldsymbol{k}_{\boldsymbol{x}^*\boldsymbol{X}}^\top\right), \quad (3)$$

here $\boldsymbol{k}_{\boldsymbol{x}^*\boldsymbol{X}} = \begin{bmatrix} \boldsymbol{k}(\boldsymbol{x}^*, \boldsymbol{x}_1) & \boldsymbol{k}(\boldsymbol{x}^*, \boldsymbol{x}_1) \end{bmatrix}^\top$ and $\boldsymbol{y} = (\boldsymbol{y}_1, \dots, \boldsymbol{y}_n)^\top$

91 where $k_{x^*X} = [k(x^*, x_1), \dots, k(x^*, x_n)]^{\top}$ and $y = (y_1, \dots, y_n)^{\top}$.

Based on the posterior distributions of f, one can 92 compute an AF, denoted by $\alpha : \mathcal{X} \to \mathbb{R}$, for a 93 new candidate x^* and evaluate how promising x^* 94 is. In BO, the next query point is often determined 95 by maximizing a selected/predefined AF, i.e., $x_{n+1} =$ 96 $\arg \max_{\boldsymbol{x} \in \mathcal{X}} \alpha (\boldsymbol{x} \mid \mathcal{D}_n)$. Most AFs are based on the 97 predictive mean and variance; for example, a com-98 monly used AF is the expected improvement (EI) 99 [1]: 100

$$\alpha_{\mathrm{EI}}\left(\boldsymbol{x} \mid \mathcal{D}_{n}\right) = \sigma(\boldsymbol{x})\varphi\left(\frac{\Delta(\boldsymbol{x})}{\sigma(\boldsymbol{x})}\right) + \left|\Delta(\boldsymbol{x})\right|\Phi\left(\frac{\Delta(\boldsymbol{x})}{\sigma(\boldsymbol{x})}\right)$$
(4)

Algorithm 1: Basic BO process.

Input: Initial dataset \mathcal{D}_0 and a trained
surrogate model; total budget N.for $n = 1, \ldots, N$ doApproximate the posterior distribution of
f using the surrogate model based on
 \mathcal{D}_{n-1} ;Find next evaluation point \boldsymbol{x}_n by
optimizing the AF;
Augment data $\mathcal{D}_n = \mathcal{D}_{n-1} \cup \{\boldsymbol{x}_n, y_n\}$,
update the surrogate model.

where $\Delta(x) = \mu(x) - f_n^*$ is the expected difference between the proposed point x and the current best solution, $f_n^* = \max_{x \in \{x_i\}_{i=1}^n} f(x)$ denotes the best function value obtained so far; $\mu(x)$ and $\sigma(x)$ are the predictive mean and predictive standard deviation at x, respectively; and $\varphi(\cdot)$ and $\Phi(\cdot)$ denote the probability density function (PDF) and the cumulative distribution function (CDF) of standard normal, respectively. Another widely applied AF for maximization problems is the **upper confidence bound** (UCB) [13]:

$$\alpha_{\text{UCB}}\left(\boldsymbol{x} \mid \mathcal{D}_{n}, \beta\right) = \mu(\boldsymbol{x}) + \beta\sigma(\boldsymbol{x}), \tag{5}$$

where β is a tunable parameter that balances exploration and exploitation. The general BO procedure can be summarized as Algorithm 1.

109 3 BKTF for Bayesian optimization

110 3.1 Bayesian hierarchical model specification

Before introducing BKTF, we first construct a *D*-dimensional grid space corresponding to the search space \mathcal{X} , where \mathcal{X}_d could be continuous, integer-valued, or categorical. We define it on *D* sets $\{S_1, \ldots, S_D\}$ and denote the whole grid by $\prod_{d=1}^D S_d$: $S_1 \times \ldots \times S_D =$ 114 $\{(s_1, \ldots, s_D) \mid \forall d \in \{1, \ldots, D\}, s_d \in S_d\}$. For dimensions with integer-valued and categorical 115 input, we consider S_d the set of corresponding discrete values. For dimensions with continuous input, 116 the coordinate set S_d is formed by m_d interpolation points c_i^d that are distributed over the bounded 117 interval $\mathcal{X}_d = [a_d, b_d]$, i.e., $S_d = \{c_i^d\}_{i=1}^{m_d}$ with $c_i^d \in \mathcal{X}_d$. The size of S_d becomes $|S_d| = m_d$, and 118 size of the entire grid space is $\prod_{d=1}^{D} |S_d|$. Note that we do not restrict S_d to be uniformly distributed. 119 We assume the underlying function f as a stochastic process that is zero-centered. We randomly 120 sample an initial dataset including n_0 input-output data pairs $\mathcal{D}_0 = \{x_i, y_i\}_{i=1}^{n_0}$, where $\{x_i\}_{i=1}^{n_0}$ 121 are located in $\prod_{d=1}^{D} S_d$, and this yields an incomplete D-dimensional tensor $\mathcal{Y} \in \mathbb{R}^{|S_1| \times \cdots \times |S_D|}$ 122 with n_0 observed points. BKTF approximates the entire data tensor \mathcal{Y} by a kernelized tensor CP

123 decomposition:

$$\boldsymbol{\mathcal{Y}} = \sum_{r=1}^{R} \lambda_r \cdot \boldsymbol{g}_1^r \circ \boldsymbol{g}_2^r \circ \cdots \circ \boldsymbol{g}_D^r + \boldsymbol{\mathcal{E}},$$
(6)

where *R* is a pre-specified tensor CP rank, $\lambda = (\lambda_1, ..., \lambda_R)^{\top}$ denote weight coefficients that capture the magnitude/importance of each rank in the factorization, $g_d^r = [g_d^r(s_d) : s_d \in S_d] \in \mathbb{R}^{|S_d|}$ denotes the *r*th latent factor for the *d*th dimension, entries in \mathcal{E} are i.i.d. white noises following $\mathcal{N}(0, \tau^{-1})$. It should be particularly noted that both the coefficients $\{\lambda_r\}_{r=1}^R$ and the latent basis functions $\{g_1^r, \ldots, g_D^r\}_{r=1}^R$ are random variables. The function approximation for $\boldsymbol{x} = (x_1, \ldots, x_D)^{\top}$ is:

$$f(\boldsymbol{x}) = \sum_{r=1}^{R} \lambda_r g_1^r(x_1) \cdots g_D^r(x_D) = \sum_{r=1}^{R} \lambda_r \prod_{d=1}^{D} g_d^r(x_d).$$
(7)

For priors, we assume $\lambda_r \sim \mathcal{N}(0,1)$ for r = 1, ..., R and use a GP prior on the latent factors for dimension d with continuous input:

$$g_d^r(x_d) \mid l_d^r \sim \mathcal{GP}\left(0, k_d^r(x_d, x_d'; l_d^r)\right), \text{ for } r = 1, \dots, R,$$
(8)

where k_d^r is a valid kernel function. In this paper, we choose a Matérn 3/2 kernel $k_d^r(x_d, x'_d; l_d^r) = \sigma^2 \left(1 + \frac{\sqrt{3}|x_d - x'_d|}{l_d^r}\right) \exp\left(-\frac{\sqrt{3}|x_d - x'_d|}{l_d^r}\right)$. We fix the kernel variance of k_d^r as $\sigma^2 = 1$, and only learn the lengthscale hyperparameters l_d^r , since the variances of the model can be captured by λ . One can 131 132 133 also exclude λ but introduce variance σ^2 as a kernel hyperparameter on one of the basis functions; 134 however, learning kernel hyperparameters is computationally more expensive than learning λ . For 135 simplicity, we can also assume the lengthscale parameters to be identical, i.e., $l_d^1 = \ldots = l_d^R = l_d$, for each dimension d. The prior distribution for the corresponding latent factor \boldsymbol{g}_d^r is then a Gaussian distribution: $\boldsymbol{g}_d^r \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{K}_d^r)$, where \boldsymbol{K}_d^r is the $|S_d| \times |S_d|$ covariance matrix computed from k_d^r . 136 137 138 We place Gaussian hyperpriors on the log-transformed kernel hyperparameters to ensure positive 139 values, i.e., $\log(l_d^r) \sim \mathcal{N}(\mu_l, \tau_l^{-1})$. For categorical input, we assume that the corresponding latent 140 basis functions $g_d^r \mid \mathbf{\Lambda}_d \sim \mathcal{N}(\mathbf{0}, \mathbf{\Lambda}_d^{-1})$ for $r = 1, \dots, R$, where the precision matrix $\mathbf{\Lambda}_d$ follows a 141 Wishart prior with an identity scale matrix and $|S_d|$ degrees of freedom, i.e., $\Lambda_d \sim \mathcal{W}(I_{|S_d|}, |S_d|)$. 142 For noise precision τ , we assume a conjugate Gamma prior $\tau \sim \text{Gamma}(a_0, b_0)$. For dimensions 143 with integer variables, we could model the covariance of the basis functions either using a kernel 144 145 matrix or with an inverse-Wishart prior, depending on specific situations.

For observations, we assume each y_i in the initial dataset \mathcal{D}_0 follows a Gaussian distribution: $y_i = \{a_i^r, (r^i)\} \{\lambda_i\} \ \tau \sim \mathcal{N}(f(\tau), \tau^{-1})$

$$_{i} \mid \{g_{d}^{r}\left(x_{d}^{i}\right)\}, \{\lambda_{r}\}, \tau \sim \mathcal{N}\left(f\left(\boldsymbol{x}_{i}\right), \tau^{-1}\right).$$

$$(9)$$

147 3.2 BKTF as a two-layer deep GP

Here we show the representation of BKTF as a two-layer deep GP. The **first** layer characterizes the generation of latent functions $\{g_d^r\}_{r=1}^R$ for dimension d. For the **second** layer, if we consider $\{g_1^r, \ldots, g_D^r\}_{r=1}^R$ as parameters and rewrite the functional decomposition in Eq. (7) as a linear function $f(\boldsymbol{x}; \{\lambda_r\}) = \sum_{r=1}^R \lambda_r \prod_{d=1}^D g_d^r(x_d)$ with $\lambda_r \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$, we can marginalize $\{\lambda_r\}$ and obtain a fully symmetric multilinear kernel/covariance function for any two data points $\boldsymbol{x} = (x_1, \ldots, x_D)^\top$ and $\boldsymbol{x}' = (x'_1, \ldots, x'_D)^\top$:

$$k\left(\boldsymbol{x}, \boldsymbol{x}'; \{g_1^r, \dots, g_D^r\}_{r=1}^R\right) = \sum_{r=1}^R \prod_{d=1}^D g_d^r\left(x_d\right) g_d^r\left(x_d'\right).$$
(10)

As can be seen, the second layer has a multilinear product kernel function parameterized by $\{g_1^r, \ldots, g_D^r\}_{r=1}^R$. There are some properties to highlight: (i) the kernel is **nonstationary** since the value of $g_d^r(\cdot)$ is location specific, and (ii) the kernel is **nonseparable** when R > 1. Therefore, this specification is very different from traditional GP surrogates, such as:

$$\begin{cases} \text{GP ARD:} & k\left(\boldsymbol{x}, \boldsymbol{x}'\right) = \prod_{d=1}^{D} k_d\left(x_d, x_d'\right), \\ & \text{kernel function is stationary and separable} \\ \text{Additive GP (2nd order):} & k\left(\boldsymbol{x}, \boldsymbol{x}'\right) = \sum_{d=1}^{D} k_d^1\left(x_d, x_d'\right) + \sum_{d=1}^{D-1} \sum_{e=d+1}^{D} k_d^2\left(x_d, x_d'\right) k_e^2\left(x_e, x_e'\right), \\ & \text{kerenl is stationary and nonseparable} \end{cases}$$

where all kernel functions are stationary with different hyperparameters (e.g., length scale and variance). Compared to the GP-based kernel specification, the multilinear kernel in Eq. (10) has a much larger set of hyperparameters and becomes more flexible and adaptive for the data. From a GP perspective, learning the hyperparameter in the kernel function in Eq. (10) will be computationally expensive; however, we can achieve efficient Bayesian inference of $\{\lambda_r, g_1^r, \ldots, g_D^r\}_{r=1}^R$ under a kernelized tensor factorization framework.

164 **3.3 Model inference**

Unlike GP, BKTF no longer enjoys an analytical posterior distribution. Based on the aforementioned prior and hyperprior settings, we adapt the MCMC updating procedure in Ref. [10, 11] to an efficient Gibbs sampling algorithm for model inference. This allows us to accommodate observations that are not located in the grid space $\prod_{d=1}^{D} S_d$. The detailed derivation of the sampling algorithm is given in Appendix A. In terms of computational cost, we note that updating g_d^r and kernel hyperparameters requires min { $\mathcal{O}(n^3), \mathcal{O}(|S_d|^3)$ } in time. Sparse GP (such as [14]) could be introduced when $n, |S_d|$ become large. See Appendix A, F for detailed discussion/comparison about computation complexity.

172 3.4 Prediction and AF computation

In each step of function evaluation, we run the MCMC sampling process K iterations for model inference, where the first K_0 samples are taken as burn-in and the last $K - K_0$ samples are used for posterior approximation. The predictive distribution for any entry f^* in the defined grid space conditioned on the observed dataset \mathcal{D}_n can be obtained by the Monte Carlo approximation $p(f^* | \mathcal{D}_n, \theta_0) \approx \frac{1}{K - K_0} \times \sum_{k=K_0+1}^{K} p\left(f^* | (\boldsymbol{g}_d^r)^{(k)}, \boldsymbol{\lambda}^{(k)}, \tau^{(k)}\right)$, where $\theta_0 = \{\mu_l, \tau_l, a_0, b_0\}$ is the set of all parameters used in hyperpriors. Although a direct analytical predictive distribution does not exist in BKTF, we can use MCMC samples to obtain the mean and variance of all points on the grid, thus offering an enumeration-based approach to define AF.

We define a Bayesian variant of the UCB as the AF by approximating the predictive mean and variance 181 (or uncertainty) in ordinary GP-based UCB with the values calculated from MCMC sampling. 182 For every MCMC sample after burn-in, i.e., $k > K_0$, we can estimate an output tensor $\tilde{\mathcal{F}}^{(\vec{k})}$ over the entire grid space using the latent factors $(\boldsymbol{g}_d^r)^{(k)}$ and the weight vector $\boldsymbol{\lambda}^{(k)}$: $\tilde{\mathcal{F}}^{(k)} = \sum_{r=1}^R \lambda_r^{(k)} (\boldsymbol{g}_1^r)^{(k)} \circ (\boldsymbol{g}_2^r)^{(k)} \circ \cdots \circ (\boldsymbol{g}_D^r)^{(k)}$. We can then compute the corresponding mean and 183 184 185 variance tensors of the $(K - K_0)$ samples $\{\tilde{\boldsymbol{\mathcal{F}}}^{(k)}\}_{k=K_{\mathbb{Q}}+1}^{K}$, and denote the two tensors by $\boldsymbol{\mathcal{U}}$ and $\boldsymbol{\mathcal{V}}$, respectively. The approximated predictive distribution at each point \boldsymbol{x} in the space becomes 186 187 $\tilde{f}(\boldsymbol{x}) \sim \mathcal{N}(u(\boldsymbol{x}), v(\boldsymbol{x}))$. Following the definition of UCB in Eq. (5), we define Bayesian UCB 188 (B-UCB) at location \boldsymbol{x} as $\alpha_{\text{B-UCB}} (\boldsymbol{x} \mid \mathcal{D}, \beta, \boldsymbol{g}_d^r, \boldsymbol{\lambda}) = u(\boldsymbol{x}) + \beta \sqrt{v(\boldsymbol{x})}$. The next search/query point can be determined via $\boldsymbol{x}_{\text{next}} = \arg \max_{\boldsymbol{x} \in \{\prod_{d=1}^{D} S_d - \mathcal{D}_{n-1}\}} \alpha_{\text{B-UCB}} (\boldsymbol{x})$. 189 190 We summarize the implementation procedure of BKTF for BO in Appendix B (see Algorithm 2). 191 192

Given the sequential nature of BO, when a new data point arrives at step n, we can start the MCMC with the last iteration of the Markov chains at step n - 1 to accelerate model convergence. The main computational and storage cost of BKTF is to interpolate and save the tensors $\tilde{\mathcal{F}} \in \mathbb{R}^{|S_1| \times \cdots \times |S_D|}$ over $(K - K_0)$ iterations for Bayesian AF estimation. This could be prohibitive when the MCMC sample size K or the dimensionality D becomes large. To avoid saving the tensors, in practice, we can simply use the maximum values of each entry over the $(K - K_0)$ iterations through iterative pairwise comparison. The number of samples after burn-in then implies the value of β in $\alpha_{\text{B-UCB}}$. We adopt this simple AF in our numerical experiments. A critical challenge in BKTF is that tensor size grows exponentially with the number of dimensions. To decrease the computational burden of enumeration-based AF, we also implement BKTF with random discretization—randomly selecting candidate points instead of reconstructing the whole space, denoted as **BKTFrandom**. BKTFrandom can be applied to functions with higher dimensions (e.g., D > 10). We discuss the comparison between BKTF and BKTFrandom in Experiments on test functions, see Section 5.1.

206 4 Related work

The key of BO is to effectively characterize the posterior distribution of the objective function 207 from a limited number of observations. The most relevant work to our study is the Bayesian 208 Kernelized Factorization (BKF) framework, which has been mainly used for modeling large-scale and 209 multidimensional spatiotemporal data with UQ. The key idea is to parameterize the multidimensional 210 stochastic processes using a factorization model, in which specific priors are used to encode spatial 211 and temporal dependencies. Signature examples of BKF include spatial dynamic factor model 212 (SDFM) [15], variational Gaussian process factor analysis (VGFA) [16], and Bayesian kernelized 213 matrix/tensor factorization (BKMF/BKTF) [10, 11, 17]. A common solution in these models is to 214 use GP prior to modeling the factor matrices, thus encoding spatial and temporal dependencies. In 215 addition, for categorical dimensions, BKTF uses an inverse-Wishart prior to modeling the covariance 216 matrix for the latent factors. A key difference among these methods is how inference is performed. 217 SDFM and BKMF/BKTF are fully Bayesian hierarchical models and they rely on MCMC for model 218 inference, where the factors can be updated via Gibbs sampling with conjugate priors; for learning 219 the posterior distributions of kernel hyperparameters, SDFM uses the Metropolis-Hastings sampling, 220 while BKMF/BKTF uses the more efficient slice sampling. On the other hand, VGFA uses variational 221 inference to learn factor matrices, while kernel hyperparameters are learned through maximum a 222 posteriori (MAP) estimation without UQ. Overall, BKTF has shown superior performance in modeling 223 multidimensional spatiotemporal processes with high-quality UQ for 2D/3D spaces [11, 17]. 224

The proposed BKTF surrogate models the objective function—as a single realization of a random 225 process—using low-rank tensor factorization with random basis functions. This basis function-226 based specification is closely related to multidimensional Karhunen-Loève (KL) expansion [18] for 227 stochastic (spatial, temporal, and spatiotemporal) processes. Empirical analysis of KL expansion is 228 also known as proper orthogonal decomposition (POD). With a known kernel/covariance function, 229 truncated KL expansion allows us to approximate the underlying random process using a set of 230 eigenvalues and eigenfunctions derived from the kernel function. Numerical KL expansion is often 231 referred to as the Garlekin method, and in practice, the basis functions are often chosen as prespecified 232 and deterministic functions [12, 19], such as Fourier basis, wavelet basis, orthogonal polynomials, 233 B-splines, empirical orthogonal functions, radial basis functions (RBF), and Wendland functions 234 (i.e., compactly supported RBF) (see, e.g., [20], [21], [22], [23]). However, the quality of UQ will be 235 undermined as the randomness is fully attributed to the coefficients $\{\lambda_r\}$; in addition, these methods 236 also require a large number of basis functions (or a large rank) to fit complex stochastic processes. 237 238 Different from methods with fixed/known basis functions, BKTF uses a Bayesian hierarchical modeling framework to better capture the randomness and uncertainty in the data, in which GP priors 239 are used to model the latent factors (i.e., basis functions are also random processes) on different 240 dimensions, and hyperpriors are introduced on the kernel hyperparameters. Therefore, BKTF becomes 241 a fully Bayesian version of multidimensional KL expansion for stochastic processes with unknown 242 covariance from partially observed data, however, without imposing any orthogonal constraint on 243 the basis functions. Following the analysis in section 3.2, BKTF is also a special case of a two-layer 244 deep Gaussian process [24, 7], where the first layer produces latent factors for each dimension, and 245 the second layer has a multilinear kernel parameterized by all latent factors. 246

247 **5 Experiments**

248 5.1 Optimization for benchmark test functions

We test the proposed BKTF model for BO on seven benchmark functions that are used for global optimization problems [25], with dimension *D* ranging from 2 to 10. All selected standard functions are multimodal; detailed descriptions are given in Appendix D. In fact, we can visually see that the standard Damavandi/Schaffer/Griewank functions in Figure 7 (see Appendix D) indeed have a



Figure 2: Optimization on benchmark test functions, where medians with 25% and 75% quartiles of 10 runs are compared.

low-rank structure. For each function, we assume the initial dataset \mathcal{D}_0 contains $n_0 = D$ observed data pairs, and we set the total number of query points to N = 80 for 4D Griewank and 6D Hartmann function, N = 200 for 10D Griewank, and N = 50 for others. We rescale the input search range to [0, 1] for all dimensions and normalize the output data using z-score normalization.

Model configuration When applying BKTF to continuous test functions, we introduce m_d interpo-257 lation points for the dth dimension of the input space. The values of m_d used for each benchmark 258 function are predefined and given in Table 1 (see Appendix D). Setting the resolution grid will require 259 certain prior knowledge (e.g., smoothness of the function); and it also depends on the available 260 computational resources and the number of entries in the tensor, which grows exponentially with m_d . 261 In practice, we find that setting $m_d = 10 \sim 100$ is sufficient for most problems. We set the CP rank 262 R = 2, and for each BO function evaluation run 400 MCMC iterations for model inference where 263 the first 200 iterations are taken as burn-in. We use Matérn 3/2 kernel as the covariance function for 264 all the test functions. 265

Note that for the 10D Griewank function, the grid-based models do not work, and only models built in continuous space and BKTFrandom can be performed. For BKTFrandom, in each evaluation we randomly select 20k points in the search space as candidates and choose the one with the best AF as the next evaluation location.

Baselines We compare BKTF with the following BO methods that use GP as the surrogate model: (1) GP α_{EI} and (2) GPgrid α_{EI} : GP as the surrogate model and EI as the AF, in continuous space $\prod_{d=1}^{D} \mathcal{X}_d$ and Cartesian grid space $\prod_{d=1}^{D} S_d$, respectively; (3) GP α_{UCB} and (4) GPgrid α_{UCB} : GP as the surrogate model with UCB as the AF where $\beta = 2$, in $\prod_{d=1}^{D} \mathcal{X}_d$ and $\prod_{d=1}^{D} S_d$, respectively; (5) additive GP: the sum of two 1st-order additive kernels per dimension as the surrogate with EI as the AF, in continuous space. This baseline has the same number of latent functions as BKTF (R = 2) but in a sum-based manner; (6) deepGP: a two-layer fully-Bayesian deep GP with EI as the AF, implemented with the *deepgp* package¹.

Similar as the setting for BKTF, we use Matérn 3/2 kernel in all GP models. Given the computational cost, we only compare deepGP on 2D functions [9]. For AF optimization in GP α_{EI} and GP α_{UCB} , we first use the DIRECT algorithm [26] and then apply the Nelder-Mead algorithm [27] to further search if there exist better solutions. We also compare with SAASBO (Sparse Axis-Aligned Subspace) [28] with Hamiltonian Monte Carlo sampling, implemented using BoTorch [29], on the 6D Hartmann and 10D Griewank functions.

Results To compare the optimization performance of different models on the benchmark functions, we define the absolute error between the global optimum f^* and the current estimated global optimum \hat{f}^* , i.e., $\left|f^* - \hat{f}^*\right|$, as the *regret*, and examine how *regret* varies with the number of function

¹https://CRAN.R-project.org/package=deepgp

evaluations. We run the optimization 10 times for every test function with a different set of initial 287 observations. The results are summarized in Figure 2. We see that for the 2D functions Branin 288 and Schaffer, BKTF clearly finds the global optima much faster than GP surrogate-based baselines. 289 For Damavandi function, where the global minimum $(f(\mathbf{x}^*) = 0)$ is located in a small sharp area 290 while the local optimum (f(x) = 2) is located at a large smooth area (see Figure 7 in Appendix D), 291 GP-based models are trapped around the local optima in most cases (i.e., regret = 2) and cannot 292 293 jump out. In contrast, BKTF explores the global characteristics of the objective function over the entire search space and reaches the global optimum within 10 iterations of function evaluations. 294 For higher dimensional Griewank and Hartmann functions, BKTF successfully arrives at the global 295 optima under the given observation budgets, while GP-based comparison methods are prone to be 296 stuck around local optima. We compare the *regret* at the last iteration in Table 2 (Appendix E.2), 297 along with the average cost of evaluations. The enumeration-based GP surrogates, i.e., GPgrid $\alpha_{\rm EI}$ 298 and GPgrid α_{UCB} , perform a little better than direct GP-based search, i.e., GP α_{EI} and GP α_{UCB} on 299 Damavandi function, but worse on others. This means that the discretization, to some extent, offers 300 possibilities for searching all the alternative points in the space, since in each function evaluation, 301 every sample in the space is equally compared solely based on the predictive distribution. Additive 302 GP is comparable with R = 2 BKTF; while the results demonstrate that BKTF can be much more 303 flexible than additive GP. As for BKTFrandom, we see that it can alleviate the curse of dimensionality 304 and be applied for higher-dimensional problems that may not be performed with a grid but at the 305 cost of more evaluation budgets, particularly it costs more iterations for lower-dimensional functions 306 compared with BKTF. 307

Overall, BKTF reaches the global optimum for every test function and shows superior performance for complex objective functions with a faster convergence rate. To intuitively compare the overall performance of different models across multiple experiments/functions, we further estimate performance profiles (PPs) [30] (see Appendix E.1), and compute the area under the curve (AUC) for quantitative analysis (see bottom right of Figure 2 and Table 2 in Appendix E.2). Our results show that BKTF obtains the best performance across all functions.

Interpretable latent space. The sampled latent functions are interpretable and imply the under-314 lying correlations of the objective function. We illustrate the learned periodic (global) patterns in 315 Appendix E.3. Effects of hyperpriors. Since we build a fully Bayesian model, the hyperparameters 316 of the covariance functions can be updated automatically from the data likelihood and hyperprior. 317 Note that in optimization scenarios where the observations are scarce, the prediction performance of 318 BKTF highly depends on the hyperprior settings, i.e., $\theta_0 = \{\mu_l, \tau_l, a_0, b_0\}$. We discuss the effects of 319 hyperpriors in Appendix E.4. Effects of rank. The only predefined model parameter is the model 320 rank, all other model parameters and hyperparameters are sampled with MCMC. We discuss the 321 effects of rank on the 2D nonstationary nonseparable function defined in Introduction (see Figure 1) 322 323 in Appendix E.5. We see that BKTF is robust to rank specification and can find the global solution efficiently with rank set as 2, 4, 6 and 8. 324

325 5.2 Hyperparameter tuning for machine learning

In this section, we evaluate the performance of BKTF for automatic machine learning (ML) tasks. 326 We compare different models to optimize the hyperparameters of two ML algorithms-random forest 327 (RF) and neural network (NN)—on classification for the MNIST database of handwritten digits² 328 and regression for the Boston housing dataset³. The tuning tasks involve both integer-valued and 329 categorical parameters, and the details are given in the Appendix G. We treat those integer-valued 330 dimensions as continuous and use a Matérn 3/2 kernel for the basis functions. Given the size of the 331 hyperparameter space, we perform BKTFrandom for classification and BKTF for regression. We 332 assume that the number of initial observations \mathcal{D}_0 equals the number of tuning hyperparameters. 333 The total budget N is 20 for the classification task and 50 for the regression. We implement RF 334 algorithms using the scikit-learn package and construct NN models by Keras with 2 hidden layers. 335 All other model hyperparameters are set as the default values. 336

Model configuration We treat all the integer hyperparameters as samples from a continuous space and then generate the corresponding Cartesian product space $\prod_{d=1}^{D} S_d$. One can interpret

²http://yann.lecun.com/exdb/mnist/

³https://www.cs.toronto.edu/~delve/data/boston/bostonDetail.html



Figure 3: Comparison of hyperparameter tuning for ML tasks: (a) classification; (b) regression.

the candidate values for each hyperparameter as the interpolation points in the corresponding input dimension. The size of the spanned space $\prod S_d$ is $91 \times 46 \times 64 \times 10 \times 11 \times 2$ and $91 \times 46 \times 13 \times 10$ for RF classifier and RF regressor, respectively; $91 \times 49 \times 31 \times 18 \times 3 \times 2$ and $91 \times 49 \times 31$ for NN classifier and NN regressor respectively. We set the tensor rank R = 4 for BKTF, set K = 400

and $K_0 = 200$ for MCMC inference, and use the Matérn 3/2 kernel for capturing correlations.

Baselines In addition to GP surrogate-based GP α_{EI} and GP α_{UCB} , we also compare with a baseline method: random search (RS), and two non-GP models: particle swarm optimization (PSO) [31] and Tree-structured Parzen Estimator (BO-TPE) [32], which are common methods for hyperparameter tuning. We exclude grid-based GP models as sampling the entire grid becomes infeasible.

Results We compare the accuracy for MNIST classification and MSE (mean squared error) for 348 Boston housing regression both in terms of the number of function evaluations and still run the 349 optimization processes ten times with different initial datasets \mathcal{D}_0 . The results obtained by different 350 351 BO models are given in Figure 3, and the final classification accuracy and regression MSE are compared in Table 5 (see Appendix H). For BKTF, we see from Figure 3 that the width between the 352 two quartiles of accuracy and error decreases as more iterations are evaluated, and the median curves 353 present better convergence rates compared to the baselines. Table 5 also shows that the proposed 354 BKTF surrogate achieves the best final mean accuracy and regression error with small standard 355 deviations. The results above demonstrate the advantage of BKTF as a surrogate. 356

357 6 Conclusion

This paper proposes to use Bayesian Kernelized Tensor Factorization (BKTF) as a new surrogate 358 model for Bayesian optimization with mixed variables (both discrete/categorical and continuous) 359 when the dimensionality is relatively small (e.g., say D < 10). Compared with traditional GP surro-360 gates, the BKTF surrogate is more flexible and adaptive to data thanks to the Bayesian hierarchical 361 specification, which provides high-quality UQ for BO tasks. The tensor factorization model behind 362 BKTF offers an effective solution to capture global/long-range correlations and cross-dimension cor-363 relations. The inference of BKTF is achieved through MCMC, which provides a natural solution for 364 365 acquisition. Experiments on both test function optimization and ML hyperparameter tuning confirm 366 the superiority of BKTF as a surrogate for BO. Moreover, the tensor factorization framework makes 367 it straightforward to adapt BKTF to handle multivariate and functional output (see e.g., [33, 34]), by directly treating the output coordinates as part of the input. A limitation of BKTF is that we restrict 368 BO to a grid search space in order to leverage tensor factorization; however, we believe that designing 369 a compatible grid space based on prior knowledge is not a challenging task. 370

There are several directions to be explored to make BKTF more scalable. Scalable GP solutions, such 371 372 as sparse GP [14] and Gaussian Markov Random Field (GMRF) [35], can be introduced to reduce the inference cost when $|S_d|$ becomes large. The current MCMC-based acquisition method requires 373 explicit reconstruction of the whole tensor, which quickly becomes infeasible when D becomes large 374 (e.g., D > 20). A natural question is whether it is possible to achieve efficient acquisition directly 375 using the basis functions and the corresponding weights without explicitly constructing the tensors. 376 This problem corresponds to finding/locating the maximum entry of a tensor given its low-rank 377 decomposition (see e.g., [36]). 378

This work aims to advance the field of probabilistic machine learning, particularly Bayesian optimization. Regardless of the model limitations, it has the potential of misuse for ML algorithms.

381 References

- [1] Roman Garnett. *Bayesian Optimization*. Cambridge University Press, 2023.
- [2] Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P Adams, and Nando De Freitas. Taking
 the human out of the loop: A review of Bayesian optimization. *Proceedings of the IEEE*,
 104(1):148–175, 2015.
- [3] Robert B Gramacy. Surrogates: Gaussian Process Modeling, Design, and Optimization for the
 Applied Sciences. Chapman and Hall/CRC, 2020.
- [4] Christopher KI Williams and Carl Edward Rasmussen. *Gaussian Processes for Machine Learning*. MIT Press, Cambridge, MA, 2006.
- [5] David K Duvenaud, Hannes Nickisch, and Carl Rasmussen. Additive Gaussian processes.
 Advances in neural information processing systems, 24, 2011.
- [6] Mickael Binois and Nathan Wycoff. A survey on high-dimensional Gaussian process modeling
 with application to Bayesian optimization. *ACM Transactions on Evolutionary Learning and Optimization*, 2(2):1–26, 2022.
- [7] Andreas Damianou and Neil D Lawrence. Deep Gaussian processes. In *International Conference* on Artificial Intelligence and Statistics, pages 207–215, 2013.
- [8] Ali Hebbal, Loïc Brevault, Mathieu Balesdent, El-Ghazali Talbi, and Nouredine Melab.
 Bayesian optimization using deep gaussian processes with applications to aerospace system
 design. *Optimization and Engineering*, 22:321–361, 2021.
- [9] Annie Sauer, Robert B Gramacy, and David Higdon. Active learning for deep gaussian process
 surrogates. *Technometrics*, 65(1):4–18, 2023.
- [10] Mengying Lei, Aurelie Labbe, Yuankai Wu, and Lijun Sun. Bayesian kernelized matrix
 factorization for spatiotemporal traffic data imputation and kriging. *IEEE Transactions on Intelligent Transportation Systems*, 23(10):18962–18974, 2022.
- [11] Mengying Lei, Aurelie Labbe, and Lijun Sun. Bayesian complementary kernelized learning for
 multidimensional spatiotemporal data. *arXiv preprint arXiv:2208.09978*, 2022.
- [12] Noel Cressie, Matthew Sainsbury-Dale, and Andrew Zammit-Mangion. Basis-function models
 in spatial statistics. *Annual Review of Statistics and Its Application*, 9:373–400, 2022.
- [13] Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3(Nov):397–422, 2002.
- [14] Joaquin Quinonero-Candela and Carl Edward Rasmussen. A unifying view of sparse approximate Gaussian process regression. *The Journal of Machine Learning Research*, 6:1939–1959, 2005.
- [15] Hedibert Freitas Lopes, Esther Salazar, and Dani Gamerman. Spatial dynamic factor analysis.
 Bayesian Analysis, 3(4):759–792, 2008.
- [16] Jaakko Luttinen and Alexander Ilin. Variational Gaussian-process factor analysis for modeling
 spatio-temporal data. Advances in Neural Information Processing Systems, 22:1177–1185,
 2009.
- [17] Mengying Lei, Aurelie Labbe, and Lijun Sun. Scalable spatiotemporally varying coefficient
 modeling with bayesian kernelized tensor regression. *arXiv preprint arXiv:2109.00046*, 2021.
- [18] Limin Wang. *Karhunen-Loeve expansions and their applications*. London School of Economics
 and Political Science (United Kingdom), 2008.
- [19] Holger Wendland. Scattered Data Approximation, volume 17. Cambridge university press,
 2004.

- Rommel G Regis and Christine A Shoemaker. A stochastic radial basis function method for the
 global optimization of expensive functions. *INFORMS Journal on Computing*, 19(4):497–509,
 2007.
- 428 [21] Gregory Beylkin, Jochen Garcke, and Martin J Mohlenkamp. Multivariate regression and
 429 machine learning with sums of separable functions. *SIAM Journal on Scientific Computing*,
 430 31(3):1840–1857, 2009.
- [22] Christopher K Wikle and Noel Cressie. A dimension-reduced approach to space-time kalman
 filtering. *Biometrika*, 86(4):815–829, 1999.
- [23] Mathilde Chevreuil, Régis Lebrun, Anthony Nouy, and Prashant Rai. A least-squares method
 for sparse low rank approximation of multivariate functions. *SIAM/ASA Journal on Uncertainty Quantification*, 3(1):897–921, 2015.
- [24] Alexandra M Schmidt and Anthony O'Hagan. Bayesian inference for non-stationary spatial
 covariance structure via spatial deformations. *Journal of the Royal Statistical Society: Series B* (*Statistical Methodology*), 65(3):743–758, 2003.
- Momin Jamil and Xin-She Yang. A literature survey of benchmark functions for global
 optimization problems. *arXiv preprint arXiv:1308.4008*, 2013.
- [26] Donald R Jones, Cary D Perttunen, and Bruce E Stuckman. Lipschitzian optimization without
 the lipschitz constant. *Journal of optimization Theory and Applications*, 79(1):157–181, 1993.
- [27] John A Nelder and Roger Mead. A simplex method for function minimization. *The computer journal*, 7(4):308–313, 1965.
- [28] David Eriksson and Martin Jankowiak. High-dimensional bayesian optimization with sparse
 axis-aligned subspaces. In *Uncertainty in Artificial Intelligence*, pages 493–503. PMLR, 2021.
- [29] Maximilian Balandat, Brian Karrer, Daniel Jiang, Samuel Daulton, Ben Letham, Andrew G Wilson, and Eytan Bakshy. Botorch: A framework for efficient monte-carlo bayesian optimization.
 Advances in neural information processing systems, 33:21524–21538, 2020.
- [30] Elizabeth D Dolan and Jorge J Moré. Benchmarking optimization software with performance
 profiles. *Mathematical programming*, 91(2):201–213, 2002.
- [31] Jun Tang, Gang Liu, and Qingtao Pan. A review on representative swarm intelligence algorithms
 for solving optimization problems: Applications and trends. *IEEE/CAA Journal of Automatica Sinica*, 8(10):1627–1643, 2021.
- [32] James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper parameter optimization. *Advances in neural information processing systems*, 24, 2011.
- [33] Dave Higdon, James Gattiker, Brian Williams, and Maria Rightley. Computer model calibration
 using high-dimensional output. *Journal of the American Statistical Association*, 103(482):570–
 583, 2008.
- [34] Shandian Zhe, Wei Xing, and Robert M Kirby. Scalable high-order gaussian process regression.
 In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2611–2620.
 PMLR, 2019.
- [35] Havard Rue and Leonhard Held. *Gaussian Markov random fields: theory and applications*.
 Chapman and Hall/CRC, 2005.
- [36] Anastasiia Batsheva, Andrei Chertkov, Gleb Ryzhakov, and Ivan Oseledets. Protes: probabilistic
 optimization with tensor sampling. *Advances in Neural Information Processing Systems*, 36:808–
 823, 2023.

468 Appendix

469 Contents (Appendix)

470	A	Mod	el inference	12
471		A.1	Sampling latent functions	12
472		A.2	Sampling kernel hyperparameters	13
473		A.3	Sampling Λ_d for latent functions on categorical inputs $\ldots \ldots \ldots \ldots \ldots$	13
474		A.4	Sampling weight vector	14
475		A.5	Sampling model noise precision	14
476	B	Algo	rithm of BKTF for BO	15
477	С	Opti	mization for the 2D nonstationary and nonseparable function	15
478		C.1	Data generation	15
479		C.2	Experimental setting	15
480		C.3	Results	15
481	D	Bene	chmark test functions	17
482	Е	Sup	plementary results on benchmark test functions	19
483		E.1	Performance profiles	19
484		E.2	Quantitative comparison	20
485		E.3	Interpretation of results	21
486		E.4	Effects of hyperpriors	22
487		E.5	Effects of rank	23
488	F	Com	parison of computational complexity	24
488 489	F G	Com Hyp	parison of computational complexity erparameter tuning for machine learning	24 25

491 A Model inference

Assume an observation dataset $\mathcal{D}_n = \{x_i, y_i\}_{i=1}^n$, we exploit an efficient element-wise Gibbs sampling algorithm for model inference.

494 A.1 Sampling latent functions

Given the Gaussian prior and Gaussian likelihood of the latent factors \boldsymbol{g}_{d}^{r} , their posterior distributions are still Gaussian. Let $y_{r}^{i} = y_{i} - \sum_{\substack{h=1 \ h\neq r}}^{R} \lambda_{h} \prod_{d=1}^{D} g_{d}^{h} \left(x_{d}^{i}\right)$ and $\boldsymbol{y}_{r} = [y_{r}^{1}, \dots, y_{r}^{n}]^{\top} \in \mathbb{R}^{n}$ for $r = 1, \dots, R$. Every \boldsymbol{y}_{r} generates a D-dimensional tensor $\boldsymbol{\mathcal{Y}}_{r} \in \mathbb{R}^{|S_{1}| \times \dots \times |S_{D}|}$ in the Cartesian product space $\prod_{d=1}^{D} S_{d}$. We define a binary tensor $\boldsymbol{\mathcal{O}}$ with the same size of $\boldsymbol{\mathcal{Y}}$ indicating the locations of the observation data, where $o(\boldsymbol{x}_{i}) = 1$ for $i \in [1, n]$, and other values are zero. The posterior of 500 \boldsymbol{g}_d^r is given by

$$p\left(\boldsymbol{g}_{d}^{r}\mid-\right) = \mathcal{N}\left(\boldsymbol{g}_{d}^{r}\mid\left[\boldsymbol{\mu}_{d}^{r}\right]^{*},\left(\left[\boldsymbol{\Lambda}_{d}^{r}\right]^{*}\right)^{-1}\right),\tag{11}$$

501 where

$$[\boldsymbol{\mu}_{d}^{r}]^{*} = \left([\boldsymbol{\Lambda}_{d}^{r}]^{*} \right)^{-1} \underbrace{\left(\tau \left(\boldsymbol{Y}_{r(d)} \circledast \boldsymbol{O}_{(d)} \right) \left(\lambda_{r} \bigotimes_{\substack{h=D\\h \neq d}}^{1} \boldsymbol{g}_{h}^{r} \right) \right)}_{\boldsymbol{a} \in \mathbb{R}^{|S_{d}|}}, \tag{12}$$

502

$$\left[\mathbf{\Lambda}_{d}^{r}\right]^{*} = \tau \operatorname{diag}\left(\mathbf{O}_{(d)}\left(\lambda_{r}\bigotimes_{\substack{h=D\\h\neq d}}^{1} \mathbf{g}_{h}^{r}\right)^{2}\right) + \left(\mathbf{K}_{d}^{r}\right)^{-1}.$$

$$(13)$$

$$\mathbf{b} \in \mathbb{R}^{\lceil S_{d} \rceil}$$

⁵⁰³ $Y_{r(d)}$ and $O_{(d)}$ are mode-*d* unfoldings of \mathcal{Y}_r and \mathcal{O} , respectively, with the size of $|S_d| \times$ ⁵⁰⁴ $(\prod_{h=1}^{D} |S_h|)$. Note that the vector term *a* in $[\mu_d^r]^*$ and *b* in $[\Lambda_d^r]^*$, which are only relevant to ⁵⁰⁵ the *n* observations and corresponding function values, can be computed element-wise instead of ⁵⁰⁶ using matrix multiplication and Kronecker product. The point-wise computation can dramatically ⁵⁰⁷ reduce the computational cost, especially for a relatively large *D*, since in such case the number of ⁵⁰⁸ observations in BO can be much smaller compared to the number of samples in the entire grid space, ⁵⁰⁹ i.e., $n \ll \prod_{d=1}^{D} |S_d|$.

510 A.2 Sampling kernel hyperparameters

We update the kernel lengthscale hyperparameters l_d^r from their marginal posteriors by integrating out the latent factors. Let $[\boldsymbol{y}_r]_d = \boldsymbol{O}_d \operatorname{vec} (\boldsymbol{Y}_{r(d)}) \in \mathbb{R}^n$, where $\boldsymbol{O}_d \in \mathbb{R}^{n \times (\prod_{d=1}^{D} |S_d|)}$ is a binary matrix obtained by removing the rows corresponding to zeros in $\operatorname{vec} (\boldsymbol{O}_{(d)})$ from $\boldsymbol{I}_{\prod_{d=1}^{D} |S_d|}$. When sampling the posteriors for kernel hyperparameters under a given d and r, their marginal likelihoods only relate to $[\boldsymbol{y}_r]_d$. The log marginal likelihood of l_d^r , for example, is:

$$\log p\left([\boldsymbol{y}_{r}]_{d} \mid l_{d}^{r}\right) \propto -\frac{1}{2}\left([\boldsymbol{y}_{r}]_{d}\right)^{\top} \boldsymbol{\Sigma}_{[\boldsymbol{y}_{r}]_{d}|l_{d}^{r}}^{-1} [\boldsymbol{y}_{r}]_{d} - \frac{1}{2}\log\det\left(\boldsymbol{\Sigma}_{[\boldsymbol{y}_{r}]_{d}|l_{d}^{r}}\right)$$

$$\propto \frac{\tau^{2}}{2} \underbrace{\left([\boldsymbol{y}_{r}]_{d}\right)^{\top} \boldsymbol{H}\left([\boldsymbol{\Lambda}_{d}^{r}]^{*}\right)^{-1} \boldsymbol{H}^{\top} [\boldsymbol{y}_{r}]_{d}}_{c} - \frac{1}{2}\log\det\left([\boldsymbol{\Lambda}_{d}^{r}]^{*}\right) - \frac{1}{2}\log\det\left(\boldsymbol{K}_{d}^{r}\right),$$
(14)

516 where
$$\boldsymbol{H} = \boldsymbol{O}_d \left(\lambda_r \bigotimes_{\substack{h=D \ h \neq d}}^1 \boldsymbol{g}_h^r \otimes \boldsymbol{I}_{(|S_d|)} \right) \in \mathbb{R}^{n \times |S_d|}$$
 and $\boldsymbol{\Sigma}_{[\boldsymbol{y}_r]_d \mid l_d^r} = \boldsymbol{H} \boldsymbol{K}_d^r \boldsymbol{H}^\top + \tau^{-1} \boldsymbol{I}_n$. The

term c in Eq. (14) is a scalar that can be computed with $\boldsymbol{u}^{\top}\boldsymbol{u}$, where $\boldsymbol{u} = (\boldsymbol{L}_d^r)^{-1}\boldsymbol{a}$ is a vector of length $|S_d|$; $\boldsymbol{L}_d^r = \text{chol}([\boldsymbol{\Lambda}_d^r]^*)$ is the Cholesky factor matrix of $[\boldsymbol{\Lambda}_d^r]^*$. This means that the complicated term c can also be calculated element-wise, and it leads to a fast learning process. With the marginal likelihood and predefined log-normal hyperpriors, we can get the marginal posteriors of the kernel hyperparameters straightforwardly; and we update them by using the slice sampling algorithm presented in [17].

523 A.3 Sampling Λ_d for latent functions on categorical inputs

For latent factors in dimensions with categorical variables/inputs, we sample the precision hyperparameter Λ_d in its prior distribution from a Wishart distribution:

$$\boldsymbol{\Lambda}_{d} \mid -\sim \mathcal{W}\left(\left(\boldsymbol{G}_{d}\boldsymbol{G}_{d}^{\top} + \boldsymbol{I}_{\mid S_{d}\mid}\right)^{-1}, \left|S_{d}\right| + R\right).$$
(15)

526 A.4 Sampling weight vector

527 Every observed data point has the following distribution:

$$y_i \sim \mathcal{N}\left(\sum_{r=1}^R \lambda_r \prod_{d=1}^D g_d^r\left(x_d^i\right) = \left[\prod_{d=1}^D g_d^1\left(x_d^i\right), \dots, \prod_{d=1}^D g_d^R\left(x_d^i\right)\right] \boldsymbol{\lambda}, \tau^{-1}\right), \ i = 1, \dots, n.$$
(16)

Let $\boldsymbol{g}(\boldsymbol{x}_i) = \left(\prod_{d=1}^{D} g_d^1\left(x_d^i\right), \dots, \prod_{d=1}^{D} g_d^R\left(x_d^i\right)\right)^\top \in \mathbb{R}^R$ and $\boldsymbol{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$; we then have $\boldsymbol{y} \sim \mathcal{N}\left(\tilde{\boldsymbol{G}}^\top \boldsymbol{\lambda}, \tau^{-1} \boldsymbol{I}_n\right)$, where $\tilde{\boldsymbol{G}} = [\boldsymbol{g}(\boldsymbol{x}_1), \dots, \boldsymbol{g}(\boldsymbol{x}_n)] \in \mathbb{R}^{R \times n}$. The posterior of $\boldsymbol{\lambda}$ follows a Gaussian distribution $\boldsymbol{y}(\boldsymbol{\lambda} \mid -) \sim \mathcal{N}\left(\boldsymbol{\mu}_{\boldsymbol{\lambda}}^*, (\boldsymbol{\Lambda}_{\boldsymbol{\lambda}}^*)^{-1}\right)$. (17)

$$p(\boldsymbol{\lambda} \mid -) \sim \mathcal{N}\left(\boldsymbol{\mu}_{\boldsymbol{\lambda}}^{*}, (\boldsymbol{\Lambda}_{\boldsymbol{\lambda}}^{*})^{-1}\right),$$
 (17)

531 where

532

$$\boldsymbol{\mu}_{\lambda}^{*} = \tau \left(\boldsymbol{\Lambda}_{\lambda}^{*} \right)^{-1} \tilde{\boldsymbol{G}} \boldsymbol{y}, \tag{18}$$

$$\Lambda_{\lambda}^{*} = \tau \tilde{G} \tilde{G}^{\top} + I_{R}.$$
⁽¹⁹⁾

533 A.5 Sampling model noise precision

534 For precision τ , we have a Gamma posterior

$$p(\tau \mid -) = \operatorname{Gamma}(\tau \mid a^*, b^*), \qquad (20)$$

535 where

536

$$a^* = a_0 + \frac{1}{2}n,\tag{21}$$

$$b^* = b_0 + \frac{1}{2} \sum_{i=1}^n \left(y_i - \sum_{r=1}^R \lambda_r \prod_{d=1}^D g_d^r \left(x_d^i \right) \right)^2.$$
(22)

537 **B** Algorithm of BKTF for BO

Algorithm 2: BKTF for BO

```
Input: Initial dataset \mathcal{D}_0.
           for n = 1 to N do
              for k = 1 to K do
                  for r = 1 to R do
                      for d = 1 to D do
                          Draw kernel lengthscale hyperparameter l_d^r or precision hyperparameter \Lambda_d;
                          Draw latent factors (\boldsymbol{q}_{d}^{r})^{(k)};
                      end for
                  end for
                  Draw model noise precision \tau^{(k)};
538
                  Draw weight vector \lambda^{(k)};
                  if k > K_0 then
                      Compute and collect 	ilde{oldsymbol{\mathcal{F}}}^{(k)}.
                  end if
              end for
              Compute mean \mathcal{U} and variance \mathcal{V} of \{\tilde{\mathcal{F}}^{(k)}\};
               Compute \alpha_{\text{B-UCB}}(\boldsymbol{x} \mid \mathcal{D}_{n-1}, \beta) as a tensor;
              Find next \boldsymbol{x}_n = \arg \max_{\boldsymbol{x}} \alpha_{\text{B-UCB}} (\boldsymbol{x} \mid \mathcal{D}_{n-1}, \beta);
              Augment the data \mathcal{D}_n = \overline{\mathcal{D}}_{n-1} \cup \{x_n, y_n\}.
           end for
```

⁵³⁹ C Optimization for the 2D nonstationary and nonseparable function

540 C.1 Data generation

The function in Figure 1(a) of the paper (see Section 1 Introduction) is a modification of the case study used in [11]. It is a 40 × 40 2D process generated in a $[1,2] \times [-1,0]$ square, with

$$\mathbf{Y}(x_1, x_2) = \left(\cos\left(4\left[f_1(x_1) + f_2(x_2)\right]\right) + \sin\left(4\left[f_1(x_2) - f_2(x_1)\right]\right) - 1\right) \\ \times \exp\left(-(x_1 - 0.5)^2 + \frac{(x_2 - 1)^2}{5}\right), \quad (23)$$

where $f(x_1) = x_1 (\sin 2x_1 + 2)$, $f(x_2) = 0.2x_2 \sqrt{99(x_2 + 1) + 4}$, $x_1 \in [1, 2]$, $x_2 \in [-1, 0]$. This is a nonstationary, nonseparable, and multimodel function, with the global maximum $f(x^*) = 0.6028$ at $x^* = (1.75, -0.55)$.

546 C.2 Experimental setting

We randomly select $n_0 = 30$ data points as the initial data and run 50 iterations (i.e., budget) of evaluation for optimization. To compare different surrogates, we run the optimization for 20 replications with different initial datasets. For the proposed BKTF surrogate, we place Matérn 3/2 kernel functions on the latent factors, set the rank R = 4, and run 1000 MCMC samples for model inference where the first 600 samples are burn-in. For comparison of BO methods, we consider typical GP surrogate with both EI and UCB ($\beta = 2$) as the AF, denoted by GP α_{EI} and GP α_{UCB} respectively, and use the same Matérn 3/2 kernel for GP surrogate.

554 C.3 Results

Figure 1(b) in Section 1 Introduction of the paper shows the medians along with the 25% and 75% quantiles of the optimization results from 20 runs. We see that GP α_{EI} and GP α_{UCB} cannot find the global optimum in most cases, and they easily get stuck in the lower left flat area which contains easily find local optima. Figure 1(c) illustrates the estimation surface of the function and the estimated AF surface from one run. It is clear that BKTF can capture global correlations with limited data. The



Figure 4: Optimization on the 2D function (Eq. 23): Posterior distributions of model hyperparameters learned by BKTF.



Figure 5: Optimization on the 2D function (Eq. 23): Mean of latent factors sampled by BKTF.

search points contain areas of almost every peak in the true function, and the peaks of its AF surface also reflect the peak area in f.

⁵⁶² Figure 4 shows the approximated posterior distributions of model parameters in BKTF. Figure 5 and

Figure 6 illustrate the posterior mean and the last 20 MCMC samples of the latent factors, respectively.

564 Samples of the latent factors in panels (a) and (b) of Figure 6 explain the uncertainty learned by

565 BKTF for this optimization problem.



Figure 6: Optimization on the 2D function (Eq. 23): The last 20 MCMC samples of the latent factors in the two dimensions learned by BKTF.

566 **D** Benchmark test functions

We summarize the characteristics of the benchmark functions tested in Table 1. Figure 7 shows the functions with 2-dimensional inputs together with the 2D Griewank function. The functional expressions and more detailed features of these test functions are provided in the following.

Function	$\mid D$	Search space	m_d	Characteristics
Branin Damavandi Schaffer	$\begin{vmatrix} 2\\2\\2 \end{vmatrix}$	$\begin{array}{c} [-5,10]\times [0,15] \\ [0,14]^2 \\ [-10,10]^2 \end{array}$	14 71 11	3 global minima, flat multimodal, global minimum located in small area multimodal, global optimum located close to local minima
Griewank	3 4 10	$\begin{array}{c} [-10,10]^3 \\ [-10,10]^4 \\ [-10,10]^{10} \end{array}$	11 11 -	multimodal, many widespread and regularly distributed local optima
Hartmann	6	$[0,1]^6$	12	multimodal, multi-input

Table 1: Summary of the studied benchmark functions.



Figure 7: Examples of the tested benchmark functions.

(1) Branin function (D = 2)

$$f(x_1, x_2) = \left(x_2 - \frac{5.1}{4\pi}x_1^2 + \frac{5}{\pi}x_1 - 6\right)^2 + 10\left(1 - \frac{1}{8\pi}\right)\cos(x_1) + 10,$$
 (24)

where $x_1 \in [-5, 10]$ and $x_2 \in [0, 15]$. It is a smooth but multimodal function with global minima $f(\boldsymbol{x}^*) = 0.3978873$ at three input points $\boldsymbol{x}^* = (-\pi, 12.275), (\pi, 2.275), (3\pi, 2.425).$

(2) Damavandi function (D = 2)

$$f(x_1, x_2) = \left[1 - \left|\frac{\sin[\pi(x_1 - 2)]\sin[\pi(x_2 - 2)]}{\pi^2(x_1 - 2)(x_2 - 2)}\right|^5\right] \left[2 + (x_1 - 7)^2 + 2(x_2 - 7)^2\right], \quad (25)$$

where $x_d \in [0, 14]$. This is a multimodal function with the global minimum $f(x^*) = 0$ at $x^* = (2, 2)$.

(3) Schaffer function (D = 2)

$$f(x_1, x_2) = 0.5 + \frac{\sin^2 \sqrt{x_1^2 + x_2^2} - 0.5}{\left[1 + 0.001 \left(x_1^2 + x_2^2\right)\right]^2},$$
(26)

where $x_d \in [-10, 10]$. The global minimum value is $f(x^*) = 0$ at $x^* = (0, 0)$. One characteristic of this function is that the global minimum is located very close to the local minima.

(4) Griewank function (D = 3, 4, 10)

$$f(\boldsymbol{x}) = 1 + \sum_{d=1}^{D} \frac{x_d^2}{4000} - \prod_{d=1}^{D} \cos\left(\frac{x_d}{\sqrt{d}}\right),$$
(27)

where $x_d \in [-10, 10]$. This is a multimodal function with global minimum $f(x^*) = 0$ at $x^* = (0, 0)$.

576 (5) Hartmann function (D = 6) A nonseparable function with multidimensional inputs and 577 multiple local minima.

$$f(\boldsymbol{x}) = -\sum_{j=1}^{4} c_j \exp\left(-\sum_{d=1}^{6} a_{jd} (x_d - b_{jd})^2\right),$$
(28)

578 where

$$\boldsymbol{A} = [a_{jd}] = \begin{bmatrix} 10 & 3 & 17 & 3.5 & 1.7 & 8\\ 0.05 & 10 & 17 & 0.1 & 8 & 14\\ 3 & 3.5 & 1.7 & 10 & 17 & 8\\ 17 & 8 & 0.05 & 10 & 0.1 & 14 \end{bmatrix}, \ \boldsymbol{c} = [c_j] = \begin{bmatrix} 1\\ 1.2\\ 3\\ 3.2 \end{bmatrix},$$

$$\boldsymbol{B} = [b_{jd}] = \begin{bmatrix} 0.1312 & 0.1696 & 0.5569 & 0.0124 & 0.8283 & 0.5586\\ 0.2329 & 0.4135 & 0.8307 & 0.3736 & 0.1004 & 0.9991\\ 0.2348 & 0.1451 & 0.3522 & 0.2833 & 0.3047 & 0.6650\\ 0.4047 & 0.8828 & 0.8732 & 0.5743 & 0.1091 & 0.0381 \end{bmatrix},$$

$$(29)$$

579 $x_d \in [0,1]$. The 6-dimensional case has 6 local minima, and the global minimum is $f(\boldsymbol{x}^*) = -3.32237$ at $\boldsymbol{x}^* = (0.20169, 0.150011, 0.476874, 0.275332, 0.311652, 0.657301).$

Note that all these minimization problems can be easily transformed as a maximization optimization problem, i.e., $\max - f(x)$.

583 E Supplementary results on benchmark test functions

584 E.1 Performance profiles

⁵⁸⁵ When computing the performance profiles (PPs), i.e., Dolan-Moré curves [30], we consider the ⁵⁸⁶ number of function evaluations to find the global optimum as the performance measure. Specifically, ⁵⁸⁷ let $t_{p,a}$ denote the number of evaluations used by method/solver *a* to reach the global solution in ⁵⁸⁸ experiment *p* (a lower value is better). The value is equal to $N_p + 100$ if the method cannot find the ⁵⁸⁹ global optimum with N_p being the given observation budget for experiment *p*. The performance ratio

$$\gamma_{p,a} = \frac{t_{p,a}}{\min\{t_{p,a} : a \in \mathcal{A}\}},\tag{30}$$

where \mathcal{A} represents the set that includes all comparing models, and the performance profile for each method is the distribution of $p(\gamma_{p,a} \le \rho)$ in terms of a factor ρ . We set $\rho = 1 : N_{\text{max}} + 1$, where $N_{\text{max}} = \max N_p$ is the largest observation budget assumed for the compared experiments. We define the problem set $\{\mathcal{P} \mid \forall p \in \mathcal{P}\}$ as the 10 runs for every function and draw the performance profiles of each model, also set \mathcal{P} as the 70 experiments in the 7 tested functions and estimate the overall performance profiles.

We show the obtained PPs across test functions (i.e., overall PPs) in Figure 2 (bottom right) for illustration (see Section 5.1 of the paper); the results of all the tested functions are shown below in Figure 8. Note we only compute PPs for the methods that compared on all test functions, i.e., deepGP and SAASBO are not considered.

Table 2 gives the AUC (area under the curve) values of these curves, where the AUC of the overall performance profiles is taken as the metric to compare the overall performances. As can be seen, BKTF obtains the best performance across all the considered functions. In addition, for most of the test functions, the AUC of grid-based baseline models is comparable with those of continuous GP-based models, suggesting that discretization of the continuous space is feasible to simplify the optimization problem.



Figure 8: Performance profiles on the standard test functions.

606 E.2 Quantitative comparison

We compare the last step *regret*, average costs of evaluations, and AUC of PPs of different methods on benchmark test functions in Table 2.

$\begin{array}{c} f(\boldsymbol{x}) \\ (D) \end{array}$		GP $\alpha_{\rm EI}$	GP α_{UCB}	GPgrid $\alpha_{\rm EI}$	GPgrid $\alpha_{\rm UCB}$	additive GP	BKTF random	BKTF	deepGP	SAAS BO
		0.01	0.01	0.31	0.24	0.05	0.00	0.00	0.02	_
	regret	±	±	±	±	±	±	±	±	
B (2)	-	0.01	0.01	0.62	0.64	0.09	0.00	0.00	0.02	
	Cost	≈ 44	≈ 42	≈ 23	≈ 36	≈ 100	≈ 47	\approx 4	≈ 47	-
	AUC	32.14	32.14	44.86	42.29	2.14	43.40	48.44	-	-
		2.00	2.00	1.60	2.00	2.00	0.60	0.00	3.00	-
	regret						\pm		±	
D (2)	G	0.00	0.00	0.80	0.00	0.00	0.92	0.00	0.007	
	Cost	-	-	-	-	-	≈ 48	≈ 5	-	-
	AUC	12.17	12.17	19.95	17.10	12.17	37.05	49.82	-	-
		0.02	0.02	0.10	0.09	0.03	0.00	0.00	0.08	-
	regret	$\begin{array}{c} \pm \\ 0.02 \end{array}$	$^{\pm}_{0.02}$	$^{\pm}$ 0.15	$^{\pm}_{0.07}$	$^{\pm}_{003}$	$^{\pm}_{000}$	$^{\pm}_{000}$	$^{\pm}_{0.04}$	
S (2)	Cost	~ 36	~ 14	> 50	> 50	~ 13	~ 54	~ 22	> 50	
		~ 30	\sim ++	21.74	20.74	\sim +3	~ 34	~ 22 10 00	/ 50	
	AUC	41.75	39.51	31.74	29.74	37.00	44.60	48.80	-	-
	regret	0.14	0.25	0.23	0.22	0.10	0.00	0.00 	-	-
C(2)		0.14	0.10	0.13	0.12	0.09	0.00^{\perp}	0.00		
G (3)	Cost	> 100	> 100	> 100	> 100	> 100	≈ 47	≈ 43	-	-
	AUC	48.90	47.69	47.69	47.69	47.69	50.44	50.78	-	-
		0.10	0.19	0.38	0.27	0.13	0.00	0.00	-	-
	regret	±	±	±	\pm	±	\pm	\pm		
G (4)		0.07	0.12	0.19	0.17	0.11	0.00	0.00		
	Cost	> 100	> 100	> 100	> 100	> 100	≈ 87	pprox 68	-	-
	AUC	79.32	77.61	77.61	77.61	78.17	80.01	80.40	-	-
		0.12	0.07	0.70	0.79	0.48	1e-5	0.00	-	0.19
	regret		±	±		±	±			±
H (6)	a		0.07	0.70	0.01	0.17	16-5	0.00		0.48
	Cost	> 100	> 100	> 100	> 100	> 100	≈ 154	≈ 60	-	> 100
	AUC	78.11	78.11	79.18	78.11	78.11	78.55	80.78	-	-
·		0.36	0.38	-	-	0.25	0.00	-	-	0.14
	regret	$\begin{bmatrix} \pm \\ 0.07 \end{bmatrix}$	010	-	-	0.30	0.00			\pm 0.10
G (10)	Cost	> 200	> 200			> 150	~ 124			> 200
		107.55	/ 200	-	-	/ 100	~ 124	-	-	> 200
	AUC	197.55	197.55	-	-	199.38	200.77	-	-	-
Overall	AUC	186.49	185.75	188.80	187.59	185.40	196.40	199.51	-	-

Table 2: Optimization on test functions: *regret* when n = N (mean \pm std.) / Average costs of evaluations / AUC of PPs.

Best results are highlighted in bold fonts. B: Branin; D: Damavandi; S: Schaffer; G: Griewank; H: Hartmann.

609 E.3 Interpretation of results

The basis functions learned by BKTF are interpretable. For example, Figure 9 shows the latent factors (r = 1) obtained at the last iteration of function evaluation in one run on 3D Griewank function. We see that BKTF can learn the periodicity (global structure) of the function benefited from the low-rank modeling. On the other hand, a stationary and separable GP cannot, other than using a specific kernel function such as the periodic kernel, which however requires strong prior knowledge to set the periodicity kernel hyperparameter.



Figure 9: Examples of latent factors learned by BKTF on 3D Griewank function.



Figure 10: Effects of hyperpriors on Branin function: Optimization with different hyperpriors.

616 E.4 Effects of hyperpriors

We compare the optimization performance on Branin function with different hyperprior settings in 617 Figure 10 as an example to illustrate the effects of hyperpriors. Specifically, we compare the optimiza-618 tion results under several hyperprior assumptions of μ_l , when τ_l^{-1} is set as 0.5. As can be seen, BKTF 619 is not able to reach the global minimum with too small or too large mean assumptions (comparable to 620 [0, 1]) on the kernel lengthscales l, for example in the cases where $\mu_l = \{\log(0.05), \log(2)\}$. In con-621 trast, it finds the global optimum after 4 iterations of function evaluations when $\mu_l = \log(0.5)$, see the 622 purple line. These imply the importance of hyperprior selection. The reason is that in the first several 623 evaluations, since the observations are rare, the prior basically determines the exploration-exploitation 624 balance and guides the search process. 625

Figure 11 shows the approximated posterior distributions for kernel hyperparameters and model noise variance when $\tau_l^{-1} = 0.5$, $\mu_l = \log (0.5)$. We see that for the re-scaled input space and normalized function output, the sampled length scales are around half of the input domain. Such settings are reasonable to capture the correlations between the observations and are also interpretable.

The effects of hyper-priors on other functions are similar, and we choose an appropriate setting relevant to the input range. The hyper-prior on τ impacts the uncertainty of the latent factors, for example a large model noise assumption allows more variances in the factors. The role of $\{a_0, b_0\}$ becomes more important when the objective function is complex that BKTF cannot well describe the function with limited observations. Generally, we select the priors that make the noise variances not quite large, such as the results of τ^{-1} shown in Figure 4 and Figure 11. An example of the uncertainty provided by BKTF is explained in Appendix C (see Figure 6).



Figure 11: Effects of hyperpriors on Branin function: Posterior probability distributions of lengthscales and model noise variance when $\tau_l^{-1} = 0.5$, $\mu_l = \log(0.5)$.

637 E.5 Effects of rank

Under a fully Bayesian treatment, kernel hyperparameters of BKTF are automatically sampled by
 MCMC with proper priors; the only selected parameter is the model rank. We test the effects of rank
 specification for the proposed BKTF surrogate on the 2D nonstationary nonseparable function defined
 in Section 1 Introduction (see Figure 1 and Appendix C). We use the same experiment settings as in
 Appendix C, i.e., 30 initial observations and 50 budget.

The results are given in Figure 12 and 13, where we compare the optimization performance of BKTF with rank $R = \{2, 4, 6, 8\}$ and two GP-based surrogate models: GP α_{EI} and GP α_{UCB} . To clearly illustrate the results, we only show the comparison on one run. In Figure 12, we compare the regret from different models, and in Figure 13 we compute and compare the mean CRPS (continuous ranked probability score) on the unobserved points.

CRPS is a widely applied metric for evaluating the performance of UQ for probabilistic models. With
 Gaussian likelihoods, CRPS can be defined as:

$$CRPS = -\frac{1}{n'} \sum_{i=1}^{n'} \sigma_i \left[\frac{1}{\sqrt{\pi}} - 2\psi \left(\frac{f_i - \hat{y}_i}{\sigma_i} \right) - \frac{f_i - \hat{y}_i}{\sigma_i} \left(2\Phi \left(\frac{f_i - \hat{y}_i}{\sigma_i} \right) - 1 \right) \right], \quad (31)$$

where n' is the number of unknown points in the defined space, i.e., $n' = \prod_{d=1}^{D} m_d - n$, \hat{y}_i and σ_i are the approximated posterior mean and std. for the *i*th data point, respectively, f_i denotes the true value for the *i*th point, and $\psi(\cdot)$ and $\Phi(\cdot)$ are the PDF (probability density function) and CDF (cumulative distribution function) of standard normal, respectively.

We see that BKTF successfully finds the global optimum with rank from 2 to 8, and obtains better (lower) CRPS values than GP baseline surrogates during the search processes. These results indicate that the proposed fully Bayesian framework is robust to the rank setting and can avoid overfitting.



Figure 12: Effects of rank specification on the test function defined in Introduction (see Appendix C, Eq. 23): Comparison of optimization performance.



Figure 13: Effects of rank specification on the test function defined in Introduction (see Appendix C, Eq. 23): Comparison of CRPS.

657 F Comparison of computational complexity

We compare the computational complexity of BKTF, BKTFrandom and the baseline methods applied on test functions in Table 3, where n is the number of observations, m_d is the number of interpolation points for the *d*th dimension, and n_{random} denotes the number of candidates we selected in BKTFrandom, which is 20k in the conducted experiments.

As can be seen, theoretically the proposed model BKTF/BKTFrandom has the lowest computational complexity.

Model	Complexity
GP $\alpha_{\rm EI}$	$\mathcal{O}\left(n^3 ight)$
GP α_{UCB}	$\mathcal{O}\left(n^{3} ight)$
GPgrid $\alpha_{\rm EI}$	$\mathcal{O}\left(\left(\prod_{d=1}^{D}m_{d} ight)^{3} ight)$
GPgrid α_{UCB}	$\mathcal{O}\left(\left(\prod_{d=1}^{D}m_{d} ight)^{3} ight)$
additive GP	$\mathcal{O}\left(n^3 ight)$
BKTF	$\min \left\{ \mathcal{O}\left(n^3\right), \mathcal{O}\left(\sum_{d=1}^D m_d^3\right) \right\}$
BKTFrandom	$\mathcal{O}\left(n_{\mathrm{random}}^{3}\right)$
deepGP	$\mathcal{O}\left(\left(\prod_{d=1}^{D}m_{d} ight)^{3} ight)$

Table 3: Comparison of model complexity.

664 G Hyperparameter tuning for machine learning

Table 4 lists all the hyperparameters in the tuning tasks.

Dataset	Algorithm	Hyperparameters	Туре	Search space
		no. of estimators	integer	[10, 100]
		max features	integer	[5, 50]
	RF classifier	min samples split	integer	[1, 04]
		min samples leaf	integer	[2, 11]
		criterion	categorical	gini entrony
MNIST		critchon	categoricai	giiii, chuopy
		neurons	integer	[10, 100]
	NN classifier	batch size	integer	[16, 64]
		epochs	integer	[20, 50]
		patience	integer	[3, 20]
		optimizer	categorical	adam, rmsprop, sgd
		activation	categorical	relu, tanh
		no. of estimators	integer	[10, 100]
	PF regressor	max depth	integer	[5, 50]
	KI legiessoi	max features	integer	[1, 13]
Boston housing		min samples split	integer	[2, 11]
		neurons	integer	[10, 100]
	NN regressor	batch size	integer	[16, 64]
		epochs	integer	[20, 50]

Table 4: Hyperparameters of the tested ML algorithms.

666 H Supplementary results on ML hyperparameter tuning

We summarize the final accuracy obtained with different BO methods for MNIST classification and the final MSE for the regression task in Table 5.

			()	00
	(a) M	NIST	(b) Bosto	on housing
	RF classifier	NN classifier	RF regressor	NN regressor
RS	96.09 ± 0.05	97.59 ± 0.37	26.10 ± 0.45	40.43 ± 3.91
PSO	95.84 ± 0.14	97.54 ± 0.51	26.07 ± 0.44	44.40 ± 6.29
GP $\alpha_{\rm EI}$	96.10 ± 0.07	98.04 ± 0.05	26.19 ± 0.45	38.46 ± 3.31
GP α_{UCB}	96.06 ± 0.05	97.46 ± 0.64	26.34 ± 0.35	36.78 ± 1.91
BO-TPE	96.10 ± 0.06	97.52 ± 0.06	26.27 ± 0.31	36.40 ± 4.72
BKTF	-	-	$\textbf{25.03} \pm 0.18$	$\textbf{30.84} \pm 1.13$
BKTFrandom	96.38 \pm 0.04	$\textbf{98.16} \pm 0.09$	-	-

Table 5: Final accuracy for (a) MNIST classification and MSE for (b) Boston housing regression.

The values are presented as mean \pm std. Best results are highlighted in bold fonts.

669 NeurIPS Paper Checklist

670	1. Claims
671	Question: Do the main claims made in the abstract and introduction accurately reflect the
672	paper's contributions and scope?
673	Answer: [Yes]
674	Justification: Please see Abstract and Section Introduction for more details.
675	Guidelines:
676	• The answer NA means that the abstract and introduction do not include the claims
677	made in the paper.
678	• The abstract and/or introduction should clearly state the claims made, including the
679 680	contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
681	• The claims made should match theoretical and experimental results, and reflect how
682	much the results can be expected to generalize to other settings.
683	• It is fine to include aspirational goals as motivation as long as it is clear that these goals
684	are not attained by the paper.
685	2. Limitations
686	Question: Does the paper discuss the limitations of the work performed by the authors?
687	Answer: [Yes]
688	Justification: We discuss the limitations in Section 6 Conclusion.
689	Guidelines:
690	• The answer NA means that the paper has no limitation while the answer No means that
691	the paper has limitations, but those are not discussed in the paper.
692	• The authors are encouraged to create a separate "Limitations" section in their paper.
693	• The paper should point out any strong assumptions and how robust the results are to
694	violations of these assumptions (e.g., independence assumptions, noiseless settings,
695	model well-specification, asymptotic approximations only holding locally). The authors
696 697	should reflect on how these assumptions might be violated in practice and what the implications would be
097	• The authors should reflect on the scope of the claims made e.g. if the approach was
698	only tested on a few datasets or with a few runs. In general empirical results often
700	depend on implicit assumptions, which should be articulated.
701	• The authors should reflect on the factors that influence the performance of the approach
702	For example, a facial recognition algorithm may perform poorly when image resolution
703	is low or images are taken in low lighting. Or a speech-to-text system might not be
704	used reliably to provide closed captions for online lectures because it fails to handle
705	technical jargon.
706	• The authors should discuss the computational efficiency of the proposed algorithms
707	and how they scale with dataset size.
708	• If applicable, the authors should discuss possible limitations of their approach to
709	address problems of privacy and fairness.
710	• While the authors might fear that complete honesty about limitations might be used by
711	reviewers as grounds for rejection, a worse outcome might be that reviewers discover
712	limitations that aren't acknowledged in the paper. The authors should use their best
/13	judgment and recognize that individual actions in layor of the community. Pavious
714 715	will be specifically instructed to not penalize honesty concerning limitations
. 10	2 Theory Assumptions and Proofs
/16	5. Theory Assumptions and Froois Ouestion: For each theoretical result, does the paper provide the full set of assumptions and
717 718	a complete (and correct) proof?
719	Answer: [Yes]

720 721	Justification: We illustrate methodology and technical details of the proposed model in Section 3 BKTF for Bayesian optimization and Appendix A Model inference.
722	Guidelines:
723	• The answer NA means that the paper does not include theoretical results
723	• All the theorems, formulas, and proofs in the paper should be numbered and cross
724	• All the theorems, formulas, and proofs in the paper should be numbered and cross-
725	 All assumptions should be clearly stated or referenced in the statement of any theorems.
726	• All assumptions should be clearly stated of referenced in the statement of any theorems.
727	• The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the suthers are appeared to provide a short
728	proof sketch to provide intuition
729	Invested to provide intuition.
730	• Inversely, any informat proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material
700	Theorems and Lemmas that the proof ralies upon should be properly referenced
/32	Theorems and Lemmas that the proof refless upon should be property referenced.
733	4. Experimental Result Reproducibility
734	Question: Does the paper fully disclose all the information needed to reproduce the main ex-
735	perimental results of the paper to the extent that it affects the main claims and/or conclusions
736	of the paper (regardless of whether the code and data are provided or hot)?
737	Answer: [Yes]
738	Justification: We illustrate detailed information for experiment implementation in Section 5
739	Experiments and Appendices C-G.
740	Guidelines:
741	• The answer NA means that the paper does not include experiments.
742	• If the paper includes experiments, a No answer to this question will not be perceived
743	well by the reviewers: Making the paper reproducible is important, regardless of
744	whether the code and data are provided or not.
745 746	• If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
747	• Depending on the contribution, reproducibility can be accomplished in various ways
748	For example, if the contribution is a novel architecture, describing the architecture fully
749	might suffice, or if the contribution is a specific model and empirical evaluation, it may
750	be necessary to either make it possible for others to replicate the model with the same
751	dataset, or provide access to the model. In general. releasing code and data is often
752	one good way to accomplish this, but reproducibility can also be provided via detailed
753	instructions for how to replicate the results, access to a hosted model (e.g., in the case
754	of a large language model), releasing of a model checkpoint, or other means that are
755	appropriate to the research performed.
756	• While NeurIPS does not require releasing code, the conference does require all submis-
757	sions to provide some reasonable avenue for reproducibility, which may depend on the
758	() If d = (1 + 1) + (1 +
759	(a) If the contribution is primarily a new algorithm, the paper should make it clear now to reproduce that algorithm
760	(b) If the contribution is primarily a new model architecture, the paper should describe
761	(b) If the contribution is primarily a new moder architecture, the paper should describe the architecture clearly and fully
762	(c) If the contribution is a new model (e.g., a large language model), then there should
764	either be a way to access this model for reproducing the results or a way to reproduce
765	the model (e.g., with an open-source dataset or instructions for how to construct
766	the dataset).
767	(d) We recognize that reproducibility may be tricky in some cases, in which case
768	authors are welcome to describe the particular way they provide for reproducibility.
769	In the case of closed-source models, it may be that access to the model is limited in
770	some way (e.g., to registered users), but it should be possible for other researchers
771	to have some path to reproducing or verifying the results.
772	5. Open access to data and code

773 774 775	Question: Does the paper provide open access to the data and code, with sufficient instruc- tions to faithfully reproduce the main experimental results, as described in supplemental material?
776	Answer: [Yes]
777	Institution: We provide the code for the 2D test function in supplementary material
770	Guidelines:
//8	The second state of the se
779	• The answer NA means that paper does not include experiments requiring code.
780	• Please see the NeurIPS code and data submission guidelines (https://nips.cc/
701	• While we encourage the release of code and data, we understand that this might not be
783	possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not
784	including code, unless this is central to the contribution (e.g., for a new open-source
785	benchmark).
786	• The instructions should contain the exact command and environment needed to run to
787 788	reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
789	• The authors should provide instructions on data access and preparation, including how
790	to access the raw data, preprocessed data, intermediate data, and generated data, etc.
791	• The authors should provide scripts to reproduce all experimental results for the new
792	proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why
793	• At submission time, to preserve anonymity, the authors should release anonymized
794 795	versions (if applicable).
796	• Providing as much information as possible in supplemental material (appended to the
797	paper) is recommended, but including URLs to data and code is permitted.
798	6. Experimental Setting/Details
799	Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
800	parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
801	results?
802	Answer: [Yes]
803 804	Justification: We specify the experimental setting and implementation details in Section 5 Experiments and Appendices C-G.
805	Guidelines:
806	• The answer NA means that the paper does not include experiments.
807	• The experimental setting should be presented in the core of the paper to a level of detail
808	that is necessary to appreciate the results and make sense of them.
809	• The full details can be provided either with the code, in appendix, or as supplemental
810	material.
811	7. Experiment Statistical Significance
812	Question: Does the paper report error bars suitably and correctly defined or other appropriate
813	information about the statistical significance of the experiments?
814	Answer: [Yes]
815	Justification: We repeat the experiments certain times and report the mean with std. results
816	in Figure 1, 2, 3, and Table 2, 5.
817	Guidelines:
818	• The answer NA means that the paper does not include experiments.
819	• The authors should answer "Yes" if the results are accompanied by error bars, confi-
820	dence intervals, or statistical significance tests, at least for the experiments that support
821	the main claims of the paper.
822	• The factors of variability that the error bars are capturing should be clearly stated (for
823	example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions)
024	run with given experimental conditions).

825 826		• The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
827		• The assumptions made should be given (e.g., Normally distributed errors).
828		 It should be clear whether the error bar is the standard deviation or the standard error.
829		of the mean.
830		• It is OK to report 1-sigma error bars, but one should state it. The authors should
831		preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
832		of Normality of errors is not verified.
833		• For asymmetric distributions, the authors should be careful not to show in tables or
834		figures symmetric error bars that would yield results that are out of range (e.g. negative
835		error rates).
836		• If error bars are reported in tables or plots, The authors should explain in the text how
837		they were calculated and reference the corresponding figures or tables in the text.
838	8.	Experiments Compute Resources
839		Question: For each experiment, does the paper provide sufficient information on the com-
840		puter resources (type of compute workers, memory, time of execution) needed to reproduce
841		the experiments?
842		Answer: [Yes]
843		Justification: All the experiments can be performed with a 16-core 2.40 GHz CPU and 32
844		GB RAM.
845		Guidelines:
846		• The answer NA means that the paper does not include experiments.
847		• The paper should indicate the type of compute workers CPU or GPU, internal cluster,
848		or cloud provider, including relevant memory and storage.
849		• The paper should provide the amount of compute required for each of the individual
850		experimental runs as well as estimate the total compute.
851		• The paper should disclose whether the full research project required more compute
852 853		than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).
854	9.	Code Of Ethics
855		Question: Does the research conducted in the paper conform, in every respect, with the
856		NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?
857		Answer: [Yes]
858		Justification: The research conducted in this paper conform with the NeurIPS code of ethics.
859		Guidelines:
860		• The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
861		• If the authors answer No, they should explain the special circumstances that require a
862		deviation from the Code of Ethics.
863		• The authors should make sure to preserve anonymity (e.g., if there is a special consid-
864		eration due to laws or regulations in their jurisdiction).
865	10.	Broader Impacts
866		Question: Does the paper discuss both potential positive societal impacts and negative
867		societal impacts of the work performed?
868		Answer: [Yes]
869		Justification: This paper presents work whose goal is to advance the field of probabilistic
870		Machine Learning, particularly Bayesian Optimization. We discussed such impacts in
871		Section 6 Conclusion.
872		Guidelines:
873		• The answer NA means that there is no societal impact of the work performed.
874		• If the authors answer NA or No, they should explain why their work has no societal
875		impact or why the paper does not address societal impact.

876 877 878 879 880 881 882 883 884 885 886 887 888 889 890 891 892 893		 Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations. The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster. The authors should consider possible harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology. If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).
895	11.	Safeguards
896 897 898	11.	Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?
899		Answer: [Yes]
900 901 902 903		Justification: This work has the potential of misuse for machine learning algorithms. How- ever the current model has certain limitations on applying for high-dimensional problems, thus such risks are low. We mentioned such risks in the last paragraph in Section 6 Conclu- sion.
904		Guidelines:
905		• The answer NA means that the paper poses no such risks.
906 907 908 909		 Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
910 911		• Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
912 913 914		• We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.
915	12.	Licenses for existing assets
916 917 918		Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?
919		Answer: [Yes]
920		Justification: We include the URLs for the datasets we used in this work.
921		Guidelines:
922		• The answer NA means that the paper does not use existing assets.
923		• The authors should cite the original paper that produced the code package or dataset.
924 925		• The authors should state which version of the asset is used and, if possible, include a URL.
926		• The name of the license (e.g., CC-BY 4.0) should be included for each asset.
927 928		• For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

929 930 931 932		• If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
933 934		• For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
935 936		 If this information is not available online, the authors are encouraged to reach out to the asset's creators.
937	13.	New Assets
938 939		Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?
940		Answer: [Yes]
941		Justification: We submit partial of the code in supplementary material and select a license
942		when submitting the paper.
943		Guidelines:
944		• The answer NA means that the paper does not release new assets.
945 946 947		• Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
948		 The paper should discuss whether and how consent was obtained from people whose asset is used.
949		• At submission time, remember to anonymize your assets (if applicable). You can either
951		create an anonymized URL or include an anonymized zip file.
952	14.	Crowdsourcing and Research with Human Subjects
953 954 955		Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?
956		Answer: [NA]
957		Justification: This paper does not involve crowdsourcing nor research with human subjects.
958		Guidelines:
959 960		• The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
961 962		• Including this information in the supplemental material is fine, but if the main contribu- tion of the paper involves human subjects, then as much detail as possible should be included in the main paper.
964 965		 According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data
966		collector.
967 968	15.	Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects
969		Question: Does the paper describe potential risks incurred by study participants, whether
970		such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
971		approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?
973		Answer: [NA]
974		Justification: This paper does not involve crowdsourcing nor research with human subjects.
975		Guidelines:
976		• The answer NA means that the paper does not involve crowdsourcing nor research with
977		numan subjects.
978 979 980		may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

981	• We recognize that the procedures for this may vary significantly between institutions
982	and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
983	guidelines for their institution.
984	• For initial submissions, do not include any information that would break anonymity (if
985	applicable), such as the institution conducting the review.