

1 **Constrained mixed-integer Gaussian mixture**  
2 **Bayesian optimization and its applications in**  
3 **designing fractal and auxetic metamaterials**

4 Anh Tran · Minh Tran · Yan Wang

5  
6 Received: date / Accepted: date

7 **Abstract** Bayesian optimization (BO) is a global optimization method that  
8 has the potential for design optimization. However, in classical BO algorithm,  
9 the variables are considered as continuous. In real-world engineering problems,  
10 both continuous and discrete variables are present. In this work, an efficient  
11 approach to incorporate discrete variables to BO is proposed. In the proposed  
12 constrained mixed-integer BO method, the sample set is decomposed into  
13 smaller clusters during sequential sampling, where each cluster corresponds  
14 to a unique ordered set of discrete variables, and a Gaussian process regres-  
15 sion (GP) metamodel is constructed for each cluster. The model prediction  
16 is formed as the Gaussian mixture model, where the weights are computed  
17 based on the pair-wise Wasserstein distance between clusters, and gradually  
18 converge to an independent GP as the optimization process advances. The  
19 definition of neighborhood can be flexibly and manually defined to account  
20 for independence between clusters, such as in the case of categorical variables.  
21 Theoretical results are provided in concert with two numerical and engineer-  
22 ing examples, and two examples for metamaterial developments, including one  
23 fractal and one auxetic metamaterials, where the effective properties depends  
24 on both the geometry and the bulk material properties.

25 **Keywords** Bayesian optimization · Gaussian process · constrained ·  
26 mixed-integer · metamaterials

27 **1 Introduction**

28 Designing materials is to identify structures at micro- and nano-scales to  
29 achieve the desirable properties. The major process of design is to establish

---

Anh Tran, Minh Tran, Yan Wang  
George W. Woodruff School of Mechanical Engineering  
Georgia Institute of Technology  
E-mail: yan.wang@me.gatech.edu

30 structure-property relationships, based on which design optimization can be  
31 performed. Simulation tools at multiple scales (from atomistic to continuum)  
32 have been developed to accelerate this process. Nevertheless, the major tech-  
33 nical challenges of efficiency and accuracy still exist. The first one is searching  
34 in high-dimensional design space to find the global optimum of material com-  
35 positions and structural configurations. The second one is the uncertainty as-  
36 sociated with the high-dimensional structure-property relationships, which are  
37 usually constructed as surrogate models or metamodels. Particularly, aleatory  
38 uncertainty can be linked to natural randomness of materials (e.g. grain sizes  
39 and grain shapes in polycrystalline materials). Epistemic uncertainty is mainly  
40 due to approximations and numerical treatments in surrogates and simulation  
41 models. Methods of searching globally for optimal and robust solutions are  
42 needed.

43 Bayesian optimization (BO) is a metamodel-based methodology to seek  
44 for the global optimal solution under uncertainty in the search space with  
45 sequential sampling. Compared to other bio-inspired global optimization al-  
46 gorithms, such as ant colony systems, particle swarm, and genetic algorithm  
47 (GA), it has the advantage of maintaining the global search history by con-  
48 structing a metamodel to approximate the objective function. Typically the  
49 metamodel is based on the Gaussian process (GP) method, and actively up-  
50 dated as more samples are collected. However, in current formulation of GP,  
51 input variables are restricted to be continuous. In real-world engineering prob-  
52 lems, input design variables and parameters can be categorical or discrete. For  
53 example, binary variables can be used to enable or disable a design feature.  
54 The number of features has integer values. Therefore extending BO method  
55 to accommodate discrete variables is an important topic for solving real-world  
56 problems.

57 Another major issue that prohibits the BO and GP framework is its lack of  
58 scalability in searching the high-dimensional space when the number of input  
59 variables is large. The required number of sample points grows exponentially as  
60  $\mathcal{O}(s^d)$  with respect to the dimension of search space  $d$ , where  $s$  is the number of  
61 sampling point for each dimension. The phenomenon is referred to as the curse-  
62 of-dimensionality in literature. As a result, the size of the covariance matrix in  
63 GP also grows exponentially with respect to the dimensionality, creating the  
64 computational bottleneck in computing the inverse of the covariance matrix.

65 In this paper, a new BO method is proposed for constrained mixed-integer  
66 optimization problems to incorporate discrete design variables into the BO  
67 algorithm. In the proposed method, the large dataset of samples is decomposed  
68 into smaller clusters, where each cluster corresponds to a unique combination  
69 of discrete variable values, which is referred to as a discrete tuple. A GP is  
70 then constructed within each cluster. During the search and metamodel update  
71 processes, the mean and variance predictions are formulated as a Gaussian  
72 mixture model, where the weighted average predictions are combined from  
73 those of neighboring clusters, based on the pair-wise distance between the main  
74 and the neighboring clusters. The neighborhood of each cluster is constructed  
75 only once during the initialization.

76 Because of the decomposition approach, the number of sampling points to  
77 construct each cluster is significantly reduced compared to the whole dataset,  
78 and the GP thus is faster to construct for each cluster. This approach, however,  
79 leads to an undesirable effect of sparsity within each GP cluster. As a result, the  
80 posterior variance might be slightly overestimated. To circumvent the sparsity  
81 effect of the decomposition approach, a weighted average scheme is adapted  
82 to "borrow" the sampling points from neighboring clusters, where the discrete  
83 tuples of the neighbors slightly differ from the discrete tuple of the original  
84 cluster. The definition of neighborhood is completely controlled by users, and  
85 neighbors can be added or removed accordingly. The unique advantage of the  
86 proposed method is that the optimization problem of both continuous and  
87 discrete variables and the acceleration of GP for high-dimensional problems  
88 are solved simultaneously. Theoretical results are provided and discussed in  
89 concert with computational metamaterials design applications.

90 In the remainder of the paper, Section 2 provides a literature review for BO  
91 methodology, its extension, such as constrained and mix-integer optimization  
92 problems, and its applications. Section 3 describes the proposed constrained  
93 mixed-integer BO algorithm using Gaussian mixture model, including theoret-  
94 ical analysis of algorithmic complexity as well as lower and upper bounds of  
95 the predictions. The methodology is demonstrated with applications in com-  
96 putational design of metamaterials. Metamaterials are an emerging class of  
97 engineered materials that exhibit interesting and desirable macroscopic prop-  
98 erties, which can be tailored, because of their engineered geometric structures  
99 rather than the material composition. In Section 4, the proposed method is  
100 verified using an analytical function that is modified based on a discrete version  
101 of Rastrigin function, an engineering example of welded beam design, where  
102 the discrete variables encode the material selection and design configuration of  
103 the beam. In the first engineering example of Section 5.1, we focus on design-  
104 ing high-strength and low-weight fractal metamaterials, where the effective  
105 material properties, such as effective Young's modulus is obtained using finite  
106 element method (FEM). In the second engineering example of Section 5.2, the  
107 method is demonstrated using an auxetic metamaterials for polymers, where  
108 the effective negative Poisson's ratio is optimized. Section 6 includes the dis-  
109 cussion of the limitations in the proposed approach, and Section 7 concludes  
110 the paper, respectively.

## 111 2 Related work

112 Here, we conduct a literature review on related BO work and its design ap-  
113 plications. In Section 2.1, the widely used acquisition functions for BO are  
114 introduced. The constrained optimization problem in BO is reviewed in sec-  
115 tion 2.2. In Section 2.3, the mixed-integer optimization problem in BO and  
116 its relate work is discussed. In Section 2.4, the applications of GP in design  
117 optimization is provided.

## 2.1 Acquisition function

BO is a metamodel-based optimization framework that uses GP as the metamodel. The major difference between BO and GP based optimization is the sampling strategy to construct the metamodel. The significant extension of BO is the implementation of a so-called acquisition function that dictates the location of the next sampling design site. This acquisition function reconciles the trade-off between exploration (navigating to the most uncertain region) and exploitation (driving the solution to the optimum) in the optimization process.

Given the objective function  $y = f(\mathbf{x})$ , the acquisition function  $a(\mathbf{x}; \{\mathbf{x}_i, y_i\}_{i=1}^N, \theta)$  depends on previous  $N$  observations or samples  $\{\mathbf{x}_i, y_i\}_{i=1}^N$  and GP hyperparameters  $\theta$ , and must be defined to strike a balance between exploration and exploitation. In exploration, the acquisition function  $a$  would lead to the next sampling point in an unknown region where the posterior variance  $\sigma^2(\mathbf{x})$  is large. In exploitation, the acquisition function  $a$  would result in the next sampling point where posterior mean  $\mu(\mathbf{x})$  is large for a maximization problem (or small for minimization). There are mainly three types of acquisition functions: probability of improvement (PI), expected improvement (EI), and upper confidence bound (UCB). They are defined as follows.

Let  $\mathbf{x}_{\text{best}} = \arg \max_{\mathbf{x}_i} f(\mathbf{x}_i)$  be the best sample achieved so far during sequential sampling for a maximization problem,  $\phi(\cdot)$  and  $\Phi(\cdot)$  be the probability density function and cumulative distribution function of the standard normal distribution respectively. The PI acquisition function [33] is defined as

$$a_{\text{PI}}(\mathbf{x}; \{\mathbf{x}_i, y_i\}_{i=1}^N, \theta) = \Phi(\gamma(\mathbf{x})), \quad (1)$$

where

$$\gamma(\mathbf{x}) = \frac{\mu(\mathbf{x}; \{\mathbf{x}_i, y_i\}_{i=1}^N, \theta) - f(\mathbf{x}_{\text{best}})}{\sigma(\mathbf{x}; \{\mathbf{x}_i, y_i\}_{i=1}^N, \theta)}, \quad (2)$$

indicates the deviation away from the best sample. The EI acquisition function [32][21] is mathematically expressed as

$$a_{\text{EI}}(\mathbf{x}; \{\mathbf{x}_i, y_i\}_{i=1}^N, \theta) = \sigma(\mathbf{x}; \{\mathbf{x}_i, y_i\}_{i=1}^N, \theta) \cdot (\gamma(\mathbf{x})\Phi(\gamma(\mathbf{x})) + \phi(\gamma(\mathbf{x}))) \quad (3)$$

Recently, Srinivas et al. [52][53] proposed a new form of UCB acquisition function,

$$a_{\text{UCB}}(\mathbf{x}; \{\mathbf{x}_i, y_i\}_{i=1}^N, \theta) = \mu(\mathbf{x}; \{\mathbf{x}_i, y_i\}_{i=1}^N, \theta) + \kappa \sigma(\mathbf{x}; \{\mathbf{x}_i, y_i\}_{i=1}^N, \theta), \quad (4)$$

where  $\kappa$  is a hyperparameter describing the exploitation-exploration balance.

## 2.2 Constrained BO

Constrained BO is a natural and important extension of the classical BO method. Constrained optimization problems based on engineering model and

simulation can be classified as two types: known and unknown constraints. The known constraints, or a priori constraints, are the ones known before the simulation, and thus can be evaluated independently without running simulations. On the other hand, the unknown constraints are the ones that are unpredictable without running the simulation, and thus can be only incorporated once the simulation is over, e.g. no solution because of numerical divergence. Generally speaking, the unknown constraints are more difficult to assess because it involves handling the classification problem, satisfied or violated, with respect to the optimization problem.

Digabel and Wild [9] summarized and provided a systematic classification and taxonomy for constrained optimization problem. Gardner et al. [11] proposed a penalized acquisition function approach to limit the searching space for the next sampling location. Gelbart et al. [12] suggested an entropy search criterion to search for the next sampling point under the formulation of the EI acquisition function. Hernández-Lobato et al. [19] [20] introduced a predictive entropy search and predictive entropy search with constraints, respectively, which maximizes the expected information gained with respect to the global maximum. Rehman and Langelaar [46] modeled constraints as a simple model and incorporated probability of feasibility measure to alternate the EI acquisition function. Li et al. [26] proposed a sequential Monte Carlo approach with radial basis function as surrogate model to solve for the constrained optimization problem.

### 2.3 Mixed-integer Bayesian optimization

The BO extension to mixed-integer problems is rather limited, partly because mixed-integer problems carry difficulties from both discrete and continuous optimization problems. Another approach is that the discrete optimization can be converted to continuous optimization, using simple rounding operation. The approach is not mathematically rigorous, but is still widely accepted in practice. Here we review several contributions in term of methodology to incorporate discrete variables.

Davis and Ierapetritou [7] combined a branch-and-bound approach with BO method to solve the mixed-integer optimization problems. Müller et al. [35,36,34] introduced three algorithms, which are Surrogate Optimization-Mixed Integer [35], Surrogate Optimization-Integer [36], and Mixed-Integer Surrogate Optimization [34], which differ in the perturbation sampling strategies and utilize GP as the surrogate model, to solve for the mixed-integer nonlinear problems. Hemker et al. [18] compared the performance of a GA, the implicit filtering algorithm, and a branch-and-bound approach formulated on BO algorithm to solve for a set of constrained mix-integer problems in groundwater management.

For mixed-integer extension for GP, van Stein et al. [54] proposed a distributed kriging approach, where the dataset is decomposed for continuous variables using  $k$ -mean algorithm, and the optimal weights are computed based

193 on the inverse posterior variance of each cluster. Gramacy et al. [15] [14] [16]  
 194 developed a treed GP that is naturally extensible to handle discrete variables.  
 195 In the case of discrete variables, the GP is one-hot encoded by the binary  
 196 combination of the discrete variables. Storlie et al. [55] developed the adap-  
 197 tive component selection shrinkage operato method (ACOSSO) extended from  
 198 Lin and Zhang [28] [27], which uses the smoothing spline ANOVA decompo-  
 199 sition to decompose the total variance to multivariate functions. Qian et al.  
 200 [42] [61] approached the mixed-integer problem from the covariance kernel of  
 201 GP, proposing the exchange correlation, the multiplicative correlation, and the  
 202 unrestricted correlation functions to handle discrete variable that is reminis-  
 203 cent of categorical regression. Swiler et al. [56] compared three above methods  
 204 and concluded that GP with special correlation kernel [42] [61] performs most  
 205 consistently among the test functions.

## 206 2.4 GP-based design optimization

207 GP, also known as kriging, has been widely applied in constructing surrogates  
 208 or metamodels for design optimization. Simpson et al. [49], Queipo et al. [43],  
 209 Martins and Lambe [30], Sóbester et al. [50], and Viana et al. [59] provided  
 210 comprehensive reviews on the use of kriging and other surrogate models for  
 211 multi-disciplinary design optimization. More recently, Li et al. [25] proposed  
 212 a kriging metamodel assisted multi-objective GA to solve multi-objective op-  
 213 timization problems. Jang et al. [22] used dynamic kriging to solve a design  
 214 optimization in fluid-solid interaction. Zhang et al. [60] also used kriging to  
 215 approximate the pump performance and optimize two objective functions with  
 216 respect to four design variables. Kim et al. [23] optimized and verified a fluid  
 217 dynamic bearings simulation using kriging approach. Kim et al. [24] applied  
 218 multi-fidelity kriging and optimized film-cooling hole arrangement. Liu et al.  
 219 [29] employed surrogate-based parallel optimization method to reduce the com-  
 220 putational time for a computational fluid dynamics problem with six design  
 221 variables. Song et al. [51] used a gradient-enhanced hierarchical kriging to  
 222 optimize drag on airfoils at a specified angle of attack. Zhou et al. [62][63]  
 223 developed a multi-fidelity kriging scheme to approximate the lift coefficient as  
 224 a function of Mach number and angle of attack in airfoils with computational  
 225 fluid dynamics analysis.

226 In the above work, design variables are all continuous. Compared to these  
 227 GP-based optimization, BO formulation provides a more generic and robust  
 228 searching procedure.

## 229 3 Proposed mixed-integer Bayesian optimization

230 The proposed mixed-integer BO based on distributed GP provides an efficient  
 231 searching method for large scale design problems, where design variables can  
 232 be either continuous or discrete. The discrete variables include both categorical

233 and integer variables, regardless of the existence of order relations. Let  $\mathbf{x} =$   
 234  $(\mathbf{x}^{(d)}, \mathbf{x}^{(c)})$  be the design variables, where  $\mathbf{x}^{(d)} \in \mathbb{D}$  are discrete variables in  
 235  $n$ -dimensional space  $\mathbb{D}$  and  $\mathbf{x}^{(c)} \in \mathbb{R}^{m-n}$  are continuous variables in  $(m-n)$ -  
 236 dimensional space  $\mathbb{R}^{m-n}$ . Together, they form a vector of design variables in  
 237 the  $m$ -dimensional space  $\mathcal{X}$ . Let  $f(\mathbf{x})$  be the objective function. The design  
 238 optimization problem solves the maximization problem

$$\mathbf{x}^* = \arg \max_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}), \quad (5)$$

239 subject to some inequality constraints

$$g_i(\mathbf{x}) \leq 0, i = 1, \dots, i_c \quad (6)$$

240 where  $i_c$  is the number of inequality constraints.

241 Here the notation for the rest of the paper is as follows.  $\mu_l(\mathbf{x})$  is used to  
 242 denote the posterior mean of the  $l^{\text{th}}$ -cluster at the query point  $\mathbf{x}$ .  $\hat{\mu}$  is the  
 243 prediction formed by Gaussian mixture model of all the clusters.  $\bar{\mu}_l$  is the  
 244 mean of the  $l^{\text{th}}$ -cluster.

245 In the proposed mixed-integer BO, the large dataset of observations is de-  
 246 composed into smaller local clusters, where each cluster is used to construct  
 247 a local GP. Because the large dataset has been decomposed and the number  
 248 of data points has reduced, the prediction within each cluster is not as ac-  
 249 curate, and can be improved by "borrowing" from neighboring dataset under  
 250 a weighted average scheme. The large dataset with continuous and discrete  
 251 variables can be decomposed to finitely many clusters, according to the tuple  
 252 of discrete variables. In each cluster, the data points share the same discrete  
 253 variable values. The classical GP approach is then applied to the dataset in  
 254 each cluster to construct a GP model.

255 Because of the decomposition scheme, the number of data points within  
 256 each cluster is reduced, compared to the number of data points of the whole  
 257 dataset. This leads to a sparser dataset within a cluster, and the posterior  
 258 variance is enlarged. To improve the prediction, the datasets from neighbor-  
 259 ing clusters are initially "borrowed" to improve the prediction on the tuple  
 260 of continuous variables  $\mathbf{x}^{(c)} \in \mathbb{R}^{m-n}$ , where the "borrowed" data points are  
 261 gradually eliminated as the optimization process converges via the weight com-  
 262 putation algorithm. On the other hand, the sparsity induced by the decom-  
 263 position scheme reduces the cost of computing the inverse of the covariance  
 264 matrix. In this weighted average scheme, the weights are computed and pe-  
 265 nalized based on the pair-wise Wasserstein distance between clusters, as well  
 266 as the posterior variance of the cluster to obtain a more accurate predictions  
 267 to aid in the convergence of the optimization process.

268 Figure 1 presents an overview of the workflow for the proposed mixed-  
 269 integer BO method in this paper. First, initial samples, typically obtained  
 270 from Monte Carlo or Latin hypercube sampling, are used to construct the  
 271 metamodel, where a local GP is associated with each individual cluster. Next,  
 272 a next sampling point is located within each cluster according to its acquisition  
 273 functional value. Then, a global sampling point for all clusters is determined

274 among the collection of all the next sampling points from each cluster. The  
 275 objective function is then called to evaluate at the global sampling location. A  
 276 local GP is updated at the cluster corresponding to the global sampling point.  
 277 A new local sampling point is located within the same cluster, and the process  
 278 repeats until some optimization criteria are met.

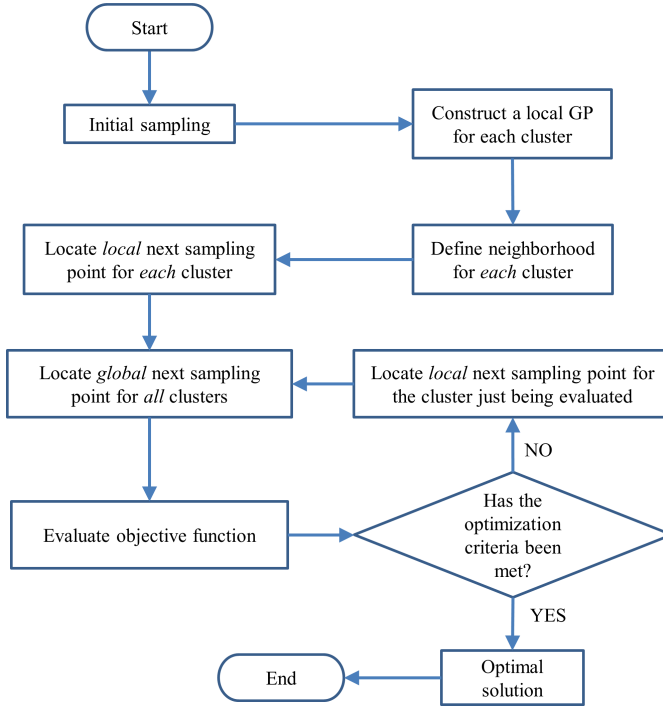


Fig. 1: Overall workflow of the proposed mixed-integer Bayesian optimization.

279 The following subsections are organized as follows. Section 3.1 briefly re-  
 280 views the GP formulation. Section 3.2 discusses the enumeration algorithm  
 281 for clusters and the discrete tuple. Section 3.3 describes the definition of cluster  
 282 neighborhood that is used to form a Gaussian mixture model. Section 3.4  
 283 details the weight computations for each individual cluster in the Gaussian  
 284 mixture model. Section 3.5 presents the computation of posterior mean and  
 285 posterior variance of the Gaussian mixture model. Section 3.6 describes the pen-  
 286 alized scheme to incorporate constraints into the acquisition function. Section  
 287 3.7 analyzes the theoretical bounds and computational cost of the proposed  
 288 mixed-integer BO method.



### 289 3.1 Gaussian process

290 We follow the notation introduced by Shahriari et al. [47] to briefly introduce  
 291 GP formulation for continuous variables.  $\text{GP}(\mu_0, k)$  is a nonparametric model  
 292 that is characterized by its prior mean  $\mu_0 : \mathcal{X} \mapsto \mathbb{R}$  and its covariance kernel  
 293  $k : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ . Define  $f_i = f(\mathbf{x}_i)$  and  $\mathbf{y}_{1:N}$  as the unknown function values  
 294 and noisy observations, respectively. In the GP formulation, it is assumed that  
 295 the  $\mathbf{f} = f_{1:N}$  are jointly Gaussian and  $\mathbf{y} = y_{1:N}$  are normally distributed given  
 296  $\mathbf{f}$ , then the prior distribution induced by the GP can be described as

$$\mathbf{f}|\mathbf{X} \sim \mathcal{N}(\mathbf{m}, \mathbf{K}), \quad \mathbf{y}|\mathbf{f}, \sigma^2 \sim \mathcal{N}(\mathbf{f}, \sigma^2 \mathbf{I}), \quad (7)$$

297 where the elements of mean vector and covariance matrix are described by  
 298  $m_i := \mu_0(\mathbf{x}_i)$  and  $K_{i,j} := k(\mathbf{x}_i, \mathbf{x}_j)$ .

299 Equation 7 describes the prior distribution induced by the GP, where  $\mathbf{X}$  is  
 300 the sampling location, and  $f$  is the objective function. In the GP formulation,  $y$   
 301 is the noise-corrupted stochastic output of  $f(\mathbf{x})$  with the variance of  $\sigma^2$ , at the  
 302 sampling location  $\mathbf{X}$ . The objective function  $f$  is assumed to be a multivariate  
 303 normal distribution function with mean  $\mathbf{m}(x)$  and covariance  $\mathbf{K}(x)$ .

304 Let  $N$  be the number of sampling locations, and  $\mathcal{D}_N = \{\mathbf{x}_i, y_i\}_{i=1}^N$  be the  
 305 set of observations. The covariance kernel  $k$  is a choice of modeling the cor-  
 306 relation between input locations  $\mathbf{x}_i$ . Covariance functions where length-scale  
 307 parameters can be inferred through maximum likelihood function is known  
 308 as automatic relevance determination kernels. One of the most widely used  
 309 kernels in this kernel family is the squared-exponential kernel,

$$\mathbf{K}_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j) = \theta_0^2 \exp\left(-\frac{r^2}{2}\right), \quad (8)$$

310 where  $r^2 = (\mathbf{x} - \mathbf{x}')^T \mathbf{\Gamma} (\mathbf{x} - \mathbf{x}')$ ,  $\mathbf{\Gamma}$  is a diagonal matrix of  $(m - n) \times (m - n)$ ,  
 311 and  $\theta_i$  is the length scale parameter.

312 The posterior Gaussian for the sequential BO is characterized by the mean

$$\mu_{N+1}(\mathbf{x}) = \mu_0(\mathbf{x}) + \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} (\mathbf{y} - \mathbf{m}), \quad (9)$$

313 and the variance

$$\sigma_{N+1}^2(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}(\mathbf{x}), \quad (10)$$

314 where  $\mathbf{k}(\mathbf{x})$  is the vector of covariance terms between  $\mathbf{x}$  and  $\mathbf{x}_{1:N}$ .

### 315 3.2 Clustering and enumeration algorithm

316 Assuming that the discrete variables are independent of each other, a clustering  
 317 and enumeration algorithm is devised to automatically decompose the large  
 318 dataset to smaller clusters based on the discrete tuple and tag a cluster with a  
 319 unique index from the enumeration scheme. For the case when some discrete  
 320 variables are dependent on others, the neighborhood can be manually changed

321 to reflect the knowledge. The set of discrete variables for each cluster are  
 322 represented as a discrete tuple where each element is a positive integer.

323 For an integer variable where order relation exists, the discrete variable can  
 324 simply be represented as a positive integer, e.g.  $1 \leq 2$ . For a categorical vari-  
 325 able where order relation does not exist, such as type of cross section (square  
 326 or circular), colors (red or blue), type of materials (aluminum or copper), con-  
 327 figuration settings, positive integers can still be used. The choice of using tuple  
 328 of positive integers as a general representation does not affect the clustering  
 329 and enumeration scheme, but would affect the construction of neighborhood  
 330 for each cluster, depending on the nature of discrete variables.

331 Suppose that the input  $\mathbf{x} = (\mathbf{x}^{(d)}, \mathbf{x}^{(c)}) = (x_1, \dots, x_n, x_{n+1}, \dots, x_m)$  in-  
 332 cludes  $n$  discrete and  $m - n$  continuous variables. If  $p_i$  is denoted as the total  
 333 number of possible values for discrete variable  $x_i, 1 \leq i \leq n$ , then the num-  
 334 ber of clusters is  $L = \prod_{i=1}^n p_i$ . Due to the complexity of possible combinations,  
 335 each cluster is assigned a unique index in such a way that the map between  
 336 their discrete variables and cluster index is one-to-one. The index is calculated  
 337 based on the total ordering of tuples. Without loss of generality, assume that  
 338 each discrete variable  $x_i$  is bounded by  $1 \leq x_i \leq p_i$ , i.e.  $x_i \in \{1, \dots, p_i\}$  for  
 339  $1 \leq i \leq n$ . Then the relation of lexicographical order, denoted as  $\prec$ , can be  
 340 defined for a pair of tuples on the set of all tuples as

$$(a_1, \dots, a_n) \prec (b_1, \dots, b_n), \quad (11)$$

341 if and only if  $\exists k : 1 \leq k \leq n : (\forall j : 1 \leq j < k : a_j = b_j)$  and  $a_k < b_k$ , and  
 342  $1 \leq a_i, b_i \leq p_i$  for all  $i$ . With the definition of lexicographical order  $\prec$ , the  
 343 cluster index  $l$  for the tuple  $(a_1, \dots, a_n)$  can now be calculated as

$$l = \sum_{i=1}^{n-1} (a_i - 1) \prod_{j=i+1}^n p_j + a_n. \quad (12)$$

344 Because the index of cluster is uniquely defined based on the tuple of discrete  
 345 variables, the tuple describing the set of discrete variables can be reconstructed  
 346 using the index of the cluster, with the quotient and remainder algorithm  
 347 recursively shown in Algorithm 1. It describes how to construct the set of  
 348 discrete variables from the cluster index  $l$ .

349 The implementation of Algorithm 1 can be based on existing functions such  
 350 as MATLAB function `ind2sub()`. Equation 12, which is a reverse operation of  
 351 Algorithm 1, can also be implemented using MATLAB function `sub2ind()`.

### 352 3.3 Construction of neighborhood

353 Consider a cluster with index  $l$ , with the tuple of discrete variables  $(a_1, \dots, a_n)$ ,  
 354 the neighbors of the  $l$ -th cluster  $\mathcal{B}(l)$  is the collection of clusters that share  
 355 most of similarity with the original cluster. Intuitively, the neighborhood is

---

**Algorithm 1** Reconstruct the tuple of discrete variables  $(x_1, \dots, x_n)$  from cluster index  $l$ .

---

**Input:** cluster index  $l$ , tuple  $(p_1, \dots, p_n)$ .

**Output:** tuple  $(a_1, \dots, a_n)$  of discrete variables

```

1: for  $i \leftarrow 1, n$  do
2:   if  $i \neq n$  then
3:      $q \prod_{j=i+1}^n p_j + r = l$                                  $\triangleright$  find quotient  $q$ , remainder  $r$ 
4:      $l \leftarrow r$ 
5:      $a_i \leftarrow q + 1$                                  $\triangleright$  assign discrete variable in order
6:   else
7:      $a_n \leftarrow r$                                  $\triangleright$  assign last discrete variable
8:   end if
9: end for
10: for  $i \leftarrow n, -1, 1$  do                                 $\triangleright$  exception if  $a_i = 0$ 
11:   if  $a_i = 0$  then
12:      $a_i \leftarrow p_i, a_{i-1} \leftarrow a_{i-1} - 1$ 
13:   end if
14: end for

```

---

356 constructed based on the belief of whether there exists a relationship between  
357 two clusters.

358 For example, for integer variables, the discrete tuples of the neighboring  
359 clusters may differ in one or a few different integer variables compared to  
360 that of the original cluster. In the same manner, for categorical variables, the  
361 discrete tuples of the neighboring clusters may differ in one or a few categorical  
362 variables compared to that of the original cluster. Based on this description,  
363 a possible choice to define the neighborhood  $\mathcal{B}(l)$  of the  $l$ -th cluster can be  
364 mathematically expressed as

$$\mathcal{B}(l) = \{(a_1^*, \dots, a_n^*) \mid d((a_1^*, \dots, a_n^*), (a_1, \dots, a_n)) \leq d_{\text{th}}\}, \quad (13)$$

365 where  $d((a_{i=1}^n, (a^*)_{i=1}^n)$  is some metric on a discrete topological space  $\mathbb{D}$ ,  
366 and  $d_{\text{th}}$  is a user-defined threshold. The metric  $d(\cdot, \cdot)$  can be any  $l_p$ -norm, for  
367 example, Manhattan distance ( $l_1$ -norm), or a counting metric of how many  
368 discrete (integer and categorical) variables are different between two tuples.  
369 It is noted that the metric  $d(\cdot, \cdot)$  does not have to strictly obey the definition  
370 of mathematical norm. In the special case when this metric is set to zero, i.e.  
371  $d((a_{i=1}^n, (a^*)_{i=1}^n) = 0$ , it means that all the clusters are considered to be  
372 completely independent of each other. The construction of neighborhood only  
373 occurs once during the initialization.

374 Furthermore, it should be emphasized that the neighboring list can be  
375 manually changed to reflect the physics-based knowledge from the users, or  
376 manually constructed to reflect the dependency of the discrete variables. In the  
377 case of categorical variables where independence is usually observed, one can  
378 simply remove the neighboring cluster from the corresponding categorical vari-

379 able, as the neighborhood can be manually changed during the initialization  
 380 phase of the optimization process.

381 It is recommended to define the neighborhood carefully, as the neighbor-  
 382 hood definition has an impact on both convergence rate, and whether the  
 383 optimization would be trapped at local optimum. The safest setting is to as-  
 384 sign  $d_{th} = 0$ , where clusters are assumed to be completely independent of  
 385 each other. Small values of  $d_{th}$ , e.g.  $d_{th} = 1$  or  $d_{th} = 2$ , might be beneficial,  
 386 depending on the specific applications. Large value is not recommended.

387 Figure 2 shows an example of constructing clusters for two discrete vari-  
 388 ables  $(x_1, x_2)$ , where  $1 \leq x_1 \leq 4$  and  $1 \leq x_2 \leq 3$ . According to Algorithm 1,  
 389 the tuple  $p$  is  $(4, 3)$ , cluster 1 is associated with  $(1,1)$ , cluster 2 is associated  
 390 with  $(1,2)$ , cluster 4 is associated with  $(2,1)$ , etc. The cluster index is denoted  
 391 as an italic number on the top right corner of the square. Consider cluster 8,  
 392 which is associated with the discrete tuple  $(3,2)$ . If the Manhattan distance  
 393 is chosen to define the neighborhood, then the choice of  $d_{th} = 0$  in Equation  
 394 13 would make every cluster the only neighbor of itself, e.g. the neighbor of  
 395 cluster 8 is cluster 8. The choice of  $d_{th} = 1$  would include clusters 5, 7, 8, 9,  
 396 11 in cluster 8's neighborhood. Similarly, the choice of  $d_{th} = 2$  would include  
 clusters 2, 4, 5, 6, 7, 8, 9, 10, 11, 12 in cluster 8's neighborhood.

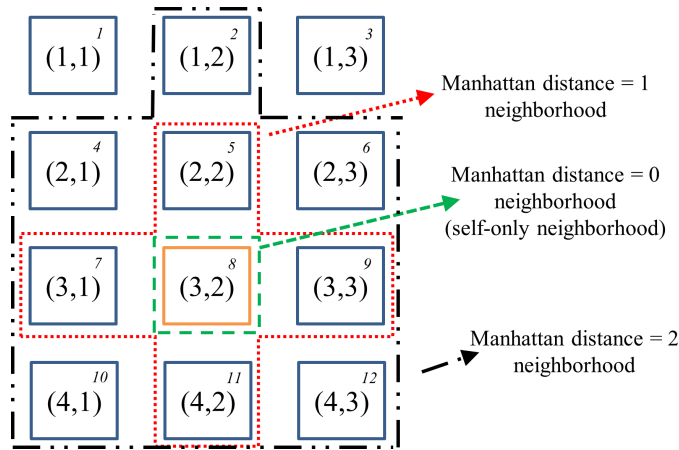


Fig. 2: An example of cluster enumeration and neighborhood definition.

397

### 398 3.4 Weight computation

399 The weight of each cluster's prediction is determined by the Wasserstein dis-  
 400 tance between the Gaussian posterior of the main cluster with that of the  
 401 neighboring clusters. Combined together, they form a Gaussian mixture model  
 402 to predict a response at a query point  $\mathbf{x}$ .

403 Consider a query point  $\mathbf{x}$  in the  $l$ -th cluster, which has the continuous  
 404 tuple  $\mathbf{x}^{(c)} = (x_{n+1}, \dots, x_m)$ . Denote the neighborhood of the  $l$ -th cluster as  
 405  $B(l) = \{l^*\}$ , where the cardinality of  $|B(l)| = k$ , i.e. there are  $k$  neighbors  
 406 in the  $l$ -th cluster neighborhood. Each of the neighboring cluster  $l^*$  can form  
 407 its own prediction  $\mathcal{N}(\mu_{l^*}, \sigma_{l^*}^2)$  from the continuous tuple, including  $\mathcal{N}(\mu_l, \sigma_l^2)$   
 408 for  $l$ -th cluster. However, the prediction must be adjusted by accounting for  
 409 the bias, i.e.  $\text{Bias}_{l^*}[\mu_{l^*}] = \mathbb{E}[\mu_{l^*} - \mu_l] = \bar{\mu}_{l^*} - \bar{\mu}_l$  as the difference between the  
 410 posterior means of two clusters, and the variance  $\sigma_{l^*}^2$ .

411 The weight  $w_{l^*}$  associated with the prediction from the  $l^*$  cluster should  
 412 be larger with smaller bias ( $\bar{\mu}_{l^*} - \bar{\mu}_l$ ) and smaller posterior variance  $\sigma_{l^*}^2$ . The  
 413 necessity of bias correction is explained later in Theorem 4. Wasserstein distance  
 414 between two univariate Gaussian  $\mathcal{N}(\mu_{l^*}, \sigma_{l^*}^2)$  and  $\mathcal{N}(\mu_l, \sigma_l^2)$  is provided  
 415 by Givens et al. [13] as

$$W_2(\mathcal{N}(\mu_{l^*}, \sigma_{l^*}^2), \mathcal{N}(\mu_l, \sigma_l^2)) = \|\mu_l - \mu_{l^*}\|^2 + \left\| \sqrt{\sigma_l^2} - \sqrt{\sigma_{l^*}^2} \right\|^2 \quad (14)$$

416 Here we propose a deterministic way to compute the numerical weights  
 417 based on the pair-wise Wasserstein distance, which eventually converges to an  
 418 independent GP as the optimization process advances. It is easy to see that  
 419 the  $W_2$ -distance of the  $l$ -th cluster's prediction to itself is zero, as  $W_2$  is a  
 420 distance. The weights are computed according to an inverse  $W_2$ -distance with  
 421 a term  $\sigma_l^2$  from the  $l$ -th cluster, as

$$w_{l^*} \propto [\sigma_l^2 + W_2(\mathcal{N}(\mu_{l^*}, \sigma_{l^*}^2), \mathcal{N}(\mu_l, \sigma_l^2))]^{-1}. \quad (15)$$

422 In Equation 15,  $w_{l^*}$  are computed based on two factors, the  $W_2$ -distance, and  
 423 the  $\sigma_l^2$  prediction of the  $l$ -th cluster. As the optimization process advances, the  
 424 posterior variance approaches zero, i.e.  $\sigma_l^2 \rightarrow 0$ . As a result, the weight scheme  
 425 converges to a single GP prediction of the corresponding  $l$ -th cluster.

### 426 3.5 Prediction using weighted average of $k$ -nearest neighboring clusters

427 We model the prediction of a query point using a Gaussian mixture distribu-  
 428 tion, where the weights are computed on the statistical Wasserstein distance.  
 429 To predict an unknown query point  $\mathbf{x} = (\mathbf{x}_d, \mathbf{x}_c) = (x_1, \dots, x_n, x_{n+1}, \dots, x_m)$ ,  
 430 we first find the cluster in which  $\mathbf{x}$  belongs to, and its neighboring clusters.  
 431 Assume that  $\mathbf{x}$  belongs to the  $l$ -th cluster, and there are  $k$ -neighboring clusters.

432 The principle for weight computation is as follows. As the bias increases,  
 433 the contributed weight of the prediction  $w_{l^*}$  from the  $l^*$ -th cluster to  $l$ -th  
 434 cluster is reduced to a smaller value. As the bias or the pair-wise distance  
 435 between clusters increases, the contributed weights also decrease. The weight  
 436 vector is normalized at every step, and eventually converges to a single GP  
 437 prediction with the weight vector of  $[0, \dots, 1, \dots, 0]$ , where 1 is located as the  
 438  $l$ -th cluster.

439 Since  $\mathbf{x}$  is located within the  $l$ -th cluster, the weight from the  $l$ -th cluster  
 440 is the highest, i.e. if  $l^* = l$ , then  $\mu_{l^*} + \bar{\mu}_l - \bar{\mu}_{l^*} = \mu_l$ , which is the GP prediction  
 441 for the  $l$ -th cluster. The posterior mean of the proposed method is written as

$$\hat{\mu} = \sum_{l^* \in \mathcal{B}(l)} w_{l^*} (\mu_{l^*} + \bar{\mu}_l - \bar{\mu}_{l^*}), \quad (16)$$

442 where the sum is taken over the list of neighboring cluster from the main cluster  
 443  $l^{\text{th}}$ .  $\bar{\mu}_l$  and  $\bar{\mu}_{l^*}$  denote the means of the  $l$ -th and  $l^*$ -th clusters, respectively.  
 444  $w^*$  denotes the weight corresponding to the  $l^*$ -th cluster, which is computed  
 445 once the discrete tuple  $\mathbf{x}^{(d)}$  of the query point  $\mathbf{x} = (\mathbf{x}^{(d)}, \mathbf{x}^{(c)})$  is determined.  
 446 The posterior variance of the proposed method is calculated as

$$\hat{\sigma}^2 = \sum_{l^* \in \mathcal{B}(l)} w_{l^*}^2 \sigma_{l^*}^2, \quad (17)$$

447 where  $\sigma_{l^*}^2$  denotes the posterior variance associated with the continuous tuple  
 448  $\mathbf{x}^{(c)}$  of the query point  $\mathbf{x} = (\mathbf{x}^{(d)}, \mathbf{x}^{(c)})$ .

449 The prediction scheme for mean  $\hat{\mu}(\mathbf{x})$  and variance  $\hat{\sigma}^2(\mathbf{x})$  for an arbitrary  
 450 location  $\mathbf{x}$  using Gaussian mixture model can be summarized in Algorithm 2.

---

**Algorithm 2** Prediction using weighted average GP from nearest neighboring clusters.

---

**Input:** location  $\mathbf{x} = (x_1, \dots, x_n, x_{n+1}, \dots, x_m)$ , mean output of each cluster  $\bar{\mu}_{(\cdot)}$

**Output:** Gaussian mixture posterior mean  $\hat{\mu}$  and posterior variance  $\sigma^2$

- 1: Find cluster index  $l$  corresponding to  $\mathbf{x}^{(d)} = (x_1, \dots, x_n)$  ▷ locate the  $l$ -th cluster
  - 2: Construct a neighborhood  $\mathcal{B}(\cdot)$  for each cluster ▷ query  $\mathbf{x}$  in all neighboring clusters
  - 3: **for**  $l^* \in \mathcal{B}(l)$  **do**
  - 4:   Compute GP posterior of the  $l^*$ -th cluster:  $\hat{\mu}_{l^*}, \sigma_{l^*}^2$
  - 5: **end for**
  - 6: Compute weight  $w_{l^*} \propto [\sigma_{l^*}^2 + W_2(\mathcal{N}(\mu_{l^*}, \sigma_{l^*}^2), \mathcal{N}(\mu_l, \sigma_l^2))]^{-1}$  ▷ pair-wise Wasserstein distance
  - 7:  $w_{l^*} \leftarrow \frac{w_{l^*}}{\sum_{l^* \in \mathcal{B}(l)} w_{l^*}}$  ▷ weight normalization
  - 8:  $\hat{\mu} \leftarrow \sum_{l^* \in \mathcal{B}(l)} w_{l^*} (\mu_{l^*} + \bar{\mu}_l - \bar{\mu}_{l^*})$  ▷ Gaussian mixture posterior mean
  - 9:  $\hat{\sigma}^2 \leftarrow \sum_{l^* \in \mathcal{B}(l)} w_{l^*}^2 \sigma_{l^*}^2$  ▷ Gaussian mixture posterior variance
  - 10: Update the average mean of the  $l$ -th cluster  $\bar{\mu}_l$
- 

### 451 3.6 Constrained acquisition function in mixed-integer Bayesian optimization

452 The acquisition function is adopted from Gardner et al. [11] for inequality  
 453 constraints, and further extended to accommodate discrete and continuous  
 454 variables to solve for the constrained mixed-integer optimization problems.

455 First, the constraint is checked using an indicator function  $\mathcal{I}(\mathbf{x})$  for all  $i_c$   
 456 constrained inequalities, as

$$\mathcal{I}(\mathbf{x}) = \begin{cases} 1 & \text{if } \forall 1 \leq i \leq i_c : g_i(\mathbf{x}) \leq 0, \\ 0 & \text{if } \exists 1 \leq i \leq i_c : 0 \leq g_i(\mathbf{x}). \end{cases} \quad (18)$$

457 The constrained acquisition function can be considered as the product of the  
 458 classical acquisition function. As a result, the acquisition function is assigned  
 459 to have zero value for infeasible region. The penalized approach can be imple-  
 460 mented directly into the auxiliary optimizer, which is used to maximize the  
 461 acquisition function in BO.

462 In distributed GP, an input  $\mathbf{x}_{\text{next}} = (x_1, \dots, x_n, x_{n+1}, \dots, x_m)$  is com-  
 463 prised of both discrete and continuous variables. For each cluster correspond-  
 464 ing to a unique set of discrete tuple  $(x_1, \dots, x_n)$ , a distinct next sampling point  
 465 associated with each cluster is located by maximizing the acquisition function  
 466 on the tuple of continuous variables  $(x_{n+1}, \dots, x_m)$  for each iteration, in the  
 467 same manner as classical BO. These next sampling points are retained within  
 468 the respective clusters. However, only the sampling point corresponding to the  
 469 maximal value of acquisition function among all clusters is chosen, and a new  
 470 sampling point within that cluster is located and updated for the correspond-  
 471 ing cluster. The sampling procedure repeats until the optimization criterion is  
 472 met. In other words, the next sampling point is chosen as

$$\mathbf{x}_{\text{next}} = \arg \max_{l^*} \arg \max_{(x_n, x_{n+1}, \dots, x_m)} a_{l^*}(\mathbf{x}; \{\mathbf{x}_i, y_i\}_{i=1}^N, \theta) \cdot \mathcal{I}(\mathbf{x}), \quad (19)$$

473 where the  $l^*$ -th cluster corresponds to the tuple of discrete variables  $(x_1, \dots, x_n)$ ,  
 474 and  $\mathcal{I}(\mathbf{x})$  is the constraint indicator function.

475 Equation 19, which describes the searching procedure for the next sam-  
 476 pling point by maximizing the penalized acquisition function, is explained as  
 477 follows. Two loops are constructed to search for the global sampling point. In  
 478 the inner loop which searches for the local sampling point within each cluster,  
 479 the penalized acquisition function is the objective function. Maximizing  
 480 this penalized acquisition function using an auxiliary optimizer yields the lo-  
 481 cal sampling point for each cluster. In the outer loop, the cluster with the  
 482 maximized acquisition function value is determined. The discrete tuple cor-  
 483 responding to the cluster index, which contains the sampling point with the  
 484 maximum value for the acquisition function, is reconstructed using Algorithm  
 485 1. In other words, the sampling location  $\mathbf{x}$  is decomposed to two parts: the  
 486 inner loop searches for the continuous tuple, whereas the outer loop yields the  
 487 discrete tuple. Theoretically, once the functional evaluation is over, only the  
 488 cluster that contains the last sampling location needs to be updated. Practi-  
 489 cally, all the clusters need to update their corresponding sampling locations  
 490  $\mathbf{x}_{\text{next}}$  after certain number of iterations, in order to avoid trapping in local  
 491 optimum.

492 The tuple of continuous variables is found by maximizing the acquisition  
 493 function, whereas the tuple of discrete variables is assigned according to the

494 cluster index. For the EI and PI acquisition functions,  $x_{\text{best}}$  is modified to  
 495 be the best point achieved so far among all clusters. For the UCB acquisi-  
 496 tion function, no modification is needed, assuming the hyperparameter  $\kappa$  is  
 497 uniform for all clusters. It is noted that the balance between exploration and  
 498 exploitation is preserved locally within each cluster, and thus is also preserved  
 499 globally for all the clusters.

### 500 3.7 Theoretical bounds and computational cost

501 Here we provide the theoretical lower and upper bounds for predictions and  
 502 algorithm complexity under the formulation of Gaussian mixture model in  
 503 Theorem 1 and Theorem 2. Theorem 3 proves that under the formulation of  
 504 the proposed method, the largest weight is associated with the main cluster.  
 505 Theorem 4 explains the necessity of translation in mean prediction so that the  
 506 expected value of the mean is the same with the expected mean in the main  
 507 cluster.

508 **Theorem 1** *The Gaussian mixture posterior mean  $\hat{\mu} = \sum_{l^* \in \mathcal{B}(l)} w_{l^*} (\mu_{l^*} + \bar{\mu}_l - \bar{\mu}_{l^*})$*

509 *is bounded by*

$$\min_{l^*} (\hat{\mu}_{l^*} + \bar{\mu}_l - \bar{\mu}_{l^*}) \leq \hat{\mu} \leq \max_{l^*} (\hat{\mu}_{l^*} + \bar{\mu}_l - \bar{\mu}_{l^*}) \quad (20)$$

510 *Proof* The proof for the posterior mean is straightforward, noting that  $w_{l^*} \geq$   
 511  $0, \forall l^*$  and  $\sum w_{l^*} = 1$ .

512 **Theorem 2** *The Gaussian mixture posterior variance  $\hat{\sigma}^2 = \sum_{l^* \in \mathcal{B}(l)} w_{l^*}^2 \sigma_{l^*}^2$  is*

513 *bounded by*

$$\left( \sum_{l^*} w_{l^*}^2 \sigma_{l^*} \right)^2 \leq \hat{\sigma}^2 \leq \max_{l^*} \sigma_{l^*}^2 \quad (21)$$

514 *Proof* For the right-hand side of the variance inequality, observe that

$$\begin{aligned} \hat{\sigma}^2 &= \sum_{l^*} w_{l^*}^2 \sigma_{l^*}^2 \leq \sum_{l^*} w_{l^*} \sigma_{l^*}^2 \quad (\text{because } w_{l^*}^2 \leq w_{l^*}) \\ &\leq \left( \sum_{l^*} w_{l^*} \right) \max_{l^*} \sigma_{l^*}^2 \leq \max_{l^*} \sigma_{l^*}^2 \quad (\text{because } \sum_{l^*} w_{l^*} = 1) \end{aligned} \quad (22)$$

515 For the left-hand side of the variance inequality, recall the Jensen's inequality:

516  $\rho \left( \frac{\sum_i a_i x_i}{\sum_i a_i} \right) \leq \frac{\sum_i a_i \rho(x_i)}{\sum_i a_i}$ , where  $\rho(\cdot)$  is a convex function. Substitute  $w_{l^*}^2 \rightarrow$

517  $a_i, \sigma_{l^*} \rightarrow x_i$ , and  $\rho(x) = x^2$  into the Jensen's inequality, we have

$$\left( \frac{\sum_{l^*} w_{l^*}^2 \sigma_{l^*}}{\sum_{l^*} w_{l^*}^2} \right)^2 \leq \frac{\sum_{l^*} w_{l^*}^2 \sigma_{l^*}^2}{\sum_{l^*} w_{l^*}^2} \quad \text{or} \quad \left( \sum_{l^*} w_{l^*}^2 \sigma_{l^*} \right)^2 \leq \left( \sum_{l^*} w_{l^*}^2 \right) \left( \sum_{l^*} w_{l^*}^2 \sigma_{l^*}^2 \right) \quad (23)$$



518 Now, note that  $\sum_{l^*} w_{l^*}^2 \leq \sum_{l^*} w_{l^*} = 1$ . We obtain the left-hand side of the  
 519 inequality.

520 **Theorem 3** *The largest weight is associated with the  $l$ -th cluster.*

521 *Proof* Based on the weight formula,

$$w_{l^*} \propto [\sigma_{l^*}^2 + W_2(\mathcal{N}(\mu_{l^*}, \sigma_{l^*}^2), \mathcal{N}(\mu_l, \sigma_l^2))]^{-1}, \quad (24)$$

522 it is easy to see that the Wasserstein distance between a cluster with itself is  
 523 zero. Thus, the right-hand side is always less than  $\sigma_{l^*}^2$ , i.e.

$$\sigma_{l^*}^2 + W_2(\mathcal{N}(\mu_{l^*}, \sigma_{l^*}^2), \mathcal{N}(\mu_l, \sigma_l^2)) \geq \sigma_{l^*}^2. \quad (25)$$

524 Inversing the last inequality completes the proof. The equality occurs when  
 525  $l^* = l$ .

526 **Theorem 4** *The expectation of the posterior mean  $\hat{\mu} = \sum_{l^* \in \mathcal{B}(l)} w_{l^*} (\mu_{l^*} + \bar{\mu}_l - \bar{\mu}_{l^*})$*

527 *is  $\bar{\mu}_l$ , i.e.  $\mathbb{E}[\hat{\mu}] = \bar{\mu}_l$ .*

528 *Proof* Take the expectation of Equation 9 for any  $l$ -th cluster over the con-  
 529 tinuous domain, and note that  $\mathbb{E}[\mathbf{y} - \mathbf{m}] = 0$ , the mean of the posterior is  
 530 recovered to the mean of the cluster, i.e.

$$\mathbb{E}[\mu_l(\mathbf{x})] = \mu_0(\mathbf{x}) = \bar{\mu}_l(\mathbf{x}). \quad (26)$$

531 Equation 26 holds for any  $l$ -th under the GP formulation. In the similar man-  
 532 ner, taking the expectation of the posterior mean  $\hat{\mu}$  from the proposed method  
 533 over the continuous domain, we arrive at

$$\begin{aligned} \mathbb{E}[\hat{\mu}] &= \sum_{l^* \in \mathcal{B}(l)} w_{l^*} \mathbb{E}[\mu_{l^*} + \bar{\mu}_l - \bar{\mu}_{l^*}] \\ &= \sum_{l^* \in \mathcal{B}(l)} w_{l^*} [\mathbb{E}[\mu_{l^*}] + \mathbb{E}[\bar{\mu}_l] - \mathbb{E}[\bar{\mu}_{l^*}]] \\ &= \sum_{l^* \in \mathcal{B}(l)} w_{l^*} [\bar{\mu}_{l^*} + \mathbb{E}[\bar{\mu}_l] - \bar{\mu}_{l^*}] \\ &= \sum_{l^* \in \mathcal{B}(l)} w_{l^*} [\bar{\mu}_l] \\ &= \bar{\mu}_l, \end{aligned} \quad (27)$$

534 where the second equality is formed by distributing the expectation operator  
 535 under linear combination rule. The third equality follows Equation 26 as de-  
 536 scribed above. The fourth equality is formed by canceling two identical terms  
 537  $\bar{\mu}_{l^*}$ .

538 A major problem of GP is its scalability, which originates from the com-  
 539 putation of the inverse of correlation matrices. The dataset decomposition has  
 540 a favorable computational aspect in which the scalability is alleviated. Here  
 541 we analyze the computational cost based on the assumption that the size of  
 542 each cluster is roughly equal. Denote the number of data points for the whole  
 543 dataset as  $N$ , and the number of clusters as  $k$ . The computational cost to  
 544 compute all covariance matrices is reduced by a factor of  $k^2$ , as  $k$  covariance  
 545 matrices are involved, and each covariance matrix has the computational com-  
 546 plexity  $\mathcal{O}\left(\frac{N}{k}\right)^3$ , thus resulting in the total cost of  $k\mathcal{O}\left(\frac{N}{k}\right)^3 = \frac{1}{k^2}\mathcal{O}(N^3)$ .  
 547 Similarly, the cost of storing covariance matrices is also reduced by a factor  
 548 of  $k$ , since  $k\mathcal{O}\left(\frac{N}{k}\right)^2 = \frac{1}{k}\mathcal{O}(N^2)$ . However, the computational cost of pre-  
 549 dicting the posterior mean  $\mu$  and posterior variance  $\sigma^2$  stays the same, since  
 550  $k\mathcal{O}\left(\frac{N}{k}\right) = \mathcal{O}(N)$ . The decomposition approach in the proposed mixed-  
 551 integer BO has a computational advantage to mitigate the scalability problem  
 552 in GP, even though it is not completely eliminated.

#### 553 4 Analytical examples

554 In this section, the proposed mixed-integer BO is compared with the genetic  
 555 algorithm (GA) with various settings. The settings for the GA are described as  
 556 follows. To verify the robustness of the proposed method, three GA settings are  
 557 chosen. In the first setting, the population size and the elite count parameters  
 558 are set to be 50 and 3, respectively. In the second setting, the population size  
 559 and the elite count parameters are set to be 150 and 10, respectively. In the  
 560 third setting, they are 1500 and 10, respectively. Other parameters are left to  
 561 be the default values in MATLAB function `ga()`.

562 In Section 4.1, a discrete modification of multi-modal Rastrigin function is  
 563 used as a benchmark function, where two variables are discrete and the other  
 564 two are continuous. In Section 4.2, a welded beam design optimization with  
 565 two discrete and four continuous variables is used to evaluate the performance  
 566 of the proposed mixed-integer BO method where discrete variables come from  
 567 the configuration and material of the beam. In Section 4.3, a pressure vessel  
 568 design optimization with four continuous variables is benchmarked. In Section  
 569 4.4, a speed reducer design optimization function with one discrete and six  
 570 continuous variables is utilized. In Section 4.5, a modification of discrete sphere  
 571 function is devised to demonstrate the proposed mixed-integer BO method  
 572 on high-dimensional optimization problems, with 5 discrete and 50 and 100  
 573 continuous variables.

## 574 4.1 Discrete Rastrigin function

575 In this example, the proposed method is applied on the discrete version of  
 576 the Rastrigin function, which is an analytical function for testing different  
 577 optimization methods. To evaluate the effectiveness of the proposed mixed-  
 578 integer BO method, the optimization performance is compared against GA  
 579 optimization performance.

### 580 4.1.1 Problem statement

581 The DACE toolbox [41] for classical GP is extended to include the pro-  
 582 posed distributed GP and Bayesian optimization. In this section, the hy-  
 583 brid Bayesian optimization is to find the global minimum on a tiled ver-  
 584 sion of Rastrigin function on 25 clusters, where each cluster corresponds to  
 585 two discrete variables. The input  $\mathbf{x} = (i, j, x, y)$  is comprised of four vari-  
 586 ables, in which the first two are discrete, and the last two are continuous,  
 587 as illustrated in Figure 3. The original two-dimensional Rastrigin function is  
 588  $f(x, y) = 20 + [x^2 - 10 \cos(2\pi x) + y^2 - 10 \cos(2\pi y)]$ , where  $-5.12 \leq x, y \leq 5.12$ .  
 589 The tiled Rastrigin function is constructed based on a tiled domain of Rastri-  
 590 grin function, where each domain is characterized by a discrete tuple  $(i, j)$ ,  
 591 and the continuous domain is translated to  $-0.75 \leq x_{\text{tiled}}, y_{\text{tiled}} \leq 0.75$  for all  
 592 clusters. Figure 3 illustrates the construction of tiled Rastrigin function, and  
 593 its relationship with the original Rastrigin function. The relationship between  
 594 the tiled and original Rastrigin can simply be described by an affine function,

$$x_{\text{orig}} = -3.50 + 1.75(i - 1) + x_{\text{tiled}}; \quad y_{\text{orig}} = -3.50 + 1.75(j - 1) + y_{\text{tiled}}, \quad (28)$$

595 where  $-0.75 \leq x_{\text{tiled}}, y_{\text{tiled}} \leq 0.75$ .

### 596 4.1.2 Numerical results

597 In this example, to find the minimum of Rastrigin function, we flip the sign of  
 598 tiled Rastrigin and use the UCB acquisition function to locate the maximum  
 599 of the negative tiled Rastrigin function. The covariance matrix adaptation evo-  
 600 lution strategy (CMA-ES) [17] method is employed to find the next sampling  
 601 point within each cluster by locating the point with the maximum acquisi-  
 602 tion function. The parameters are set as follows:  $\kappa = 5$ ,  $d_{\text{penalty}} = 10^{-4}$ ,  
 603  $N_{\text{shuffle}} = 15$ , where  $N_{\text{shuffle}}$  is the number of steps which CMA-ES is reactiv-  
 604 ated with different initial position to search for the next sampling point on  
 605 each local GP in order to avoid trapping in the local minima. To construct  
 606 the initial GP response surface, 5 random data points are sampled from each  
 607 cluster.

608 Because the global minimum of the original Rastrigin function is at  $(x =$   
 609  $0, y = 0)$  with the functional evaluation  $f(0, 0) = 0$ , the hybrid Bayesian  
 610 optimizer on the tiled Rastrigin function is expected to converge to cluster 13,  
 611 as illustrated in Figure 3. The neighbor list of cluster 13 includes clusters 8,

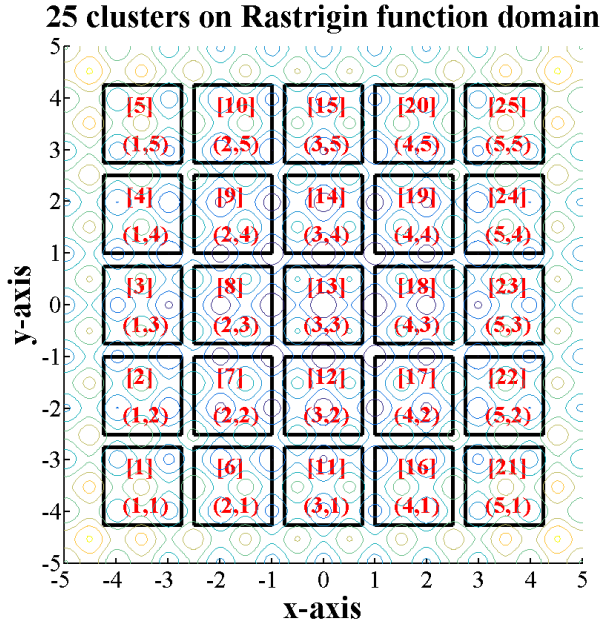


Fig. 3: Tiled Rastrigin function comprising of 25 clusters, where each cluster correspond to a square of dimension  $1.50 \times 1.50$  and a tuple  $(i, j)$ . The cluster index is denoted within the square bracket  $[\cdot]$ , whereas the tuple is within the parenthesis  $(\cdot, \cdot)$  in each square.

612 12, 13, 14, and 18. Figure 4 compares the numerical performance between the  
 613 proposed mixed integer BO and the GA with three different settings.

614 Figure 4 presents the performance of the proposed method (solid line)  
 615 with five different settings, and the GA method (dash line) with three dif-  
 616 ferent settings. For the proposed mixed integer BO, the threshold distance  
 617  $d_{th}$  is changed. The proposed mixed integer BO performs best with small  $d_{th}$   
 618 parameter, which measures the dissimilarity between discrete tuples.

#### 619 4.2 Welded beam design problem

620 To verify the result of the proposed method, an analytical engineering model  
 621 for welded beam design is adapted from Deb and Goyal [8], Gandomi and  
 622 Yang [10], Rao [44], Datta and Figueira [6], as shown in Figure 5, with some  
 623 slight modifications.

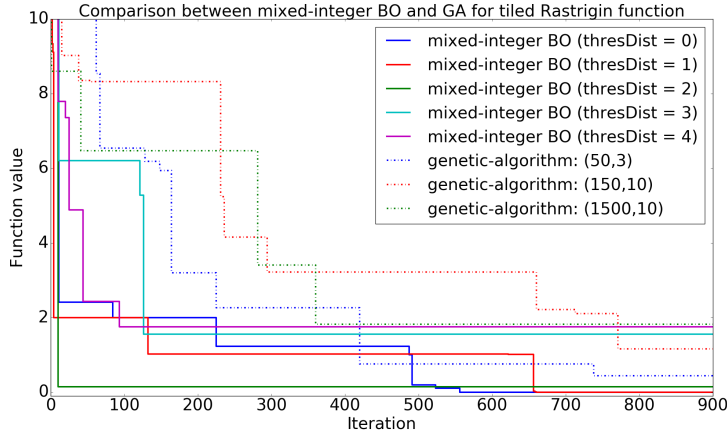


Fig. 4: Performance comparison between the GA and the proposed mixed-integer BO for the tiled Rastrigin function.

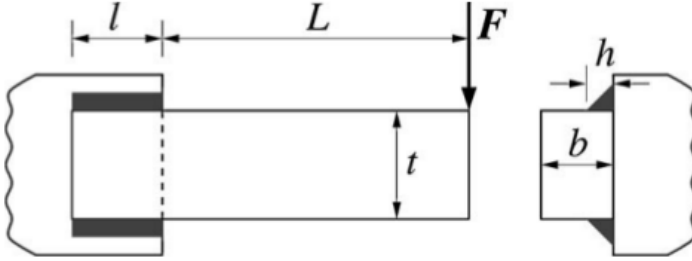


Fig. 5: Welded beam design problem [6].

#### 624 4.2.1 Problem statement

625 The low-carbon steel (C-1010) beam is welded to a rigid base to support a designated  
 626 load  $F$ . The thickness of the weld  $h$ , the length of the welded joint  $l$ , the  
 627 width of the beam  $t$  and the thickness of the beam  $b$  are the design continuous  
 628 variables. Two different welding configurations can be used, four-sided welding  
 629 and two-side welding [8]. The bulk material of the beam can be steel, cast iron,  
 630 aluminum, or brass, which is associated with different material properties. The  
 631 stress, deflection, and buckling conditions are derived from Ravindran et al.  
 632 [45], where the constant parameters are as follows.  $L = 14\text{inch}$ ,  $\delta_{\max} = 0.25$   
 633 inch, and  $F = 6,000\text{lb}$ . The input  $\mathbf{x}$  is comprised of  $(w, m, h, l, t, b)$ , where  
 634  $w$  and  $m$  are discrete variables, and  $h, l, t, b$  are continuous variables. We  
 635 note that  $h, t, b$  are commonly considered as discrete variables in multiples of  
 636 0.0625 in, as well as continuous variables, bounded between lower and upper  
 637 bounds.

638 Under this formulation, the objective is to minimize

$$f(w, m, h, l, t, b) = (1 + C_1)(wt + l)h^2 + C_2tb(L + l) \quad (29)$$

639 subject to the five inequality constraints:

$$\text{shear stress}(\tau) : g_1 = 0.577\sigma_d - \tau(\mathbf{x}) \geq 0 \quad (30a)$$

$$\text{bending stress in the beam}(\sigma) : g_2 = \sigma_d - \sigma(\mathbf{x}) \geq 0 \quad (30b)$$

$$\text{buckling load on the bar}(P_c) : g_3 = b - h \geq 0 \quad (30c)$$

$$\text{deflection of the beam} : g_4 = P_c(\mathbf{x}) - F \geq 0 \quad (30d)$$

$$\text{side constraints} : g_5 = \delta_{\max} - \delta(\mathbf{x}) \geq 0 \quad (30e)$$

640 where

$$\sigma(\mathbf{x}) = \frac{6FL}{t^2b}, \delta(\mathbf{x}) = \frac{4FL^3}{Et^3b}, P_c(\mathbf{x}) = \frac{4.013tb^3\sqrt{EG}}{6L^2} \left(1 - \frac{t}{4L}\sqrt{\frac{E}{G}}\right) \quad (31a)$$

$$\tau = \sqrt{(\tau')^2 + (\tau'')^2 + 2\tau'\tau''\cos\theta}, \tau' = \frac{F}{A}, \tau'' = \frac{F(L + 0.5l)R}{J} \quad (31b)$$

$$w = 0 : \begin{cases} A = \sqrt{2}hl \\ J = \sqrt{2}hl \left[ \frac{(h+t)^2}{4} + \frac{l^2}{12} \right] \\ R = \frac{1}{2}\sqrt{l^2 + (h+t)^2} \\ \cos\theta = \frac{l}{2R} \end{cases}, \quad (31c)$$

$$w = 1 : \begin{cases} A = \sqrt{2}h(t+l) \\ J = \sqrt{2}hl \left[ \frac{(h+t)^2}{4} + \frac{l^2}{12} \right] + \sqrt{2}ht \left[ \frac{(h+l)^2}{4} + \frac{t^2}{12} \right] \\ R = \max \left\{ \frac{1}{2}\sqrt{l^2 + (h+t)^2}, \frac{1}{2}\sqrt{t^2 + (h+l)^2} \right\} \\ \cos\theta = \frac{l}{2R} \end{cases} \quad (31d)$$

641 where  $w$  is the binary variable to model the type of weld,  $w = 0$  is used for two-  
 642 sided welding and  $w = 1$  is used for four-sided welding.  $C_1(m)$ ,  $C_2(m)$ ,  $\sigma_d(m)$ ,  
 643  $E(m)$ ,  $G(m)$  are material-dependent parameters [8][10] listed in Table 1. The  
 644 lower and upper bounds of the problem are  $0.0625 \leq h \leq 2$ ,  $0.1 \leq l \leq 10$ ,  
 645  $2.0 \leq t \leq 20.0$ , and  $0.0625 \leq b \leq 2.0$  [6].

#### 646 4.2.2 Numerical results

647 Here, the input vector is encoded as  $\mathbf{x} = (w, m, h, l, t, b)$ , where  $w \in \{0, 1\}$ ,  
 648 where  $w = 0$  and  $w = 1$  correspond to the two-sided and four-sided welding,  
 649 respectively;  $m \in \{1, 2, 3, 4\}$  corresponds to steel, cast iron, aluminum, and  
 650 brass, respectively.

Table 1: Material-dependent parameters and constants in the welded beam design problem.

Constants	Description	steel	cast iron	aluminum	brass
$C_1$	cost per volume of the welded material (\$/in <sup>3</sup> )	0.1047	0.0489	0.5235	0.5584
$C_2$	cost per volume of the bar stock (\$/in <sup>3</sup> )	0.0481	0.0224	0.2405	0.2566
$\sigma_d$	design normal stress of the bar material (psi)	$30 \cdot 10^3$	$8 \cdot 10^3$	$5 \cdot 10^3$	$8 \cdot 10^3$
$E$	Young's modulus of bar stock (psi)	$30 \cdot 10^6$	$14 \cdot 10^6$	$10 \cdot 10^6$	$16 \cdot 10^6$
$G$	shear modulus of bar stock (psi)	$12 \cdot 10^6$	$6 \cdot 10^6$	$4 \cdot 10^6$	$6 \cdot 10^6$

651 In this example, there are 8 clusters, because there are two choices for  $w$   
652 and four choices for  $m$ . The neighborhood  $\mathcal{B}(\cdot)$  is considered as universal, i.e.  
653 the neighborhood for each cluster includes all clusters, such that they are all  
654 aware of others. The bounds for hyper-parameters  $\theta$  for the GP in each cluster  
655 are set as follows.  $\underline{\theta} = (0.1, 0.1, 0.1, 0.1)$ .  $\bar{\theta} = (20.0, 20.0, 20.0, 20.0)$ . Every four  
656 iterations, the sampling point location in each cluster is computed again to  
657 avoid trapping in local minima. CMA-ES [17] is used as an auxiliary optimizer  
658 for maximizing the acquisition function. There are two random sampling points  
659 in each cluster to initialize the GP construction. The EI acquisition function  
660 is used.

661 Figure 6 shows the convergence plot of the cost function in the welded beam  
662 design, where the circle, cross, triangle, and square corresponds to steel, cast  
663 iron, aluminum, brass, respectively. The optimal cost value  $f(\mathbf{x})$  evolves at it-  
664 erations 0, 1, 2, 3, 5, and 132, with the values of 20.1995, 5.0605, 3.7949, 3.2436,  
665 1.7420, 1.6297, respectively, with the last one being four-sided welded. Com-  
666 pared to Datta and Figureira [6], where the optimal value is  $f(\mathbf{x}) = 1.9553$ ,  
667 our obtained result  $f(\mathbf{x}) = 1.6297$  is smaller, because in our formulation  $h$ ,  $t$ ,  
668 and  $b$  are continuous variables, in contrast to Datta and Figureira [6] with  $h$ ,  $t$ ,  
669  $b$  as discrete variables. Furthermore, the convergence occurs relatively fast, as  
670 the optimization algorithm exploits the most promising cluster by maximizing  
671 the acquisition function. This behavior can be explained by the fact that in  
672 this welded beam design example, different materials have significantly differ-  
673 ent cost objective functional value, which aids the optimization convergence.

674  
675 To further demonstrate the effectiveness of the proposed method, we com-  
676 pare with GA. Two versions of the proposed method are used. In the first  
677 version, every cluster are considered as independent, leaving no neighbor in  
678 the neighborhood, whereas in the second version, all the clusters are considered  
679 as neighbors.

680 The performance comparison is presented in Figure 7, showing that both  
681 variants of the mixed-integer BO clearly outperforms the GA in the welded  
682 beam design problem. The solution obtained from the GA is  $[0, 1, 0.24920115,$   
683  $5.30060037, 7.12520087, 0.25345267]$ , where the objective function is evaluated  
684 at 2.04016262. On the other hand, from the first variant (none is neighbor)  
685 of the proposed method, the solution obtained is  $[1, 1, 0.16934934, 5.61720010,$   
686  $4.90884889, 0.27985016]$ , where the objective function is evaluated at 1.68206763.  
687 From the second variant (all are neighbors) of the proposed method, the solu-

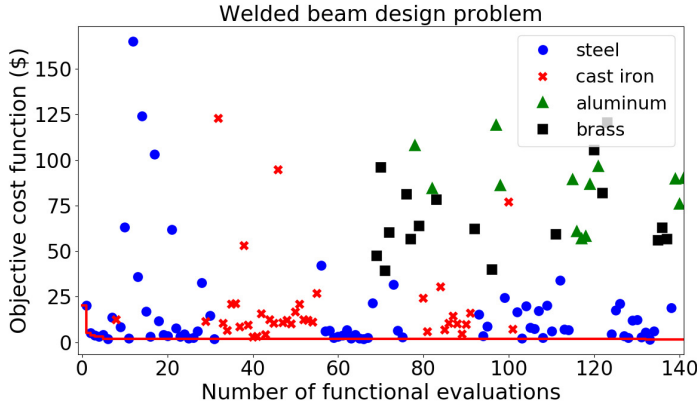


Fig. 6: Convergence plot of the cost function in the welded beam design, with all clusters are neighbors, showing different combinatorial of discrete and categorical variables are attempted.

tion obtained is  $[1, 1, 0.16934934, 5.61720010, 4.90884889, 0.27985016]$ , where  
the objective function is evaluated at 1.66457625. The convergence plots of  
these two variants are very similar. The asymptotic value using the second  
variant is slightly better than that using the first variant. However, we note  
that as the optimization process advances, the prediction converges to a single  
GP prediction, and thus both variants are similar at the later stage of  
search. The proposed mixed integer method clearly outperforms the GA in all  
settings.

#### 4.3 Pressure vessel design problem

Here, the proposed mixed-integer BO method is applied to solve the pressure  
vessel design optimization problem. The objective of this problem is to  
minimize the cost of a storage tank with  $3 \cdot 10^3$  psi internal pressure shown in  
Figure 8, where the minimum volume is  $750 \text{ ft}^3$ . The shell is made by joining  
two hemispheres and forming the longitudinal cylinder with another weld. The  
design variables are listed as follows.  $x_1$  is the thickness of the hemisphere.  $x_2$   
is the shell thickness.  $x_3$  is the inner radius of the hemisphere.  $x_4$  is the length  
of the cylinder.

The objective function that accounts for the cost is

$$f(\mathbf{x}) = 0.6224x_1x_3x_4 + 1.7781x_2x_3^2 + 3.1661x_1^2x_4 + 19.84x_1^2x_3, \quad (32)$$

where the imposed constraints are

$$g_1(\mathbf{x}) = -x_1 + 0.0193x_3 \leq 0, \quad g_2(\mathbf{x}) = -x_2 + 0.009541x_3 \leq 0, \quad (33a)$$

$$g_3(\mathbf{x}) = -\pi x_3^2 x_4^2 - \frac{4}{3} x_3^3 + 1296000 \leq 0, \quad g_4(\mathbf{x}) = x_4 - 240 \leq 0, \quad (33b)$$



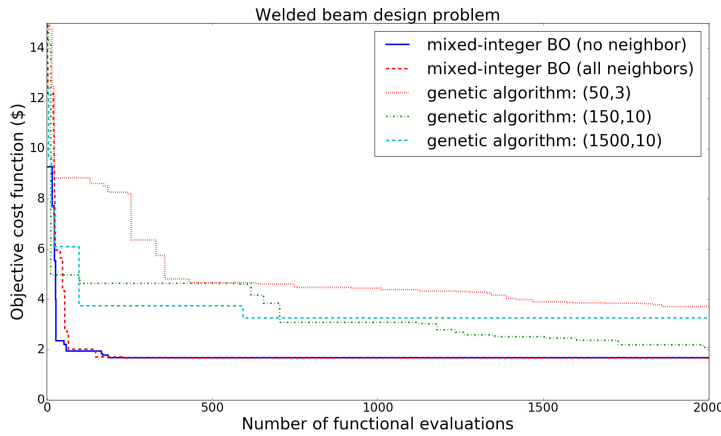


Fig. 7: Performance comparison between the GA and the proposed mixed-integer BO for the welded beam design.

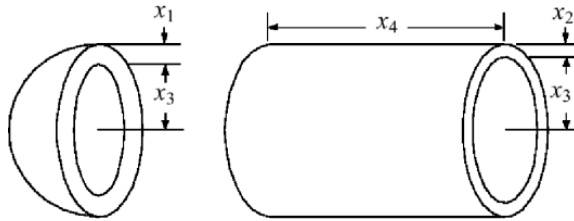


Fig. 8: Pressure vessel design optimization problem [4].

707 and  $0.00625 \leq x_1, x_2 \leq 0.61875$ ,  $10.0 \leq x_3, x_4 \leq 200.0$ . All variables are  
 708 considered as continuous in this example.

709 Figure 9 shows the performance comparison between the proposed mixed-  
 710 integer BO and the GA with various settings in terms of number of func-  
 711 tional evaluations. Again, the BO clearly shows its advantage in term of con-  
 712 vergence speed for continuous variables. The optimal input is  $[0.193114320,$   
 713  $0.0954997100, 10, 76.2478356]$ , where the corresponding objective functional  
 714 value is 125.02822748.

#### 715 4.4 Speed reducer design problem

716 Figure 10 shows the design optimization problem of a speed reducer [4].  
 717 Seven design variables are described as follows.  $x_1$  is the face width.  $x_2$  is  
 718 the module of teeth.  $x_3$  is the number of teeth on pinion.  $x_4$  is the length  
 719 of the first shaft between bearings.  $x_5$  is the length of the second shaft be-  
 720 tween bearings.  $x_6$  is the diameter of the first shaft.  $x_7$  is the diameter of

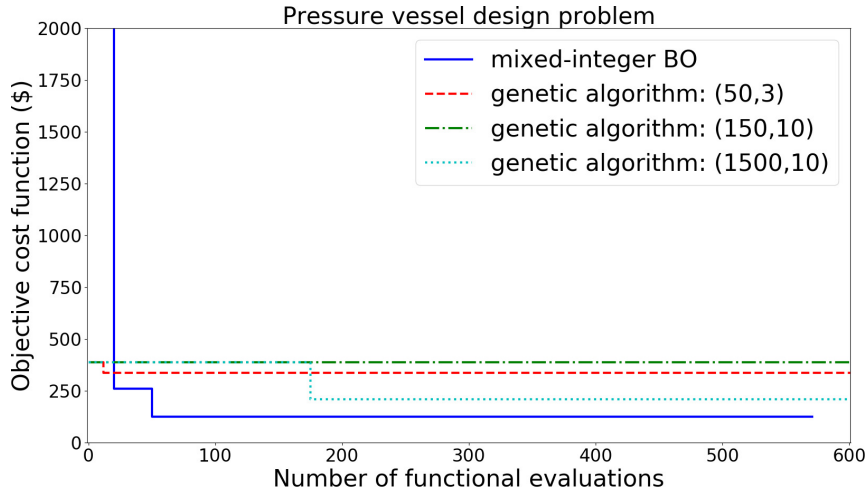


Fig. 9: Performance comparison between the GA and the proposed mixed-integer BO for the pressure vessel design.

721 the second shaft.  $x_3$  is the discrete variable, whereas the rest of the variables  
 722 are continuous. The problem is 7-dimensional, one discrete and six contin-  
 723 uous. With the formulation of the problem, there are 12 local GPs corresponding to 12 discrete values of  $x_3$ . In iteration 148, the mixed-integer BO

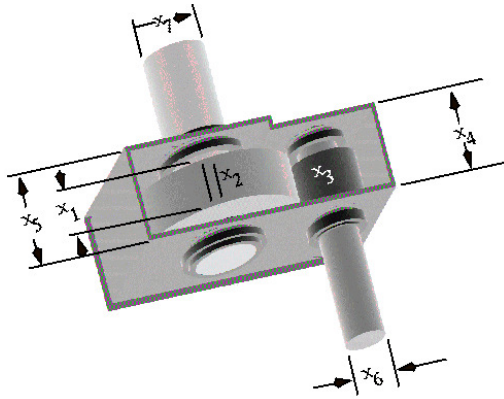


Fig. 10: Speed reducer design optimization problem [4] from NASA.

724 converges to the global minimum of  $f(\mathbf{x}^*) = 2996.29614837$ , where  $\mathbf{x}^* =$   
 725  $[3.50000447, 0.7, 17, 7.30566156, 7.8, 3.35022572, 5.28668406]$ . The result is com-  
 726 parable with Cagina et al. [4], where particle swarm optimization is employed,  
 727 yielding the optimal  $f(\mathbf{x}^*) = 2996.348165$ , where  $\mathbf{x}^* = 3.5, 0.7, 17, 7.3, 7.8, 3.350214, 5.286683]$ .  
 728  
 729

730 To evaluate the effect of initial sample size, the mixed-integer BO is per-  
 731 formed with different number of initial samples. Figure 11 shows the conver-  
 732 gence plot of the GA and the mixed-integer BO, each with various settings. In  
 733 terms of the number of functional evaluations, the mixed-integer BO clearly  
 734 shows the advantages with faster convergence, compared to the GA. The effect  
 735 of initial samples is also shown in Figure 11. It is observed that the proposed  
 736 mixed-integer BO converges relatively fast after the initial sampling stage.  
 737 Thus, for low-dimensional problems, it may not be necessary to sample exten-  
 738 sively at the initial sampling stage. The balance between exploration and  
 739 exploitation is well-tuned by the acquisition function, which is GP-UCB [53]  
 in this case.

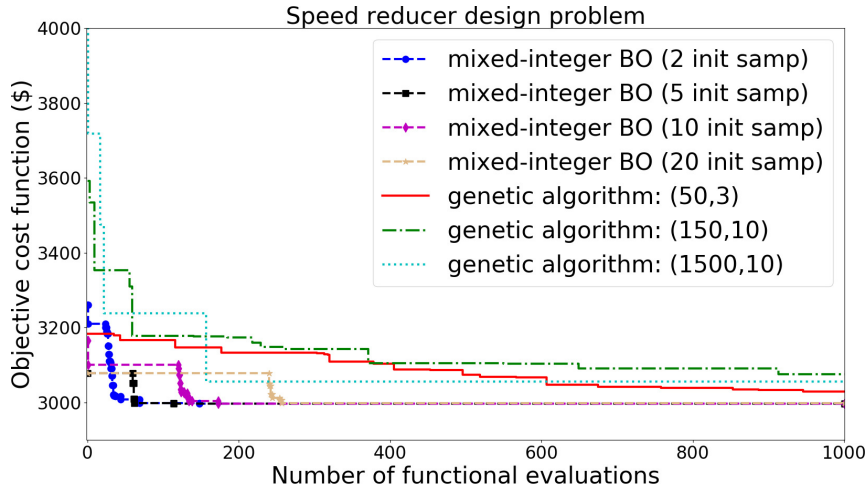


Fig. 11: Performance comparison between the GA and the proposed mixed-integer BO with different initial samples for the speed reducer design.

740

#### 741 4.5 High-dimensional discrete sphere function

742 To evaluate the performance of the proposed mixed-integer BO in high-dimensional  
 743 problems, two discrete sphere functions with 5-dimensional discrete variables  
 744 and 50-dimensional and 100-dimensional continuous variables, respectively, are  
 745 used to benchmark. The discrete sphere function is

$$f(\mathbf{x}^{(d)}, \mathbf{x}^{(c)}) = f(x_1, \dots, x_n, x_{n+1}, \dots, x_m) = \prod_{i=1}^n |x_i| \left( \sum_{j=n+1}^m x_j^2 \right) \quad (34)$$

746 where  $1 \leq x_i \leq 2$  ( $1 \leq i \leq n$ ) are  $n$  integer variables and  $-5.12 \leq x_j \leq$   
 747  $5.12(n+1 \leq j \leq m)$  are  $m-n$  continuous variables. Again, GA is used to

748 compare against the proposed mixed-integer BO method. The global optimal  
 749 of this function is  $f(\mathbf{x}^*) = 0$ , where  $\mathbf{x}^* = [1, 1, 1, 1, 1, 0, \dots, 0]$ . The number of  
 750 clusters in this example is  $2 \times 2 \times 2 \times 2 \times 2 = 32$ , where each cluster corresponds  
 751 to a local GP. Figure 12 shows the convergence plot of the proposed mixed-  
 752 integer BO with different number of initial samples and GA with different  
 753 settings for the (50+5)D discrete spherical function, where 5 variables are  
 754 discrete and 50 variables are continuous.

755 As seen in Figure 12, the proposed mixed-integer BO quickly identifies the  
 756 discrete tuple (1, 1, 1, 1, 1) that corresponds to the minimal response, with re-  
 757 spect to the discrete tuple. The rest of the convergence plot focuses on the  
 758 optimization of the continuous variables. The GA with population size of 50  
 759 and elite count of 3 performs on par with the proposed mixed-integer BO,  
 760 whereas other GA settings converge much slower. The mixed-integer BO with  
 761 2 initial samples converges relatively fast at the beginning. However, the con-  
 762 vergence at the later stage stagnates over a long period. On the contrary, the  
 763 mixed-integer with 20 initial samples converge very fast right after the initial  
 764 sampling stage. One of the reasons is that the local GP is able to approximate  
 765 the objective function more accurately with more initial samples, compared  
 to the one with less initial samples.

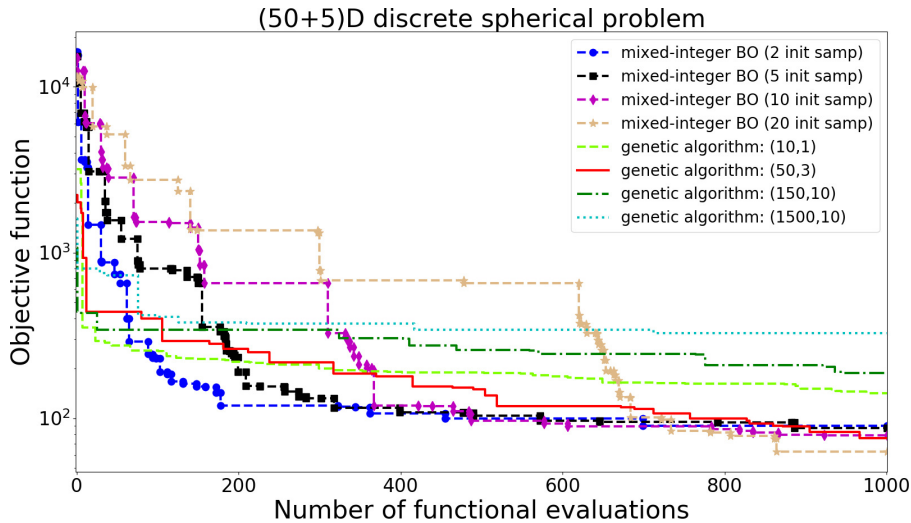


Fig. 12: Performance comparison between the GA and the proposed mixed-integer BO with different initial samples for (50+5)D discrete spherical function.

766

767

768

769

770

Similarly, Figure 13 shows the convergence plot of the proposed mixed-integer BO with a different number of initial samples and GA with different settings for (100+5)D discrete spherical function, where 5 variables are discrete. The mixed-integer with 2 initial samples converges poorly, whereas other

771 variants perform better. One of the reasons is that with the low initial sample  
 772 size, the discrete tuple is incorrectly identified as  $(1,1,1,1,2)$ , as opposed to  
 773  $(1,1,1,1,1)$ . The other variants of the proposed mixed-integer BO are able to  
 774 identify the correct tuple immediately after the initial sampling stage. Thus,  
 it may be beneficial to have sufficient number of initial samples.

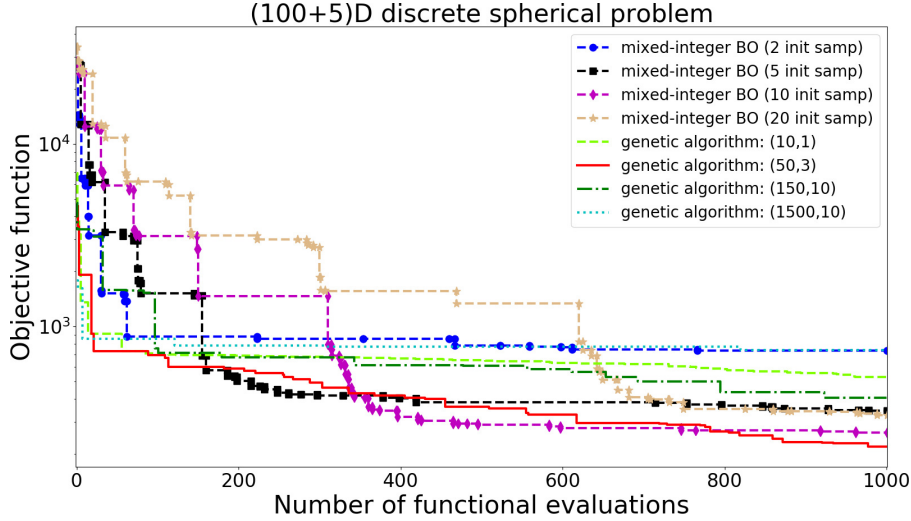


Fig. 13: Performance comparison between the GA and the proposed mixed-integer BO with different initial samples for  $(100+5)$ D discrete spherical function.

775

## 776 5 Metamaterials design examples

777 In this section, we demonstrate the applicability of the proposed method to the  
 778 design of metamaterials, in which properties can be tailored depending on the  
 779 geometric design of the structures. In Section 5.1, a mechanical metamaterial  
 780 is considered, where the objective is to design a low-weight and high-strength  
 781 unit cell. In Section 5.2, an auxetic metamaterial unit cell is considered. The  
 782 proposed BO method is applied to minimize the negative Poisson's ratio.

### 783 5.1 An example of designing high-strength low-weight fractal metamaterials

784 Motivated by the recent experimental work of Meza et al. [31] in designing  
 785 high-strength and low-weight metamaterials at nano-scale for ceramic systems  
 786 where the effective mechanical strength can be enhanced by hierarchical structure.  
 787 We demonstrate the proposed methodology in searching for high-strength

788 and low-weight metamaterials for multiple classes of materials. Particularly,  
 789 our metamaterials are constructed with fractal geometry. Fractal geometry has  
 790 the special property of self-similarity at different length scales. A parametric  
 791 design and optimization approach for fractal metamaterials is demonstrated  
 792 here. In this example, the goal is to maximize the effective strength of the  
 793 structure. The effective strength is defined as the ratio between the effective  
 794 Young modulus and the volume of material with the assumption of homoge-  
 795 nized material for the bulk properties. The material selection, including Ashby  
 796 chart, is formulated as an inequality constraint to limit the searching space of  
 797 materials.

### 798 5.1.1 Parametric design of fractal truss structures

799 Mathematically, fractals can be constructed iteratively using the so-called ite-  
 800 rated function systems (IFSs). An IFS is a finite set of contraction mappings  
 801  $\{f_i\}_{i=1}^N$  on a complete metric space  $X$  [1]. Starting from an initial set  $\mathcal{P}_0$ , the  
 802 fractal can be constructed iteratively as  $\mathcal{P}_{k+1} = \cup_{i=1}^N f_i(\mathcal{P}_k)$ . Geometrically,  
 803 the IFSs  $f_i$  can be expressed in terms of rotation, translation, scaling, and  
 804 other set topological operations, such as complement, union or intersect.

805 In this example, the fractal truss structures are constructed from the 2D  
 806 profiles shown in Figure 14c. They are based on the square shape, even though  
 807 in principle they can be constructed from any arbitrary polygon such as trian-  
 808 gle and hexagon. Figure 14c presents the first three levels of IFS construction.  
 809 The IFSs are inspired by the projection of Keplerian 3D fractals onto its cor-  
 810 responding 2D plane. Here, the IFS operators include the translation matrix  
 811  $T = \text{diag}\{\pm d/2, \pm d/2, 1\}$  and the scaling matrix  $S = \text{diag}\{1/2, 1/2, 1\}$ . The  
 812 rotation is not considered. Physically, the first four IFSs simply scale the design  
 813 of previous fractal level by  $1/2$ , and translate them to the northwest, north-  
 814 east, southwest, and southeast, respectively. The fifth IFS scales the design  
 815 of previous fractal level by one half, and deletes other features that overlaps  
 816 within the region.

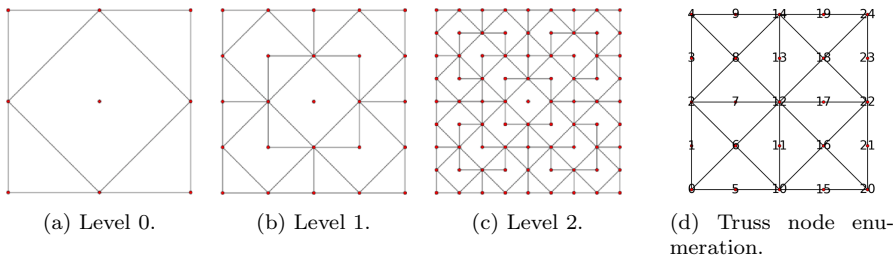


Fig. 14: Truss design parameters on the unit square: (a) (b) (c) Iterated function systems of truss designs on unit square, level 0-2, respectively. (d) Truss options on fractal level 0 unit square.

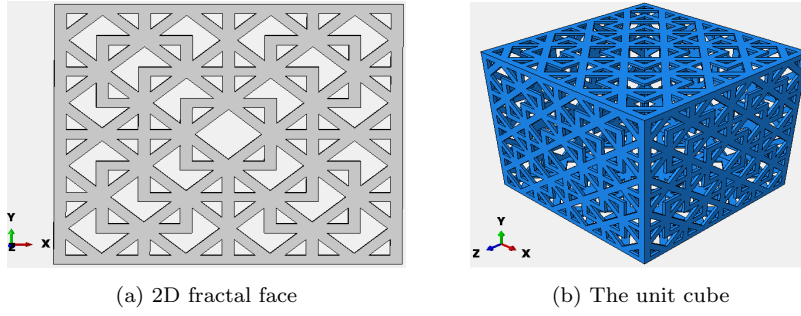


Fig. 15: Design of fractal unit cube. (a) The 2D fractal profile with a fractal level of 2 and only inner square truss option enabled. (b) The unit cube is composed of six identical fractal faces, and each face is designed by truss options, thickness, and extrusion depth

817 Figure 14 illustrates the square basis with three design options: (1) diagonal  
 818 truss, (2) inner square truss, and (3) perpendicular truss. The diagonal truss  
 819 option enables edges connecting nodes 4, 8, 12, 16, 20 and nodes 0, 6, 12, 18,  
 820 24. The inner square truss option enables edges connecting nodes 2, 6, 10, 15,  
 821 22, 18, 14, 8. The perpendicular truss option enables edges connecting nodes 2,  
 822 7, 12, 17, 22 and nodes 10, 11, 12, 13, 14. In the example of Figure 14c, only the  
 823 inner square truss option is enabled. In the construction process, the options  
 824 are enabled by setting the truss control parameters to 0 or 1, respectively. The  
 825 fundamental adjacency matrix of fractal level 0 is built to indicate whether a  
 826 pair of nodes are connected. With the design of level 0 unit cell, the IFSs are  
 827 applied recursively to create the more complicated geometry at the desired  
 828 level. Once the profile is constructed, additional offset operations are applied  
 829 to generate thickness of the 2D truss elements for a full 3D structure. Figure  
 830 15a shows a complete 2D fractal face. With the square face defined, a complete  
 831 3D fractal unit cell is built with six of the faces, as shown in Figure 15b.

### 832 5.1.2 Constitutive material model and the finite element analysis

833 A general anisotropic material has 21 independent elastic constants to de-  
 834 scribe the stress-strain ( $\sigma$ - $\varepsilon$ ) relationship. To simplify the materials constitu-  
 835 tive model, we assume isotropic and linear elastic materials behavior at small  
 836 strain regime, where  $\sigma$ - $\varepsilon$  relationship for bulk material properties can be ob-  
 837 tained via Young's modulus  $E$  and Poisson's ratio  $\nu$ , i.e.

$$\sigma_{ij} = \frac{E}{1 + \nu} \left( \varepsilon_{ij} + \frac{\nu}{1 - 2\nu} \varepsilon_{kk} \delta_{ij} \right), \quad (35)$$

838 where  $i, j$  can be either  $x, y$ , or  $z$ , and  $\delta_{ij}$  is the Kronecker delta of  $i$  and  $j$ .  
 839 The material properties  $E$  and  $\nu$ , as well as materials  $\rho$ , are taken as inputs  
 840 to describe the linear elastic regime in the FEM simulation to obtain stress.

841 In simulations, we are concerned with an uniaxial compression. Therefore,  
 842 to simplify the terminology, we refer to the component of effective stiffness  
 843 tensor in the loading direction as effective Young's modulus. It is noteworthy  
 844 that the effective stiffness tensor of the designed fractal truss structure is not  
 845 the same as the bulk material stiffness tensor. Two displacement boundary  
 846 conditions are imposed on the unit cube. One is the fixed boundary condition  
 847 for both translation and rotation, and the other is the constant displacement  
 848 on the opposite side of the cube. The stress is obtained by taking the maximum  
 849 nodal stress in the active direction. The effective Young's modulus is calculated  
 850 as the ratio of the maximal nodal stress  $\sigma_{33}$  at the designated engineering strain  
 851  $\varepsilon = 0.01$ . The quadratic tetrahedral element (C3D10 in ABAQUS) is utilized  
 852 for the FEM simulation. The total number of elements is between 5,000 and  
 853 10,000. The exact number varies with respect to the finite element simulation.  
 854 The size of the cube is around 1mm ( $10^{-3}$ m).

855 The dimension of the design space is 9, in which 4 discrete and 5 continuous  
 856 variables are combined to create an input  $\mathbf{x} = (x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9)$ .  
 857 The discrete variables include fractal level, the diagonal, inner square, and  
 858 perpendicular truss options. The fractal level  $x_1$  is an integer of either 0, 1, or  
 859 2, whereas each of the truss options  $x_2, x_3, x_4$  is a binary variable from design  
 860 space, taking a value of 0 or 1. The continuous variables include thickness  
 861  $x_5 = t$  of the truss, the extrusion depth  $x_6 = et$  of the unit face, the materials  
 862 bulk density  $x_7 = \rho$ , bulk elastic Young's modulus  $x_8 = E$ , and bulk Poisson's  
 863 ratio  $x_9 = \nu$ .

864 Three constraints are imposed as follows. Thickness and extrusion depth  
 865 are limited to a constant that is related to the fractal level to preserve the  
 866 fractal geometry of the structure. The higher the fractal level is, the smaller  
 867 is the constant. Similarly, the material bulk density, Young's modulus, and  
 868 Poisson's ratio are bounded within a physical limit, where values are taken  
 869 from Table 3.1 of Bower [3] for woods, copper, tungsten carbide, silica glass,  
 870 and alloys. As a result, the imposed constraints are

$$\underline{T} \leq x_5 \leq \bar{T}, \quad x_6 \geq \bar{T}, \quad (36a)$$

$$x_5 \leq 7 \cdot x_6, \quad x_6 \leq 7 \cdot x_5, \quad (36b)$$

871 where  $\underline{T} = 10^{-6}$  is the threshold for manufacturability, and  $\bar{T}$  is the threshold  
 872 for the truss thickness as

$$\bar{T} = \begin{cases} \frac{1}{2 \cdot 2^{x_1+1}}, & \text{if } x_3 = x_4 = 2, \\ \frac{1}{2 \cdot 2^{x_1}}, & \text{otherwise.} \end{cases} \quad (37)$$

873 We expect the simulations to converge on the high-strength and low-density  
 874 type of materials. However, Ashby chart indicates a high correlation between  
 875 compressive strength and density among all types of materials. To circumvent  
 876 this problem, another constraint is introduced to limit the search region, based  
 877 on the upper bound of longitudinal wave speed as  $\sqrt{E/\rho} = \sqrt{x_8/x_7} \leq 10^{4.25}$   
 878 m/s.



## 879 5.1.3 Simulation and results

880 Figure 16 shows an example of von Mises stress during the uniaxial compression  
 881 of the architected metamaterial cell, as described in Section 5.1.2. In the  
 882 simulation settings and its post-process, only  $\sigma_{zz}$  is concerned.

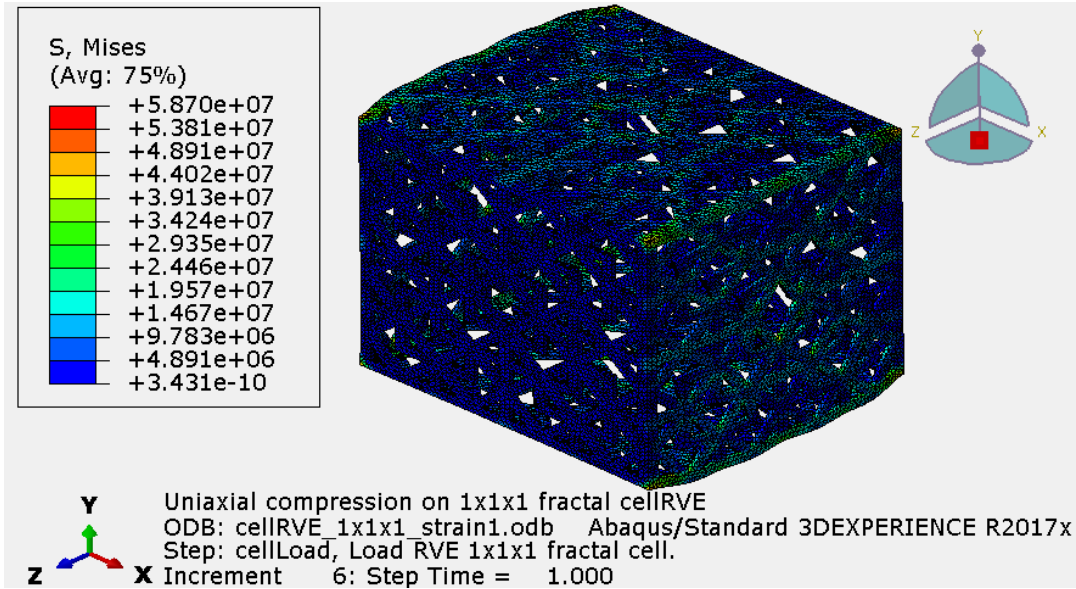


Fig. 16: An example of von Mises stress of the structure under loading condition.

883 The lower bounds of continuous variables  $(x_5, x_6, x_7, x_8, x_9)$  are  $(2 \cdot 10^{-6}, 2 \cdot$   
 884  $10^{-6}, 0.4 \cdot 10^{+3}, 9 \cdot 10^{+9}, 0.16)$ . The lower bounds of  $x_7, x_8, x_9$  correspond to the  
 885 density of wood, bulk Young's modulus of wood, and Poisson's ratio of silica  
 886 glass, respectively. The upper bounds of continuous variables  $(x_5, x_6, x_7, x_8, x_9)$   
 887 are  $(0.5 \cdot 10^{-3}, 0.5 \cdot 10^{-3}, 8.9 \cdot 10^{+3}, 650 \cdot 10^{+9}, 0.35)$ . The upper bounds of  
 888  $x_7, x_8, x_9$  correspond to the density of copper, bulk Young's modulus of tung-  
 889 sten carbide, and Poisson's ratio of a general alloy, respectively.

890 To initialize the optimization process, two random inputs are sampled to  
 891 construct the GP model for each cluster. The number of clusters in this exam-  
 892 ple is  $2 \times 2 \times 2 \times 3 = 24$ . The EI acquisition is used to locate the next sampling  
 893 location  $\mathbf{x}$ . The CMA-ES [17] is used as an auxiliary optimizer to maximize  
 894 the penalized acquisition function. The optimization process is carried out for  
 895 170 iterations, as shown in Figure 17. At iteration 0, 1, 2, 11, 14, 26, 148, better  
 896 objective function values of 1.9723, 2.7827, 10.4725, 12.1207, 22.1071, 23.3766,  
 897  $36.8316 \cdot 10^6$  GPa/kg, are identified, respectively. The relatively fast conver-  
 898 gence plot demonstrates the effectiveness of the proposed BO method for the

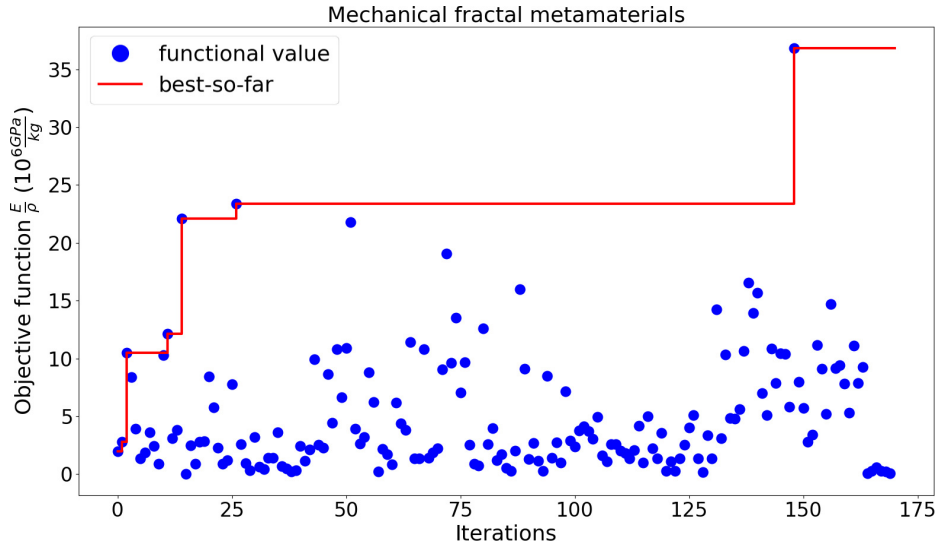


Fig. 17: Convergence plot of the objective function, which is the ratio between the effective Young’s modulus and the weight of the cell, i.e.  $E_{\text{eff}}/m$ .

899 mix-integer optimization problems. Due to the expensive computational cost  
 900 of the FEM simulation, the number of iterations is limited to 200.

## 901 5.2 Design optimization of fractal auxetic metamaterials

902 In the second example, we study the auxetic metamaterial with application  
 903 in flexible and stretchable devices. Inspired by the experimental work of Cho  
 904 et al. [5] in designing auxetic metamaterials using fractal cut, and its sub-  
 905 sequent numerical and experimental work by Tang et al. [57] in developing  
 906 shape-programmable materials, we use auxetic metamaterials to demonstrate  
 907 the proposed BO methodology. The goal of this example is to minimize the  
 908 effective Poisson’s ratio, which is negative and evaluated through a FEM sim-  
 909 ulation.

### 910 5.2.1 Parametric design of auxetic metamaterials

911 Here, a parametric design of the unit cell, where the fractal level is fixed at 2,  
 912 is devised. The cut motif  $\alpha$  and  $\beta$  for one level of the auxetic cell is shown in  
 913 Figure 18. Basically, this cut motif controls the free rotational hinges of the  
 914 architected structure, such that the deformation energy dissipates through  
 915 rotational motion, rather than translational motion. The principle of cut de-  
 916 sign is based on the connectivity of the rotating units, where the connectivity  
 917 depends on the cut patterns, which in turn determines the maximum stretch-  
 918 ability of the designed specimen. For further details about the fractal cut and

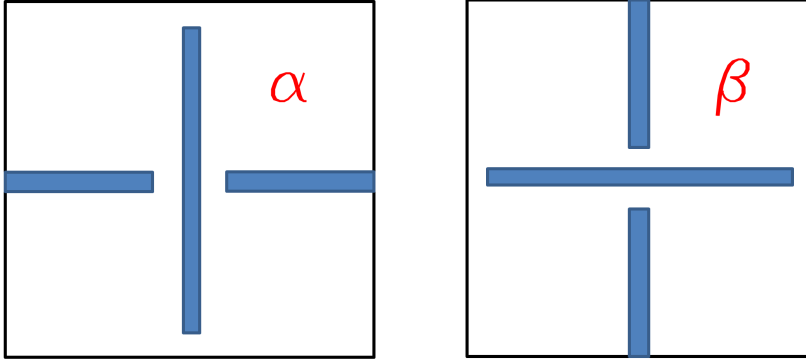


Fig. 18: Cut motif  $\alpha$  and  $\beta$  in designing auxetic metamaterials by fractal cuts.

its rotating mechanisms, readers are referred to the work of Cho et al. [5] and Tang et al. [57]. To create a fractal cut, a simple IFS is imposed on the cut to create subsequent level, with the scaling ratio of  $1/2$ , and is then translated to four corners.

To tailor the negative Poisson's ratio, the shape of the cut is modeled as splines, where the coordinates of the control points are considered as inputs. The choice of  $\alpha$  and  $\beta$  cut is formulated using discrete variables. The dimension of this problem is 18, in which 2 discrete and 16 continuous variables are used. The parametric input  $\mathbf{x}$  includes  $x_1, x_2$  as discrete variables, which takes value of either 1 ( $\alpha$ -motif) or 2 ( $\beta$ -motif) for level 1 and level 2 cuts, respectively. The first 4 continuous variables  $x_3, x_4, x_5, x_6$  are used to describe the shape of the large center cut of level 1. The next 4 continuous variables  $x_7, x_8, x_9, x_{10}$  describe the shape of two small side cuts of level 1. In the same manner, the last 8 continuous variables are used to model the large center cut and two small side cuts of level 2. Figure 19 shows an example of the parametric design implementation of the designed auxetic metamaterials in the ABAQUS environment. The solid dots represent the control points of the cut. (Color is available on the electronic version. The blue solid dots denote the level 1 control points, whereas the red solid dots denote the level 2 control points.)

### 5.2.2 Constitutive material model and the finite element analysis

The study of Tang et al. [57] has demonstrated that the effective Poisson's ratio  $\nu_{\text{eff}}$  is indeed a function of strain  $\epsilon$ . In this work, we assume that the base material is natural rubber reinforced by carbon-black. Mooney-Rivlin constitutive model is used to describe the hyperelastic material behavior, where the suitable energy function  $W$  is expressed as

$$W = C_{10}(\bar{I}_1 - 3) + C_{01}(\bar{I}_2 - 3) + \frac{1}{D_1}(J - 1)^2, \quad (38)$$

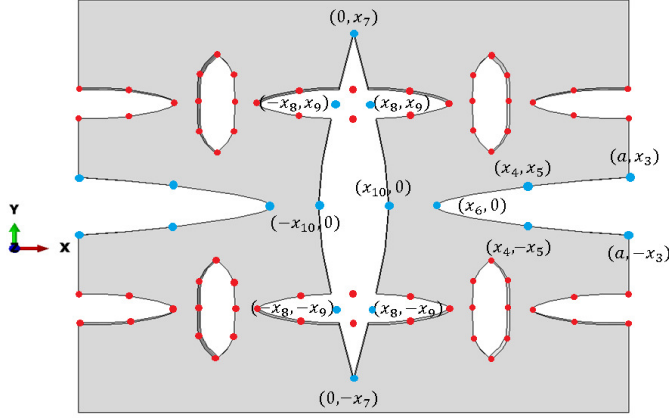


Fig. 19: An implemented example of auxetic metamaterials by fractal cuts. The solid dots present the control points of the cut. (Color is available on the electronic version. Blue dots correspond to level 1, whereas red dots correspond to level 2.)

944 where  $J$  is the elastic volume ratio,  $I_1, I_2, I_3$  are the three invariants of Green  
 945 deformation tensor defined in term of principal stretch ratios  $\lambda_1, \lambda_2, \lambda_3$ , i.e.

$$I_1 = \sum_{i=1}^3 \lambda_i^2, \quad I_2 = \sum_{i,j=1;i \neq j}^3 \lambda_i \lambda_j, \quad I_3 = \prod_{i=1}^3 \lambda_i, \quad (39)$$

946 and  $\bar{I}_1 = I_1 J^{-2/3}$ ,  $\bar{I}_2 = I_2 J^{-4/3}$ . The materials parameter is adopted from  
 947 Shahzad et al. [48], where  $C_{10} = 0.3339\text{MPa}$ ,  $C_{01} = -3.37 \cdot 10^{-4}$ , and  $D_1 =$   
 948  $1.5828 \cdot 10^{-3}$ .

949 The initial size of the square is 20 cm  $\times$  20 cm, and the thickness of  
 950 the specimen is 1mm. The specimen is then deformed in a uniaxial tension  
 951 configuration in  $y$ -direction, where the displacement is fixed at 10 cm in one  
 952 direction. The configuration for the simulation is plane-strain configuration,  
 953 where displacement in the extrusion direction ( $z$ -direction) is fixed as zero.

954 In the deformed configuration, we extract the displacement in  $x$ -direction  
 955 to infer the engineering transverse strain, and compute the effective Poisson's  
 956 ratio as the ratio between transverse and longitudinal engineering strain.

957 The element used in this FEM simulation is the eight-node brick element  
 958 (C3D8R, C3D6, and C3D4). The FEM is developed in the ABAQUS environ-  
 959 nment. The number of elements for each simulation is approximately 5,000.

960 In this example, several constraints are imposed on the design variables,  
 961 which are

$$x_5 \leq 0.010 - t, \quad x_8 \leq x_4 - t, \quad x_{16} \leq x_{12} - t \quad (40a)$$

$$0 \leq x_6 \leq x_8, \quad 0 \leq x_7 \leq x_5, \quad x_4 \leq x_2 \leq 0.010, \quad 0 \leq x_3 \leq x_1 \quad (40b)$$

962 where  $\underline{t} = 0.0015$  m is the smallest thickness of the specimen. Two other  
 963 constraints include the implementation of convexity for the large center cut of  
 964 level 1 and level 2. Figure 20 presents an example of deformed configuration  
 965 after the simulation converges.

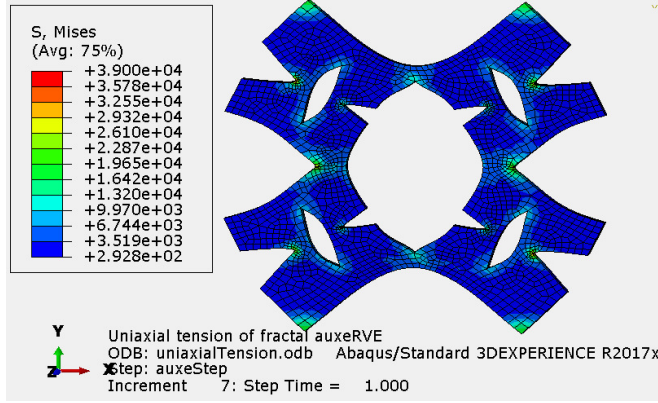


Fig. 20: An example of uniaxial tension simulation of plane-strain configuration in designing auxetic metamaterials using fractal cut.

### 966 5.2.3 Simulation and results

967 The lower bounds of the continuous variables are  $(0.25; 3.5; 0.50; 1.75; 8.0;$   
 968  $0.25; 4.0; 0.50; 0.25; 3.5; 0.50; 1.75; 4.0; 0.25; 3.0; 0.50) \cdot 10^{-3}$ . The upper bounds  
 969 of the continuous variables are  $(2.00; 6.5; 1.75; 3.00; 9.5; 1.50; 8.0; 1.75; 2.00;$   
 970  $6.5; 1.75; 3.00; 5.5; 1.50; 4.0; 1.75) \cdot 10^{-3}$ .

971 Two random initial sampling points are created within each cluster. Be-  
 972 cause the fractal level is fixed at 2, where each fractal level corresponds to one  
 973 cut motif  $\alpha$  or  $\beta$ , 4 clusters are created during the initialization. The initial  
 974 hyper-parameters  $\theta_i$  for all  $i$  are set at 0.2. The lower and upper bounds for  
 975 the hyper-parameters  $\theta_i$  for all  $i$  are  $(0.01, 20)$ .

976 The optimization process is carried out for 790 iterations. Figure 21 shows  
 977 the convergence plot of the optimization process, where the best objective  
 978 function value  $\nu_{\text{eff}}$  is updated in iterations 0, 4, 24, 26, 30, 45, 63, 66, 69,  
 979 78, 81, 84, 513, 582, 647, with the value of -0.6603, -0.6605, -0.6628, -0.6628,  
 980 -0.6902, -0.6941, -0.7143, -0.7410, -0.7517, -0.7576, -0.7627, -0.7784, -0.7785,  
 981 -0.7802, -0.7804, respectively. The proposed BO shows relatively fast conver-  
 982 gence for mid-level dimensionality  $d = 16$ , thus demonstrating the effectiveness  
 983 in tackling mix-integer nonlinear optimization problems.

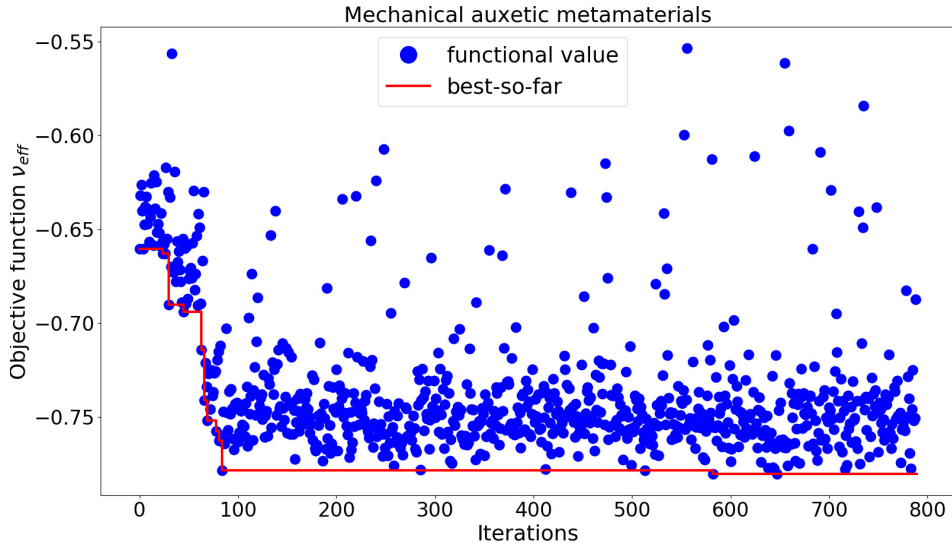


Fig. 21: Convergence plot of the objective function, which is the effective Poisson’s ratio  $\nu_{\text{eff}}$ . The best objective function value is updated at iterations 0, 4, 24, 26, 30, 45, 63, 66, 69, 78, 81, 84, 513, 582, 647, sequentially.

## 984 6 Discussion

985 One of the advantages of the proposed BO algorithm is its extension to incor-  
 986 porate discrete variables for nonlinear mixed-integer optimization problems.  
 987 The discrete variables include both categorical and integer variables, thus can  
 988 be applied with or without the notion of order. The neighborhood of each  
 989 cluster is built once during the initialization of the process, and can be cus-  
 990 tomized to adapt to specific user-defined requirements. Additionally, because  
 991 the neighborhood can be modified and/or defined manually, the independence  
 992 between clusters can be achieved by removing the corresponding clusters. Such  
 993 independence is quite common in the case of categorical variables. However,  
 994 the optimization performance of the proposed method does not depend on the  
 995 enumeration of the clusters. We emphasize that if the cluster is ceased to exist,  
 996 then it can be manually removed, and the cluster indices can be reenumerated  
 997 manually by a slight modification of Equation 12 and Algorithm 1.

998 The weight computation scheme is devised in such a way that asymptoti-  
 999 cally, the weight prediction converges to a single GP prediction, by imposing  
 1000 a weight vector which has 0 everywhere, except for a single 1 that corresponds  
 1001 to the corresponding cluster. It is recommended to choose the neighbors care-  
 1002 fully. One way to do so is to set a small threshold discrete distance  $d_{\text{th}}$ , which  
 1003 measures the dissimilarity between clusters based on the discrete tuples, e.g.  
 1004  $d_{\text{th}} \leq 1$ , and manually remove clusters that are known to be independent  
 1005 beforehand at the end of initialization. The safest setting is  $d_{\text{th}} = 0$ , which

1006 assumes clusters are completely independent of each other. This setting has  
1007 some negative effect on the convergence rate, but would eventually reach the  
1008 global optimal solution, and would not be trapped at local optima.

1009 The initial sample size plays a role in the performance of the proposed  
1010 mixed-integer BO method. It has been shown that for some low-dimensional  
1011 problems, the initial sample size does not affect the optimization performance.  
1012 However, for high-dimensional problems, the initial sample size does impact  
1013 the optimization performance. Too many initial samples at the beginning  
1014 would prevent the optimization from quick convergence. However, with moder-  
1015 ate amount of initial samples, and thus a more accurate local GP, the mixed-  
1016 integer BO converges faster, compared with fewer initial samples. As a general  
1017 rule of thumb, the total initial sample size is recommended at between  $5d$  and  
1018  $10d$ , where  $d$  is the dimension of the problem, including both discrete and  
1019 continuous variables.

1020 Here the scalability of GP for high-dimensional problems is alleviated, but  
1021 not completely eliminated. It is noted that the decomposition and weighted  
1022 average approach has been adopted [37, 39, 38, 40, 54, 58] for continuous vari-  
1023 ables. The decomposition method for continuous variables is typically referred  
1024 to as local GP. This approach is promising in tackling the scalability problem.  
1025 Particularly, in one of our previous studies [58], we have shown that the local  
1026 GP is computationally one-order cheaper, compared to the classical GP, while  
1027 maintaining a reasonable approximation error. Nevertheless, further research  
1028 is required to develop an efficient and robust decomposition scheme for both  
1029 discrete and continuous variables.

1030 One of the limitations in the proposed approach is the scalability with re-  
1031 spect to discrete variables. Because of the decomposition scheme, the number  
1032 of the clusters is the number of the combinatorial possibilities, i.e. the product  
1033 of the number of choices for each discrete variable, and thus resulting in the  
1034 sparsity problem in each cluster. To mitigate the undesirable sparsity effect,  
1035 a Gaussian mixture model that combines all the predictions from neighboring  
1036 clusters is used to exploit some useful information from the neighborhood.  
1037 As mentioned previously, the mixed-integer optimization problem, in general,  
1038 is difficult, because it combines the difficulties for both discrete and continu-  
1039 ous optimization. Particularly, some discrete and combinatorial optimization  
1040 problems are NP-complete, such as the traveling salesman problem, knapsack  
1041 problem, and graph coloring problem, to name a few. Another extension is to  
1042 model the weights as stochastic variables, so that the metaheuristic method-  
1043 ologies can be applied [2].

1044 The clustering and enumeration algorithm described in Algorithm 1 is  
1045 based on the assumption of the independence of discrete variables. Algorithm  
1046 1 does not work if the discrete variables are dependent. However, in the case  
1047 that discrete variables are dependent on each other, manual neighborhood  
1048 definition of clusters can be introduced manually, and the proposed BO algo-  
1049 rithm is functional with the demonstrated efficiency. However, the users must  
1050 declare the neighborhood of each cluster manually. Strictly speaking, the com-  
1051 putational efficiency of the proposed algorithm only depends on the number

of clusters, not the number of discrete variables. If all discrete variables are completely independent of each other, as demonstrated in the above examples, then the number of clusters is equal to the product of the number of choices for each discrete variable, i.e.  $L = \prod p_i$ .

Another practical limitation for the proposed BO algorithm for engineering models and simulations is its sequential nature of sampling and search. Each run of simulations usually demands a considerable amount of computational time. In practice, for high-fidelity and dedicated simulations, one should resort to multi-fidelity or batch-parallel BO for further improvement.

## 7 Conclusion and Future Work

In this paper, we propose a new BO algorithm to solve the nonlinear constrained mixed-integer design optimization problems. In this algorithm, the large dataset is decomposed according to the discrete tuples, in which each discrete tuple corresponds to a unique GP model. The prediction for mean and variance is formulated as a Gaussian mixture model, in which the weights are computed based on the pair-wise Wasserstein distance between clusters. Constraints, which are formulated as a set of inequalities, are included during the optimization process. Theoretical bounds and algorithmic complexity are provided to demonstrate the computational efficiency compared to the classical GP.

The proposed algorithm is demonstrated with two fractal metamaterials design examples, where the mechanical properties are tailored by the hierarchically designed architect. In the first example, the algorithm is used to search for the fractal metamaterial with high-strength and low-density properties, where material selection is considered. In the second example, the algorithm is utilized to design an auxetic metamaterial for flexible and stretchable devices, where the effective Poisson's ratio is chosen as the objective function. For both computational materials design examples, constraints are imposed to limit the design space. The proposed algorithm shows a promising performance in solving engineering problems, where high dimensionality is often an issue.

While several limitations exist, such as scalability for discrete and continuous variables, further research extensions can be made to improve the current methodology, including metaheuristic methodologies for stochastic combinatorial optimization.

## Acknowledgments

The research was supported in part by the National Science Foundation under grant number CMMI-1306996. Authors thank Prof. Hongyuan Zha (Georgia Tech) for numerous helpful conversations about Bayesian optimization. This research was supported in part through research cyberinfrastructure resources



and services provided by the Partnership for an Advanced Computing Environment (PACE) at the Georgia Institute of Technology, Atlanta, Georgia, USA. The authors are grateful to two anonymous reviewers for their constructive feedback.

## References

1. Barnsley, M.F.: *Fractals everywhere*. Academic press (2014)
2. Bianchi, L., Dorigo, M., Gambardella, L.M., Gutjahr, W.J.: A survey on metaheuristics for stochastic combinatorial optimization. *Natural Computing* **8**(2), 239–287 (2009)
3. Bower, A.F.: *Applied mechanics of solids*. CRC press (2011)
4. Cagnina, L.C., Esquivel, S.C., Coello, C.A.C.: Solving engineering optimization problems with the simple constrained particle swarm optimizer. *Informatica* **32**(3) (2008)
5. Cho, Y., Shin, J.H., Costa, A., Kim, T.A., Kunin, V., Li, J., Lee, S.Y., Yang, S., Han, H.N., Choi, I.S., et al.: Engineering the shape and structure of materials by fractal cut. *Proceedings of the National Academy of Sciences* **111**(49), 17390–17395 (2014)
6. Datta, D., Figueira, J.R.: A real-integer-discrete-coded particle swarm optimization for design problems. *Applied Soft Computing* **11**(4), 3625–3633 (2011)
7. Davis, E., Ierapetritou, M.: A kriging based method for the solution of mixed-integer nonlinear programs containing black-box functions. *Journal of Global Optimization* **43**(2-3), 191–205 (2009)
8. Deb, K., Goyal, M.: A combined genetic adaptive search (GeneAS) for engineering design. *Computer Science and informatics* **26**, 30–45 (1996)
9. Digabel, S.L., Wild, S.M.: A taxonomy of constraints in simulation-based optimization. arXiv preprint arXiv:1505.07881 (2015)
10. Gandomi, A.H., Yang, X.S.: Benchmark problems in structural optimization. In: *Computational optimization, methods and algorithms*, pp. 259–281. Springer (2011)
11. Gardner, J.R., Kusner, M.J., Xu, Z.E., Weinberger, K.Q., Cunningham, J.P.: Bayesian optimization with inequality constraints. In: *ICML*, pp. 937–945 (2014)
12. Gelbart, M.A., Snoek, J., Adams, R.P.: Bayesian optimization with unknown constraints. arXiv preprint arXiv:1403.5607 (2014)
13. Givens, C.R., Shortt, R.M., et al.: A class of wasserstein metrics for probability distributions. *The Michigan Mathematical Journal* **31**(2), 231–240 (1984)
14. Gramacy, R.B., Lee, H.K.: Gaussian processes and limiting linear models. *Computational Statistics & Data Analysis* **53**(1), 123–136 (2008)
15. Gramacy, R.B., Lee, H.K.H.: Bayesian treed Gaussian process models with an application to computer modeling. *Journal of the American Statistical Association* **103**(483), 1119–1130 (2008)
16. Gramacy, R.B., Taddy, M., et al.: Categorical inputs, sensitivity analysis, optimization and importance tempering with tgp version 2, an R package for treed Gaussian process models. *Journal of Statistical Software* **33**(6), 1–48 (2010)
17. Hansen, N., Müller, S.D., Koumoutsakos, P.: Reducing the time complexity of the derandomized evolution strategy with covariance matrix adaptation (CMA-ES). *Evolutionary computation* **11**(1), 1–18 (2003)
18. Hemker, T., Fowler, K.R., Farthing, M.W., von Stryk, O.: A mixed-integer simulation-based optimization approach with surrogate functions in water resources management. *Optimization and Engineering* **9**(4), 341–360 (2008)
19. Hernández-Lobato, J.M., Gelbart, M., Hoffman, M., Adams, R., Ghahramani, Z.: Predictive entropy search for Bayesian optimization with unknown constraints. In: *International Conference on Machine Learning*, pp. 1699–1707 (2015)
20. Hernández-Lobato, J.M., Gelbart, M.A., Adams, R.P., Hoffman, M.W., Ghahramani, Z.: A general framework for constrained bayesian optimization using information-based search. *Journal of Machine Learning Research* (2016)
21. Huang, D., Allen, T.T., Notz, W.I., Zeng, N.: Global optimization of stochastic black-box systems via sequential kriging meta-models. *Journal of Global Optimization* **34**(3), 441–466 (2006)

- 1146 22. Jang, H.L., Cho, H., Choi, K.K., Cho, S.: Reliability-based design optimization of fluid–  
1147 solid interaction problems. Proceedings of the Institution of Mechanical Engineers, Part  
1148 C: Journal of Mechanical Engineering Science **228**(10), 1724–1742 (2014)
- 1149 23. Kim, K., Lee, M., Lee, S., Jang, G.: Optimal design and experimental verification of  
1150 fluid dynamic bearings with high load capacity applied to an integrated motor propulsor  
1151 in unmanned underwater vehicles. Tribology International **114**, 221–233 (2017)
- 1152 24. Kim, Y., Lee, S., Yee, K., Rhee, D.H.: High-to-low initial sample ratio of hierarchical  
1153 kriging for film hole array optimization. Journal of Propulsion and Power (2017)
- 1154 25. Li, M., Li, G., Azarm, S.: A kriging metamodel assisted multi-objective genetic algo-  
1155 rithm for design optimization. Journal of Mechanical Design **130**(3), 031401 (2008)
- 1156 26. Li, X., Gong, C., Gu, L., Jing, Z., Fang, H., Gao, R.: A reliability-based optimization  
1157 method using sequential surrogate model and Monte Carlo simulation. Structural and  
1158 Multidisciplinary Optimization pp. 1–22 (2018)
- 1159 27. Lin, Y., Zhang, H.H.: Component selection and smoothing in multivariate nonparamet-  
1160 ric regression. The Annals of Statistics **34**(5), 2272–2297 (2006)
- 1161 28. Lin, Y., Zhang, H.H.: Component selection and smoothing in smoothing spline analysis  
1162 of variance models. Annals of Statistics **34**(5), 2272–2297 (2006)
- 1163 29. Liu, J., Song, W.P., Han, Z.H., Zhang, Y.: Efficient aerodynamic shape optimization of  
1164 transonic wings using a parallel infilling strategy and surrogate models. Structural and  
1165 Multidisciplinary Optimization **55**(3), 925–943 (2017)
- 1166 30. Martins, J.R., Lambe, A.B.: Multidisciplinary design optimization: a survey of archi-  
1167 tectures. AIAA journal **51**(9), 2049–2075 (2013)
- 1168 31. Meza, L.R., Das, S., Greer, J.R.: Strong, lightweight, and recoverable three-dimensional  
1169 ceramic nanolattices. Science **345**(6202), 1322–1326 (2014)
- 1170 32. Mockus, J.: On Bayesian methods for seeking the extremum. In: Optimization Techn-  
1171iques IFIP Technical Conference, pp. 400–404. Springer (1975)
- 1172 33. Mockus, J.: The Bayesian approach to global optimization. System Modeling and Op-  
1173 timization pp. 473–481 (1982)
- 1174 34. Müller, J.: MISO: mixed-integer surrogate optimization framework. Optimization and  
1175 Engineering **17**(1), 177–203 (2016)
- 1176 35. Müller, J., Shoemaker, C.A., Piché, R.: SO-MI: A surrogate model algorithm for com-  
1177 putationally expensive nonlinear mixed-integer black-box global optimization problems.  
1178 Computers & Operations Research **40**(5), 1383–1400 (2013)
- 1179 36. Müller, J., Shoemaker, C.A., Piché, R.: SO-I: a surrogate model algorithm for expensive  
1180 nonlinear integer programming problems including global optimization applications.  
1181 Journal of Global Optimization **59**(4), 865–889 (2014)
- 1182 37. Nguyen-Tuong, D., Peters, J.: Local Gaussian process regression for real-time model-  
1183 based robot control. In: Intelligent Robots and Systems, 2008. IROS 2008. IEEE/RSJ  
1184 International Conference on, pp. 380–385. IEEE (2008)
- 1185 38. Nguyen-Tuong, D., Peters, J.R., Seeger, M.: Local Gaussian process regression for real  
1186 time online model learning. In: Advances in Neural Information Processing Systems,  
1187 pp. 1193–1200 (2009)
- 1188 39. Nguyen-Tuong, D., Seeger, M., Peters, J.: Model learning with local Gaussian process  
1189 regression. Advanced Robotics **23**(15), 2015–2034 (2009)
- 1190 40. Nguyen-Tuong, D., Seeger, M., Peters, J.: Real-time local Gaussian process model learn-  
1191 ing. In: From Motor Learning to Interaction Learning in Robots, pp. 193–207. Springer  
1192 Berlin Heidelberg (2010)
- 1193 41. Nielsen, H.B., Lophaven, S.N., Søndergaard, J.: DACE, a MATLAB Kriging toolbox,  
1194 vol. 2. Citeseer (2002)
- 1195 42. Qian, P.Z.G., Wu, H., Wu, C.J.: Gaussian process models for computer experiments  
1196 with qualitative and quantitative factors. Technometrics **50**(3), 383–396 (2008)
- 1197 43. Queipo, N.V., Haftka, R.T., Shyy, W., Goel, T., Vaidyanathan, R., Tucker, P.K.:  
1198 Surrogate-based analysis and optimization. Progress in Aerospace Sciences **41**(1), 1–28  
1199 (2005)
- 1200 44. Rao, S.S.: Engineering optimization: theory and practice. John Wiley & Sons (2009)
- 1201 45. Ravindran, A., Reklaitis, G.V., Ragsdell, K.M.: Engineering optimization: methods and  
1202 applications. John Wiley & Sons (2006)

- 1203 46. Rehman, S.u., Langelaar, M.: Expected improvement based infill sampling for global  
1204 robust optimization of constrained problems. *Optimization and Engineering* **18**(3),  
1205 723–753 (2017)
- 1206 47. Shahriari, B., Swersky, K., Wang, Z., Adams, R.P., de Freitas, N.: Taking the human  
1207 out of the loop: A review of Bayesian optimization. *Proceedings of the IEEE* **104**(1),  
1208 148–175 (2016)
- 1209 48. Shahzad, M., Kamran, A., Siddiqui, M.Z., Farhan, M.: Mechanical characterization and  
1210 FE modelling of a hyperelastic material. *Materials Research* **18**(5), 918–924 (2015)
- 1211 49. Simpson, T.W., Mauery, T.M., Korte, J.J., Mistree, F.: Kriging models for global ap-  
1212 proximation in simulation-based multidisciplinary design optimization. *AIAA journal*  
1213 **39**(12), 2233–2241 (2001)
- 1214 50. Sóbester, A., Forrester, A.I., Toal, D.J., Tresidder, E., Tucker, S.: Engineering design ap-  
1215 plications of surrogate-assisted optimization techniques. *Optimization and Engineering*  
1216 **15**(1), 243–265 (2014)
- 1217 51. Song, C., Song, W., Yang, X.: Gradient-enhanced hierarchical kriging model for aerody-  
1218 namic design optimization. *Journal of Aerospace Engineering* **30**(6), 04017072 (2017)
- 1219 52. Srinivas, N., Krause, A., Kakade, S.M., Seeger, M.: Gaussian process optimization in  
1220 the bandit setting: No regret and experimental design. *arXiv preprint arXiv:0912.3995*  
1221 (2009)
- 1222 53. Srinivas, N., Krause, A., Kakade, S.M., Seeger, M.W.: Information-theoretic regret  
1223 bounds for Gaussian process optimization in the bandit setting. *IEEE Transactions*  
1224 *on Information Theory* **58**(5), 3250–3265 (2012)
- 1225 54. van Stein, B., Wang, H., Kowalczyk, W., Bäck, T., Emmerich, M.: Optimally weighted  
1226 cluster kriging for big data regression. In: *International Symposium on Intelligent Data*  
1227 *Analysis*, pp. 310–321. Springer (2015)
- 1228 55. Storlie, C.B., Bondell, H.D., Reich, B.J., Zhang, H.H.: Surface estimation, variable se-  
1229 lection, and the nonparametric oracle property. *Statistica Sinica* **21**(2), 679 (2011)
- 1230 56. Swiler, L.P., Hough, P.D., Qian, P., Xu, X., Storlie, C., Lee, H.: Surrogate models for  
1231 mixed discrete-continuous variables. In: *Constraint Programming and Decision Making*,  
1232 pp. 181–202. Springer (2014)
- 1233 57. Tang, Y., Yin, J.: Design of cut unit geometry in hierarchical kirigami-based auxetic  
1234 metamaterials for high stretchability and compressibility. *Extreme Mechanics Letters*  
1235 **12**, 77–85 (2017)
- 1236 58. Tran, A., He, L., Wang, Y.: An efficient first-principles saddle point searching method  
1237 based on distributed kriging metamodels. *ASCE-ASME Journal of Risk and Uncertainty*  
1238 *in Engineering Systems, Part B: Mechanical Engineering* **4**(1), 011006 (2018)
- 1239 59. Viana, F.A., Simpson, T.W., Balabanov, V., Toropov, V.: Special section on multidis-  
1240 ciplinary design optimization: Metamodeling in multidisciplinary design optimization:  
1241 How far have we really come? *AIAA Journal* **52**(4), 670–690 (2014)
- 1242 60. Zhang, Y., Hu, S., Wu, J., Zhang, Y., Chen, L.: Multi-objective optimization of double  
1243 suction centrifugal pump using kriging metamodels. *Advances in Engineering Software*  
1244 **74**, 16–26 (2014)
- 1245 61. Zhou, Q., Qian, P.Z., Zhou, S.: A simple approach to emulation for computer models  
1246 with qualitative and quantitative factors. *Technometrics* **53**(3), 266–273 (2011)
- 1247 62. Zhou, Q., Wang, Y., Choi, S.K., Jiang, P., Shao, X., Hu, J.: A sequential multi-fidelity  
1248 metamodeling approach for data regression. *Knowledge-Based Systems* (2017)
- 1249 63. Zhou, Q., Wang, Y., Choi, S.K., Jiang, P., Shao, X., Hu, J., Shu, L.: A robust opti-  
1250 mization approach based on multi-fidelity metamodel. *Structural and Multidisciplinary*  
1251 *Optimization* pp. 1–23 (2017)