
Accelerating the Discovery of Rare Materials with Bounded Optimization Techniques

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Discovering a rare material within a vast search space exhibits a Needle-in-a-
2 Haystack challenge. This challenge of finding a rare material, *i.e.*, the "needle",
3 inside a vast search space, *i.e.*, the "haystack", arises when there is an extreme
4 imbalance of optimum conditions relative to the size of the search space. For ex-
5 ample, only 0.82% out of 146k total materials in the open-access Materials Project
6 database have a negative Poisson's ratio, a rare material property. However, current
7 state-of-the-art optimization algorithms are not designed with the capabilities to
8 find solutions to these challenging multidimensional Needle-in-a-Haystack prob-
9 lems, resulting in slow convergence to a global optimum or pigeonholing into a
10 local minimum. In this paper, we present a Zooming Memory-Based Initialization
11 algorithm, entitled ZoMBI, that builds on conventional Bayesian optimization
12 principles to quickly and efficiently optimize Needle-in-a-Haystack problems in
13 both less time and fewer experiments by addressing the common convergence
14 and pigeonholing issues. ZoMBI actively extracts knowledge from the previously
15 best performing evaluated experiments to iteratively zoom in the sampling search
16 bounds towards the global optimum "needle" and then prunes the memory of
17 low-performing historical experiments to accelerate compute times by reducing
18 the algorithm time complexity from $O(n^3)$ to $O(1)$, as the number of experiments
19 sampled increases. Additionally, ZoMBI implements two custom acquisition func-
20 tions that use active learning to further guide the sampling of new experiments
21 towards the global optimum. We validate the algorithm's performance on two
22 real-world 5-dimensional Needle-in-a-Haystack material property optimization
23 datasets: discovery of auxetic Poisson's ratio materials and discovery of high ther-
24 moelectric figure of merit materials. The ZoMBI algorithm demonstrates compute
25 time speed-ups of 400x compared to traditional Bayesian optimization as well as
26 efficiently discovering materials in under 100 experiments that are up to 3x more
27 highly optimized than those discovered by current state-of-the-art algorithms.

28 1 Introduction to Rare Material Discovery

29 Current optimization algorithms perform well on low-dimensional problems that are smooth and have
30 wide basins of attraction. Examples of smooth manifolds with wide basins of attraction within material
31 science include process- and recipe-optimization problems such as tuning perovskite manufacturing
32 variables to achieve higher efficiency [1], optimizing microfluidics flow parameters to achieve ideal
33 droplet formation [2], optimizing silver nanoparticle recipes for optical properties [3], and tuning
34 perovskite compositions with physics-based constraints to maximize stability [4]. Optimization
35 techniques like Bayesian optimization (BO) are well-suited to model these simple manifolds using
36 a Gaussian Process (GP) surrogate [5, 6, 7, 8, 9]. However, the performance of this BO with a GP
37 breaks down as the manifold complexity increases. Material property optimization problems that

38 have high technological significance, such as discovering materials with rare properties or materials
 39 with a specific combination of properties, have search space manifolds that more closely resemble a
 40 *Needle-in-a-Haystack* [10], shown in Figure 1(b), rather than a smooth or convex space. This Needle-
 41 in-a-Haystack (NiaH) problem arises when only few optimum conditions exist within the entire
 42 search space, resulting in an extreme imbalance. Interpolating the parameter space of an imbalanced
 43 search space with an estimation function, such as a GP, results in smoothing over the optimum or
 44 over-predicting the properties of the materials found near the optimum [11, 12, 13]. Examples of
 45 NiaH materials optimization problems include discovering auxetic materials (*i.e.*, materials that
 46 have a highly negative Poisson’s ratio, ν) for energy absorptive medical devices or protective armor
 47 [14, 15, 16] and discovering materials that have a combination of high electrical conductivity and low
 48 thermal conductivity (*i.e.*, a highly positive thermoelectric figure of merit, ZT) used from improving
 49 sensor technology to enable ubiquitous solid-state cooling [17, 18, 19]. Both of these rare material
 50 optimization problems are examples where an extreme data balance exists in the search space because
 51 only a fraction of the total number of materials exhibit these rare properties [14, 20, 21, 22, 23]. This
 52 NiaH optimization challenge of extremely imbalanced search spaces is largely applicable to many
 53 fields, not just materials science, including the fields of ecological resource management [24], fraud
 54 detection [25, 26], and rare diseases [27, 26].

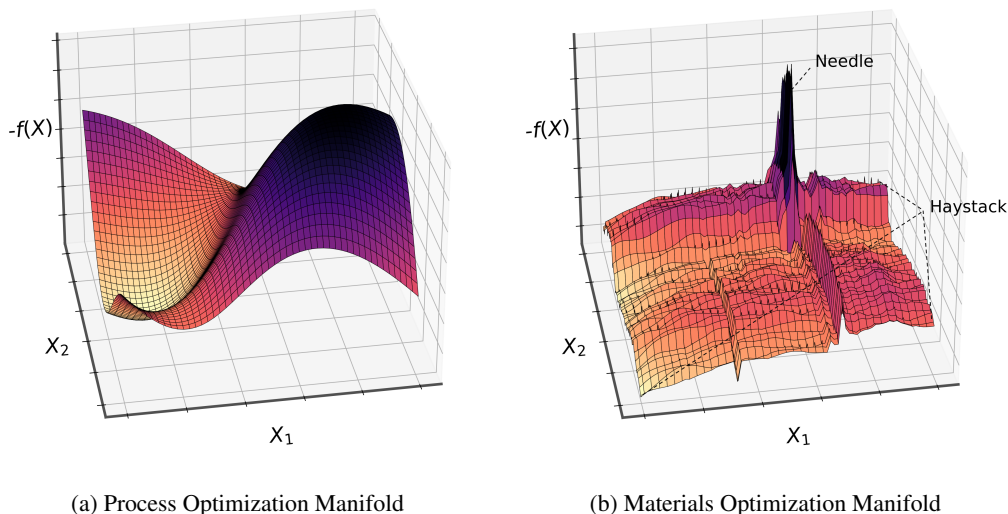


Figure 1: **Archetype Manifolds in Materials Science Optimization.** (a) Smooth and wide basin of attraction landscape that is common to process optimization problems. This 2D projected manifold is adapted from the 6D perovskite process optimization problem by Liu *et al.* [1]. (b) Rough and narrow basin of attraction landscape that is typical of material property optimization problems. This 2D projected manifold is obtained from the 5D negative Poisson’s ratio optimization problem presented in this paper [20, 21].

55 Several challenges exist for the current landscape of computational tools that inhibit effective
 56 optimization of these complex NiaH problems. Firstly, the "needle" makes up only a small percentage
 57 of the total manifold search space, resulting in a weak correlation between the measured input
 58 parameters and the target property of interest, inhibiting discovery of the region containing the needle
 59 [28, 29, 11]. This challenge requires the development of an algorithm that can more quickly determine
 60 the plausible region of the manifold where the needle exists. The second challenge for algorithms,
 61 such as BO, to optimize NiaH manifolds is in the nature of the acquisition function to pigeonhole
 62 sampling into local minima because of the narrowness of the needle’s basin of attraction [30, 31].
 63 Standard BO acquisition functions, including expected improvement (EI) [32] and lower confidence
 64 bound (LCB) [7, 12], are static sampling techniques that only adjust sampling based on the output of
 65 the surrogate model, which enacts smoothing of the needle [11, 5, 6]. To overcome this challenge,
 66 active learning-based tuning of the acquisition function hyperparameters can be implemented to
 67 improve the sampling quality and avoid pigeonholing. We design two active learning acquisition
 68 functions, LCB Adaptive and EI Abrupt, further discussed in the Appendix (sections A.1 & A.2).

69 Lastly, there exists a computing challenge for NiaH problems where, typically, several thousands of
70 samples must be observed to find an optimum when using an algorithm that is poorly-suited to tackle
71 NiaH manifolds [10]. The compute time of BO using a GP surrogate scales with the complexity
72 $O(n^3)$, where n is the number of experiments sampled, hence, the compute time of traditional BO
73 blows up as more data is required to find the optimum [33, 34, 35, 5, 6, 36, 37]. To solve this
74 computing challenge, an algorithm must be designed that both efficiently optimizes the space in as
75 few experiments as possible and reduces the effect of compounding compute times over the length of
76 the optimization procedure.

77 1.1 Related Literature on Fast and Bounded Optimization

78 In recent literature, algorithms have been developed to address some of these challenges individually,
79 but not all of them together. The first class of solutions bound the search space using a trust region
80 approach to sample regions with higher probability of containing the optimum. Uber AI develop
81 TuRBO [38] that compiles a set of independent model runs, using separate GP surrogate models to
82 compute a new, smaller search region, narrowed in on the target optimum. Regis develops TRIKE
83 [39] that utilizes maximization of the EI acquisition function to bound a trust region containing the
84 global optimum. Diouane *et al.* develop TREGO [40], which interleaves sampling between global
85 and local search regions, where the local search regions are defined by the single best historical
86 experiment sampled. Although these methods offer solutions to one of the three challenges presented,
87 each method has its downfalls when optimizing NiaH problems. For example, TuRBO requires the
88 computation of several GP model runs, which increases compute time and also does not guarantee
89 that the needle will be resolved due to interpolation effects; TRIKE is inflexible to the use of other
90 acquisition functions as it locks the user in to only using EI, which may pigeonhole into local
91 minima; TREGO uses only the best sampled experiment to define its search regions, which will yield
92 inconsistent or sub-optimal results when the needle consists of a fractional region of the manifold
93 and single point is unlikely to land in its basin of attraction.

94 The second class of solutions to the challenges presented in this paper are designed to decrease the
95 computing time required to run an optimization procedure. A common method for reducing the
96 compute time of BO with a GP surrogate is to introduce a sparse GP [5, 41, 36]. A sparse GP uses a
97 small subset of pseudo data, often denoted as m , to reduce the GP time complexity from $O(n^3)$ to
98 $O(nm^2)$ [42]. However, the process of selecting a useful subset requires minimizing the Kullback-
99 Leibler divergence between the sparse GP and true posterior GP, which is often a computationally
100 intensive procedure of using variational inference [43]. In addition to sparse GPs, new algorithms
101 have been developed in literature to improve the compute time of optimization in various ways.
102 Van Stein *et al.* develop MiP-EGO [44], which parallelizes the function evaluations of efficient
103 global optimization (EGO) to discover optima faster and in fewer experiments using derivative-free
104 computation [45]. Joy *et al.* [46] use directional derivatives to accelerate hyperparameter tuning
105 by 100x and achieve higher accuracy than the FABOLAS baseline by Klein *et al.* [47]. Zhang *et al.*
106 develop FLASH [48] to achieve optimization speed-ups of 50% by using a linear parametric
107 model to guide algorithm search within high-dimensional spaces. Snoek *et al.* [13] design a neural
108 network-based parametric model that reduces the overall time complexity of BO to $O(n)$ compared
109 to the complexity of $O(n^3)$ of standard BO with a GP surrogate model. These existing methods
110 from literature within the class of solutions for accelerating compute time are generally introducing
111 external models necessary to perform optimization, such as neural networks, variational inference,
112 or parameteric models. While these external models do speed-up compute time, they often lack the
113 predictive capabilities to capture the weak correlation between measured input parameters and the
114 target property of interest in NiaH problems. We illustrate this mechanic later in the paper when
115 comparing the optimization results on two materials science NiaH problems of a fast algorithm
116 MiP-EGO with that of TuRBO, an algorithm better suited for discovering optima within narrow basins
117 of attraction.

118 Although these methods from existing literature address some of the challenges in optimizing NiaH
119 problems, none of them have been designed specifically to quickly and efficiently discover a needle-
120 like optimum within a haystack of sub-optimal points, resulting in all of them falling short of full
121 solution. Therefore, in this paper, we design an algorithm that addresses all three of the challenges
122 faced when optimizing NiaH problems by (1) zooming in the manifold search bounds iteratively and
123 independently for each dimension based on m number of best memory points to quickly converge to
124 the plausible region containing the global optimum needle, (2) anti-pigeonholing into local minima

125 by using actively learned acquisition function hyperparameters to tune the exploitation-to-exploration
 126 ratio, (3) relieving compute utilization by pruning the low-performing memory points not being
 127 used to zoom in the search bounds. The proposed algorithm, entitled [Zo]oming [M]emory-[B]ased
 128 [I]nitialization (ZoMBI), combines these three contributions into a method that efficiently optimizes
 129 NiaH problems quickly. In essence, this process of scanning broadly and then focusing in on points
 130 of interest based on memory was inspired by the way we humans solve similar problems, but stands
 131 in contrast to the way standard BO methods with static acquisition functions solve problems. We
 132 demonstrate the performance of this algorithm on two NiaH materials science datasets: (1) discovery
 133 of materials with negative Poisson’s ratio and (2) discovery of materials with both high electrical
 134 conductivity and low thermal conductivity. The performance of the proposed ZoMBI algorithm is
 135 compared against standard BO with static acquisition functions and two state-of-the-art (SoTA)
 136 algorithms, one from each of the two classes of partial NiaH solutions: (1) TuRBO (bounded search
 137 space) and (2) MiP-EGO (faster compute).

138 2 Methodology: Bounded & Memory-Pruning Optimization

139 The [Zo]oming [M]emory-[B]ased [I]nitialization ZoMBI algorithm has two key features: (1) iterative
 140 inward bounding of proceeding search spaces using the m number of best-performing memory
 141 points within the prior search space and (2) iterative pruning of low-performing historical search
 142 space memory. The newly computed search space bounds are unique for each dimension, such
 143 that optimum basin of attraction of complex, non-convex NiaH manifolds can be discovered. This
 144 algorithm leverages these two key features to guide the acquisition of new data towards more optimal
 145 regions while only fitting the surrogate within the suggested optimum region to resolve more detail
 146 of the space of interest, as shown in Figure 2. This process subsequently reduces the compute time
 147 significantly compared to the compute of a GP in a standard BO procedure, as shown in Figure 3.

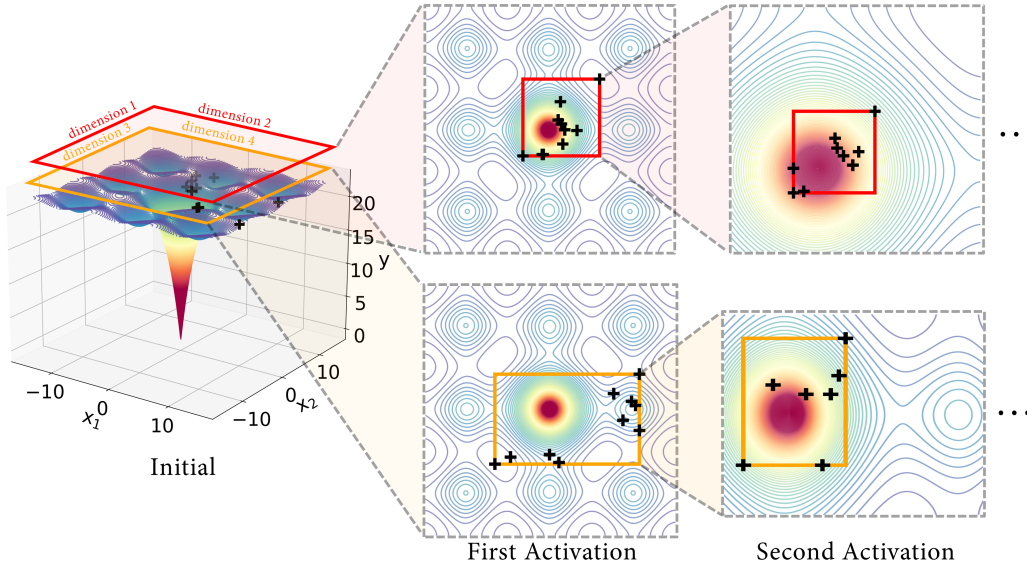


Figure 2: **Zooming Search Bounds.** For every activation of ZoMBI, the search bounds are zoomed inward based on the prior best-performing memory points. A 4D Ackley function manifold is projected in 2D. The bounding regions of each 2D slice are illustrate by the red and orange boxes. The ϕ number forward experiments sampled are illustrated as black markers. The global optimum is indicated by the red region of the heatmap.

148 We define m as the number of retained memory points during an activation of ZoMBI. The m memory
 149 points are saved to memory while all other data are erased from memory. These are the historical
 150 data points that achieve the m lowest (for minimization) target values, y , and they are used to zoom
 151 in the search bounds. Using these memory points, the multi-dimensional upper and lower bounds
 152 of the zoomed search space are computed for each dimension, d . Let $\mathbf{X} := \{X_1, X_2, \dots, X_n\}$ be a set
 153 of data points, where $X_j \in \mathbb{R}^d$. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be the objective function. We first assume that the

154 points in \mathbf{X} are in general position so that $f(\mathbf{X})$ contains unique elements. Then, for each $m \leq n$
 155 define $\mathbf{X}^{(m)} = \{X_{\pi(1)}, \dots, X_{\pi(m)}\}$ where π is a permutation on $\{1, \dots, n\}$ so that $\{f(X_{\pi(j)})\}$ is
 156 in ascending order. If $f(\mathbf{X})$ contains repeated elements, we may first remove the points with repeated
 157 f values and apply the definition above. Then, for each d , the bounds are defined as:

$$\mathcal{B}_d^l = \min_{X \in \mathbf{X}^{(m)}} \{[X]_d\}$$

$$\mathcal{B}_d^u = \max_{X \in \mathbf{X}^{(m)}} \{[X]_d\},$$
(1)

158 where \mathcal{B}_d^l and \mathcal{B}_d^u computed lower and upper bounds for each dimension, d , respectively. The bounds
 159 $[\mathcal{B}_d^l, \mathcal{B}_d^u]$ constrain the proceeding acquisition of new data as well as the computation of a GP, such
 160 that sampling cannot occur outside of the bounded region. This constraining process operates
 161 independently for each dimension, such that each dimension has a unique lower and upper bound.
 162 To initialize the algorithm with data from the constrained space, i data points are sampled from the
 163 bounded region using Latin Hypercube Sampling (LHS). LHS splits a d -dimensional space into $i * d$
 164 equally spaced strata, where i is the number of points to sample uniformly over d dimensions with
 165 low variability, unlike random sampling that has high sampling variability [50]. A GP surrogate
 166 model is retrained on these i LHS points sampled from the constrained space and then for every
 167 proceeding experiment sampled from the space, denoted as a forward experiment, the surrogate model
 168 is retrained. Thus, the GP is only being trained on information within the constrained region and as
 169 the constrained region iteratively zooms inward and decreases in hypervolume, so does the region
 170 computed by the GP. This process allows for more information to be resolved within regions plausibly
 171 containing the global optimum basin of attraction. Up to ϕ forward experiments are sampled in serial,
 172 where $\{X_i\} \cup \{X_\phi\} \subseteq \{X_n\}$. These forward experiments are sampled by maximizing an acquisition
 173 value, $a \in [0, 1]$, computed by a user-selected acquisition function from one of the four functions EI,
 174 EI Abrupt, LCB, and LCB Adaptive. Once $i + \phi$ number of experiments are sampled, the bounds are
 175 re-constrained using the m best performing experiments, i new experiments are sampled from the
 176 zoomed-in space using LHS, and then the memory is pruned. The process of collecting ϕ forward

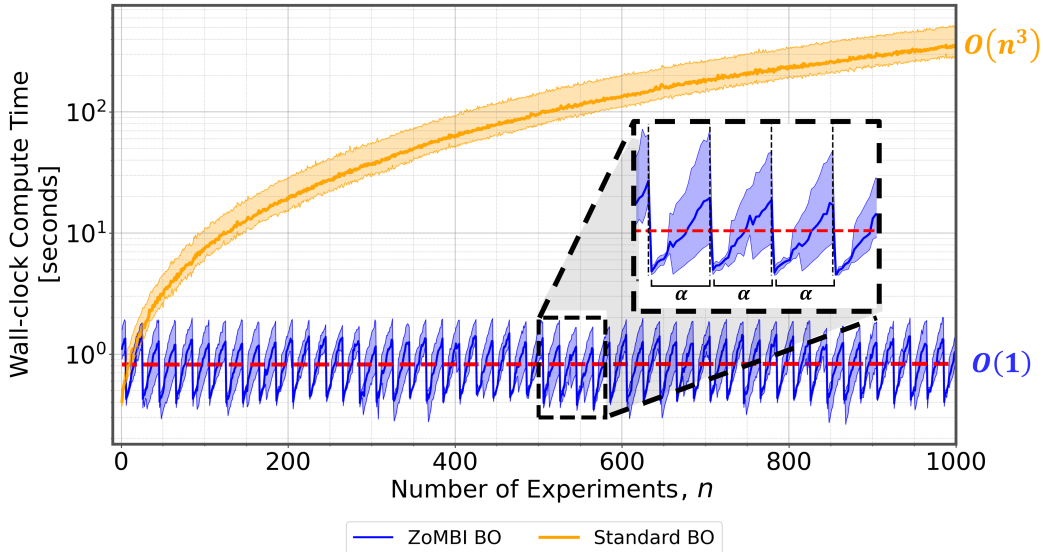


Figure 3: **Wall-clock Compute Time.** The compute time per experiment is illustrate for traditional BO with a GP surrogate (orange) and for ZoMBI with a GP surrogate (blue) with the y -axis in log-scale. Five independent trials of each method were run to optimize a 5D Ackley function with a narrow basin of attraction using an NVIDIA Tesla Volta V100 GPU [49]. The averages of the trials are shown as solid orange and blue lines while the shaded regions indicate the maximum and minimum compute times bounds. The red dashed line indicates the trend of the ZoMBI compute times. The measured compute time includes the time to compute the GP surrogate model and the time to acquire an experiment from the surrogate.

177 experiments is repeated. A complete constraining-resetting iteration is denoted as an activation, α .
178 This iterative zooming and pruning process over several α significantly speeds up compute time.

179 3 Rare Material Discovery Results

180 3.1 Compute Time

181 In this section, we assess the compute time of the developed algorithm in comparison that of standard
182 BO methods. As more experiments are amassed and committed to memory to run traditional BO
183 by computing a GP and an acquisition value, the compute time increases polynomially, following
184 the $O(n^3)$ time complexity of GP matrix multiplication [33, 5, 6, 51, 36, 37]. This complexity is
185 unfavorable as it leads to compounding compute times as more experiments are run. Therefore, we
186 implement a memory pruning feature into the ZoMBI algorithm that iteratively selects which prior
187 data points to keep and which ones to prune from the memory during each activation, α . Via memory
188 pruning, the number of experiments used to train the GP surrogate varies between $[i, i + \phi]$ for every α ,
189 rather than being proportional to n . This is computationally favorable because $\{X_i\} \cup \{X_\phi\} \subseteq \{X_n\}$.
190 Thus, for a single α , the time complexity is $O((i + \phi)^3)$. However, since ϕ resets back to zero after
191 each α , a non-increasing sawtooth pattern in compute time is exhibited, hence, as $\alpha, n \rightarrow \infty$, the
192 complexity approaches $O(1)$. Figure 3 illustrates that the sawtooth compute time pattern maps to
193 the resetting interval of ϕ , which trends towards a constant, non-increasing value over many α and
194 n . After collecting 1000 experiments, the compute time of traditional BO trend towards > 400
195 seconds, whereas after 1000 experiments, the compute time ZoMBI trends towards a constant 1
196 second. Therefore, the memory pruning feature of ZoMBI accelerates the optimization compute time
197 by over 400x at $n = 1000$ and achieves further relative acceleration as n increases. The memory
198 pruning mechanic of ZoMBI drives fast compute times without sacrificing the ability to converge on
199 rare materials, demonstrated in the following sections (3.2 & 3.3).

200 3.2 Poisson’s Ratio

201 We demonstrate the ability of the ZoMBI algorithm to optimize Needle-in-a-Haystack problems on two
202 real-world datasets. The first dataset consists of 146k materials and the objective is to find the material
203 with the minimum negative Poisson’s ratio, ν . The second dataset consists of 1k materials and the
204 objective is to find the material with the maximum thermoelectric merit, ZT , *i.e.*, a material with
205 high electrical conductivity and low thermal conductivity. Both of these datasets are 5-dimensional
206 and are obtained from the open-access Materials Project database [20].

207 The ν dataset exhibits a Needle-in-a-Haystack problem due to very few materials having negative
208 ν values [14, 20, 21, 15]. A positive $\nu > 0$, describes a material that expands when a compressive
209 load is applied to the orthogonal direction [52, 53]. Conversely, a negative $\nu < 0$ describes a material
210 that contracts rather than expands when compressed in the orthogonal direction, denoted as an
211 auxetic material [14, 23] – a rare phenomenon that occurs in only 0.82% of materials within the
212 Materials Project database [20, 21]. Auxetic materials with highly negative Poisson’s ratios have
213 energy absorptive properties, which are ideal materials for wearable medical devices and protective
214 armor that must absorb the energy of large impacts to keep bones from shifting or to inhibit the
215 penetration of the protective layer [15, 16]. Thus, for this NiaH problem, the objective is to discover
216 the material with the lowest ν value. Figure A.2 illustrates the spread of ν values within the raw
217 dataset as a histogram as well as a manifold generated by a Random Forest (RF) regression on the
218 raw dataset using 500 trees. The search space generated by the RF is noisy and non-convex with
219 narrow basins of attraction containing each optimum, resulting in a challenging NiaH optimization
220 problem. The ground truth "needle" materials with the lowest ν values are Li_2NbF_6 with $\nu \approx -1.7$
221 and Na_2CO_3 with $\nu \approx -1.2$.

222 Figure 4 illustrates the performance of ZoMBI in discovering the lowest ν -value material, compared
223 to the SoTA TuRBO and MiP-EGO algorithms. The ZoMBI algorithm is run with each of the four
224 acquisition functions: LCB, LCB Adaptive, EI, and EI Abrupt. In under 100 evaluated experiments,
225 LCB and LCB Adaptive discover one of the needles within the dataset (Li_2NbF_6) and, similarly, EI
226 Abrupt discovers the other needle (Na_2CO_3). The distribution of ν values for the final experiment
227 across all ensemble runs is illustrated for each method to highlight the sampling density and general
228 rate of success. LCB Adaptive and EI Abrupt are the first two implementations of ZoMBI to discover
229 a $\nu < 0$ material because of their ability to actively tune their sampling hyperparameters. After

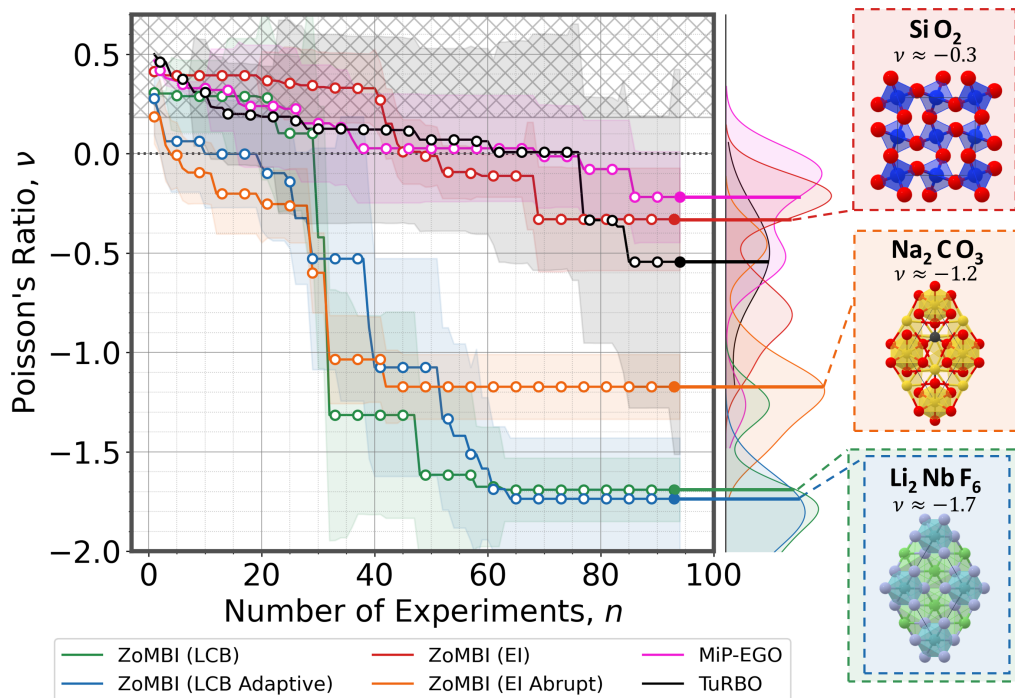


Figure 4: **Discovery of Rare Negative Poisson’s Ratio Materials.** The optimization objective is to find the material with the minimum Poisson’s ratio, ν_{\min} , in 100 experiments. The green, blue, red, and orange lines indicate the median best running evaluated sample of ZoMBI using the LCB, LCB Adaptive, EI, and EI Abrupt acquisition functions, respectively. The pink and black lines indicate the median best running evaluated sample of the SoTA methods, MiP-EGO and TuRBO, respectively. The median for each method is taken over the best 12 independent model runs. The shaded regions indicate the variance between model runs. The crosshatched region indicates the space discovered by standard BO methods, without the use of ZoMBI. The dashed black line indicates the $\nu = 0$ inflection point. The distribution of the final sampled ν value for each method at the 100th experiment is shown as a kernel density estimation with a 0.5 smoothing factor. The materials formulae and unit cells that have the closest evaluated ν value discovered by each ZoMBI method at the end of the 100 experiments are illustrated.

230 30 experiments, the ZoMBI search bounds have zoomed inward enough for the explorative LCB
 231 acquisition function to discover a region of the manifold containing highly negative ν material,
 232 eventually leading to the global minimum needle. These three implementations of ZoMBI: LCB,
 233 LCB Adaptive, and EI Abrupt, have a steep drop in the discovered ν value, allowing these methods
 234 to discover an optimum fast, in fewer experiments than both SoTA methods. Overall, LCB and
 235 LCB Adaptive implementations of ZoMBI discover the most optimum minimum $\nu \approx -1.7$, while
 236 the SoTA algorithms TuRBO and MiP-EGO only discover $\nu \approx -0.55$ and $\nu \approx -0.20$, respectively.
 237 These results demonstrate that with proper selection an acquisition function, ZoMBI achieves better
 238 performance and a higher success rate than SoTA on optimizing this real-world materials science
 239 NiaH problem.

240 3.3 Thermoelectrics

241 The ZT dataset exhibits a Needle-in-a-Haystack problem, similar to the ν dataset because very
 242 few materials have high ZT values [20, 10]. However, rather than ZT being a directly measurable
 243 mechanical material property like Poisson’s ratio, ZT must be computed using a combination of
 244 several thermal and electrical material properties [54]:

$$ZT = \frac{S^2 \sigma}{\kappa} T, \quad (2)$$

245 where S is the Seebeck coefficient, σ is electrical conductivity, T is the average temperature, and
 246 κ is thermal conductivity. The ZT is computed for each material in the Materials Project database
 247 using BoltzTraP [55]. Of the initial 146k materials, 1k of them have the required thermal and
 248 electrical properties to compute a ZT value. ZT is a common figure of merit used to describe
 249 the thermal-to-electrical or electrical-to-thermal conversion efficiency of thermoelectric materials
 250 [56, 57, 58, 59]. A higher ZT indicates that the material is better able to convert a thermal gradient
 251 into an electrical current [54]. Materials with large ZT values have a range of applications from
 252 usage as solid-state cooling devices to being used as sensors that when heated up, will produce an
 253 electrical signal [17, 18, 19]. For this NiaH problem, the objective is to discover the material with
 254 the highest ZT value. Figure A.3 illustrates the spread of ZT values within the raw dataset as a
 255 histogram, as well as a manifold generated by an RF regression on the raw dataset using 500 trees.
 256 Similar to the ν manifold, the ZT manifold is noisy and non-convex with narrow basins of attraction
 257 [20, 55]. The ground truth "needle" materials with the highest ZT values are $\text{Na}_4\text{Al}_3\text{Ge}_3\text{IO}_{12}$ with
 258 $ZT \approx 1.4$ and $\text{Sr}_4\text{Al}_6\text{SO}_{12}$ with $ZT \approx 1.9$.

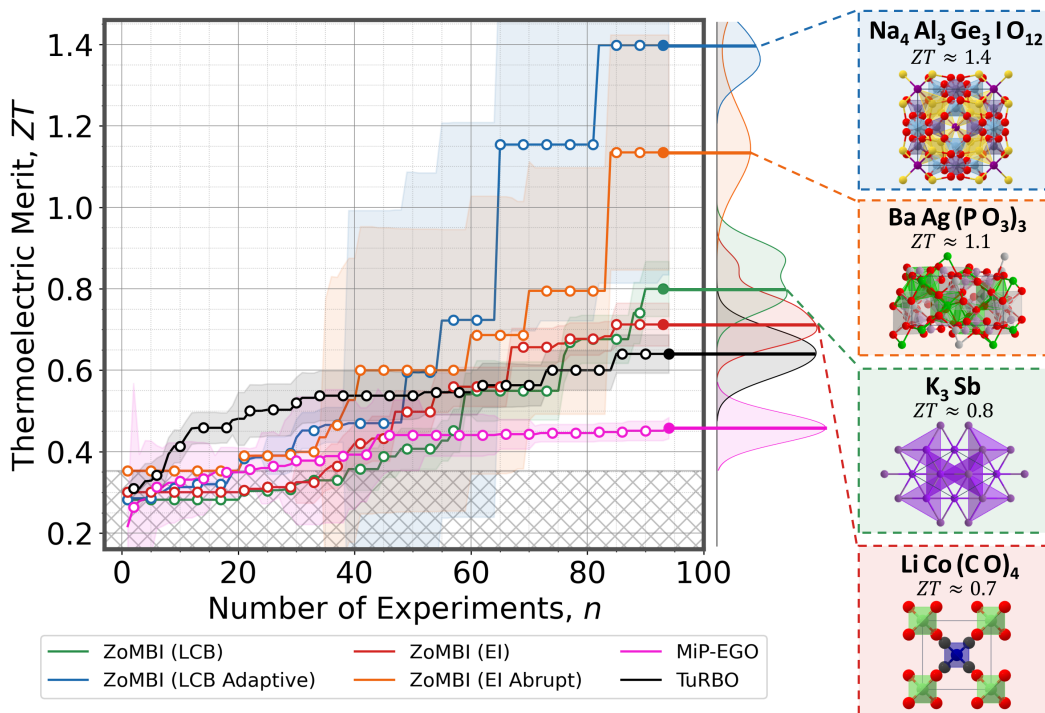


Figure 5: **Discovery of Rare Positive Thermoelectric Merit Materials.** The optimization objective is to find the material with the maximum thermoelectric merit, ZT_{\max} , in 100 experiments. The green, blue, red, and orange lines indicate the median best running evaluated sample of ZoMBI using the LCB, LCB Adaptive, EI, and EI Abrupt acquisition functions, respectively. The pink and black lines indicate the median best running evaluated sample of the SoTA methods, MiP-EGO and TuRBO, respectively. The median for each method is taken over the best 12 independent model runs. The shaded regions indicate the variance between model runs. The crosshatched region indicates the space discovered by standard BO methods, without the use of ZoMBI. The distribution of the final sampled ZT value for each method at the 100th experiment is shown as a kernel density estimation with a 0.5 smoothing factor. The materials formulae and unit cells that have the closest evaluated ZT value discovered by each ZoMBI method at the end of the 100 experiments are illustrated.

259 Figure 5 illustrates the performance of ZoMBI in discovering the highest ZT -value material, compared
 260 to the SoTA TuRBO and MiP-EGO algorithms. Initially, we see TuRBO outperform all other algorithms,
 261 but then it is unable to accelerate its sampling towards the needle basins of attraction. Similarly,

262 MiP-EGO gets trapped in a local minimum and is unable to escape. Conversely, after 50 evaluated
263 experiments, ZoMBI LCB Adaptive and EI Abrupt supersede TuRBO and quickly discover high ZT
264 materials, illustrating the advantage of active learning acquisition functions. Although the active
265 learning acquisition functions prove to be more successful than the SoTA algorithms, none of the
266 tested algorithms are able to discover the maximum global needle, $\text{Sr}_4\text{Al}_6\text{SO}_{12}$, only the second best
267 needle, $\text{Na}_4\text{Al}_3\text{Ge}_3\text{IO}_{12}$. This result is likely due to the data imbalance being too extreme that far out
268 on the tail of the ZT dataset, in turn, generating an RF manifold complexity too high, even for ZoMBI.
269 Hence, indicating that there are limitations in the manifold complexity that ZoMBI can optimize, and
270 further illustrating that convergence on the global optimum needle is not guaranteed using this method.
271 However, for the ZT dataset, the LCB Adaptive implementation of ZoMBI discovers the second
272 best needle, $\text{Na}_4\text{Al}_3\text{Ge}_3\text{IO}_{12}$ with $ZT \approx 1.4$, while the SoTA algorithms TuRBO and MiP-EGO only
273 discover $ZT \approx 0.65$ and $ZT \approx 0.45$, respectively. Thus, LCB Adaptive demonstrates the highest
274 performing optimization results across both of the real-world NiaH datasets, discovering the most
275 optimal materials the fastest for both the ν and ZT datasets.

276 4 Summary & Conclusion

277 In this paper, we proposed the [Zo]oming [M]emory-[B]ased [I]nitialization (ZoMBI) algorithm
278 that builds on the principles of Bayesian optimization to accelerate the discovery of rare materials
279 by two-fold, firstly by requiring fewer experiments to achieve a better optimum than state-of-the-
280 art, and secondly by pruning the memory of low-performing historical experiments to speed-up
281 compute time. The ZoMBI algorithm exceeds state-of-the-art performance on optimizing Needle-in-
282 a-Haystack datasets by (1) using the values of the m best performing previously sampled memory
283 points to iteratively zoom in the search bounds of the manifold uniquely on each dimension and
284 (2) implementing two custom acquisition functions, LCB Adaptive and EI Abrupt, that actively
285 learn information about the manifold during optimization to tune the sampling of new experimental
286 conditions from a surrogate. The main contributions of this algorithm solve three fundamental
287 challenges of optimizing non-convex Needle-in-a-Haystack problems: (1) the challenge of locating
288 the hypervolume region of the manifold containing the narrow global optimum basin of attraction
289 [28, 29, 11] is alleviated by introducing iterative search bounds based on learned knowledge of the
290 manifold; (2) unwanted pigeonholing into local minima [30, 31, 5, 6] is avoided by both the zooming
291 mechanics of ZoMBI as well as using the two acquisition functions developed in his paper, LCB
292 Adaptive and EI Abrupt, that tune their hyperparameters through active learning; (3) the challenge
293 of polynomially increasing compute times of BO using a GP surrogate [33, 34, 35, 5, 6, 36, 37] is
294 addressed by actively pruning the retained memory of the algorithm after each activation, α , in turn,
295 reducing the time complexity from $O(n^3)$ to $O(1)$ as $\alpha, n \rightarrow \infty$. By developing the ZoMBI algorithm
296 to solve these challenges, it becomes possible to quickly and efficiently find optimal solutions to
297 complex Needle-in-a-Haystack problems in fewer experiments. Hence, this tool can be applied to rare
298 material discovery, a class of data imbalanced Needle-in-a-Haystack problems, to enable widespread
299 discovery of new materials with important technical applications from designing high-performance
300 medical devices to engineering ubiquitous solid-state cooling systems.

301 References

- 302 [1] Zhe Liu, Nicholas Rolston, Austin C. Flick, Thomas W. Colburn, Zekun Ren, Reinhold H.
303 Dauskardt, and Tonio Buonassisi. Machine learning with knowledge constraints for process
304 optimization of open-air perovskite solar cell manufacturing. *Joule*, 6(4):834–849, 2022.
- 305 [2] Alexander E. Siemenn, Evyatar Shaulsky, Matthew Beveridge, Tonio Buonassisi, Sara M.
306 Hashmi, and Iddo Drori. A Machine Learning and Computer Vision Approach to Rapidly
307 Optimize Multiscale Droplet Generation. *ACS Applied Materials & Interfaces*, 14(3):4668–
308 4679, 2022.
- 309 [3] Flore Mekki-Berrada, Zekun Ren, Tan Huang, Wai Kuan Wong, Fang Zheng, Jiaxun Xie, Isaac
310 Parker Siyu Tian, Senthilnath Jayavelu, Zackaria Mahfoud, Daniil Bash, Kedar Hippalgaonkar,
311 Saif Khan, Tonio Buonassisi, Qianxiao Li, and Xiaonan Wang. Two-step machine learning
312 enables optimized nanoparticle synthesis. *npj Computational Materials* 2021 7:1, 7(1):1–10,
313 2021.

- 314 [4] Shijing Sun, Armi Tiihonen, Felipe Oviedo, Zhe Liu, Janak Thapa, Yicheng Zhao, Noor Titan P.
315 Hartono, Anuj Goyal, Thomas Heumueller, Clio Batali, Alex Encinas, Jason J. Yoo, Ruipeng
316 Li, Zekun Ren, I. Marius Peters, Christoph J. Brabec, Mounqi G. Bawendi, Vladan Stevanovic,
317 John Fisher, and Tonio Buonassisi. A data fusion approach to optimize compositional stability
318 of halide perovskites. *Matter*, 4(4):1305–1322, 2021.
- 319 [5] Edward Snelson and Zoubin Ghahramani. Sparse Gaussian Processes using Pseudo-inputs. In
320 Y Weiss, B Schölkopf, and J Platt, editors, *Advances in Neural Information Processing Systems*,
321 volume 18. MIT Press, 2005.
- 322 [6] Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine*
323 *Learning*. The MIT Press, 2005.
- 324 [7] Eric Brochu, Vlad M. Cora, and Nando de Freitas. A Tutorial on Bayesian Optimization
325 of Expensive Cost Functions, with Application to Active User Modeling and Hierarchical
326 Reinforcement Learning. 2010.
- 327 [8] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical Bayesian Optimization Of
328 Machine Learning Algorithms. pages 1–12, 2001.
- 329 [9] Qiaohao Liang, Aldair E. Gongora, Zekun Ren, Armi Tiihonen, Zhe Liu, Shijing Sun, James R.
330 Deneault, Daniil Bash, Flore Mekki-Berrada, Saif A. Khan, Kedar Hippalgaonkar, Benji
331 Maruyama, Keith A. Brown, John Fisher, and Tonio Buonassisi. Benchmarking the perfor-
332 mance of Bayesian optimization across multiple experimental materials science domains. *npj*
333 *Computational Materials* 2021 7:1, 7(1):1–10, 2021.
- 334 [10] Yoolhee Kim, Edward Kim, Erin Antono, Bryce Meredig, and Julia Ling. Machine-learned
335 metrics for predicting the likelihood of success in materials discovery. *npj Computational*
336 *Materials*, 6(131), 2020.
- 337 [11] Ioan Andricioaei and John E Straub. Finding the needle in the haystack: Algorithms for
338 conformational optimization. *Computers in Physics*, 10:449, 1996.
- 339 [12] Matthias Seeger. Gaussian processes for machine learning. *International journal of neural*
340 *systems*, 14(2):69–106, 2004.
- 341 [13] Jasper Snoek, Oren Ripped, Kevin Swersky, Ryan Kiros, Nadathur Satish, Narayanan Sundaram,
342 Md Mostofa Ali Patwary, Prabhat, and Ryan P. Adams. Scalable Bayesian optimization using
343 deep neural networks. *32nd International Conference on Machine Learning, ICML 2015*,
344 3:2161–2170, 2015.
- 345 [14] John Dagdelen, Joseph Montoya, Maarten De Jong, and Kristin Persson. Computational
346 prediction of new auxetic materials. *Nature Communications*, 8(1):1–8, 2017.
- 347 [15] Krishna Kumar Saxena, Raj Das, and Emilio P. Calius. Three Decades of Auxetics Re-
348 search Materials with Negative Poisson’s Ratio: A Review. *Advanced Engineering Materials*,
349 18(11):1847–1870, 2016.
- 350 [16] Q Liu. Literature Review: Materials with Negative Poisson’s Ratios and Potential Applications
351 to Aerospace and Defense. Technical report, Australian Government Department of Defense,
352 2006.
- 353 [17] Wael A Salah and Mai Abuhelwa. Review of Thermoelectric Cooling Devices Recent Applica-
354 tions. *Journal of Engineering Science and Technology*, 15(1):455–476, 2020.
- 355 [18] Ran He, Gabi Schierning, and Kornelius Nielsch. Thermoelectric Devices: A Review of Devices,
356 Architectures, and Contact Optimization. *Advanced Materials Technologies*, 3(4):1700256,
357 2018.
- 358 [19] Jun Mao, Gang Chen, and Zhifeng Ren. Thermoelectric cooling materials. *Nature Materials*,
359 20(4):454–461, 2020.

- 360 [20] Anubhav Jain, Shyue Ping Ong, Geoffroy Hautier, Wei Chen, William Davidson Richards,
361 Stephen Dacek, Shreyas Cholia, Dan Gunter, David Skinner, Gerbrand Ceder, and Kristin A.
362 Persson. Commentary: The Materials Project: A materials genome approach to accelerating
363 materials innovation. *APL Materials*, 1(1):011002, 2013.
- 364 [21] Maarten De Jong, Wei Chen, Thomas Angsten, Anubhav Jain, Randy Notestine, Anthony
365 Gamst, Marcel Sluiter, Chaitanya Krishna Ande, Sybrand Van Der Zwaag, Jose J. Plata, Cormac
366 Toher, Stefano Curtarolo, Gerbrand Ceder, Kristin A. Persson, and Mark Asta. Charting the
367 complete elastic properties of inorganic crystalline compounds. *Scientific Data*, 2(1):1–13,
368 2015.
- 369 [22] Amir Yeganeh-Haeri, Donald J. Weidner, and John B. Parise. Elasticity of α -Cristobalite: A
370 Silicon Dioxide with a Negative Poisson’s Ratio. *Science*, 257(5070):650–652, 1992.
- 371 [23] Rod Lakes and K. W. Wojciechowski. Negative compressibility, negative Poisson’s ratio, and
372 stability. *Physica Status Solidi (B) Basic Research*, 245(3):545–551, 2008.
- 373 [24] Lisa J Rew, Bruce D Maxwell, Frank L Dougher, and Richard Aspinall. Searching for a needle
374 in a haystack: evaluating survey methods for non-indigenous plant species. *National Park
375 Biological Invasions*, 8:523–539, 2006.
- 376 [25] Wei Wei, Jinjiu Li, Longbing Cao, Yuming Ou, Jiahang Chen, W Wei, J Li, L Cao, Y Ou, and
377 J Chen. Effective detection of sophisticated online banking fraud on extremely imbalanced data.
378 *World Wide Web*, 16(4):449–475, 2012.
- 379 [26] Neil G Marchant and Benjamin I P Rubinstein. Needle in a Haystack: Label-Efficient Evaluation
380 under Extreme Class Imbalance. *KDD ’21, August 14–18, 2021, Virtual Event, Singapore*,
381 page 11, 2021.
- 382 [27] Mohammed Khalilia, Sounak Chakraborty, and Mihail Popescu. Predicting disease risks from
383 highly imbalanced data using random forest. *BMC Medical Informatics and Decision Making*,
384 11(1):1–13, 2011.
- 385 [28] Koby Crammer and Gal Chechik. A Needle in a Haystack: Local One-Class Optimization.
386 *Proceedings of the 21st International Conference on Machine Learning, Banff, Canada*, 2004.
- 387 [29] Haixiang Liu, Yuanming Hu, Bo Zhu, Wojciech Matusik, and Eftychios Sifakis. Narrow-
388 Band Topology Optimization on a Sparsely Populated Grid. *ACM Transactions on Graphics*,
389 37(6):1–14, 2018.
- 390 [30] Helena E. Nusse and James A. Yorke. Basins of Attraction. *Science*, 271(5254):1376–1380,
391 1996.
- 392 [31] George Datsoris and Alexandre Wagemakers. Effortless estimation of basins of attraction.
393 *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 32(2):023104, 2022.
- 394 [32] Philipp Hennig and Christian J. Schuler. Entropy search for information-efficient global
395 optimization. *Journal of Machine Learning Research*, 13:1809–1837, 2012.
- 396 [33] Belyaev Mikhail, Burnaev Evgeny, and Kapushev Yermek. Exact Inference for Gaussian
397 Process Regression in case of Big Data with the Cartesian Product Structure. 2014.
- 398 [34] Cheng Li, Sunil Gupta, Santu Rana, Vu Nguyen, Svetha Venkatesh, and Alistair Shilton. High
399 Dimensional Bayesian Optimization Using Dropout. *Proceedings of the 26th International
400 Joint Conference on Artificial Intelligence, IJCAI*, 2017.
- 401 [35] Zi Wang, Chengtao Li, Stefanie Jegelka, and Pushmeet Kohli. Batched High-dimensional
402 Bayesian Optimization via Structural Kernel Learning. *Proceedings of the 34th International
403 Conference on Machine Learning, Sydney, Australia, PMLR*, 70, 2017.
- 404 [36] Thang D Bui, Josiah Yan, and Richard E Turner. A Unifying Framework for Gaussian Process
405 Pseudo-Point Approximations using Power Expectation Propagation. *Journal of Machine
406 Learning Research*, 18:1–72, 2017.

- 407 [37] Gongjin Lan, Jakub M Tomczak, Diederik M Roijers, and A E Eiben. Time Efficiency in
408 Optimization with a Bayesian-Evolutionary Algorithm. 2020.
- 409 [38] David Eriksson, Michael Pearce, Jacob R Gardner, Ryan Turner, and Matthias Poloczek.
410 Scalable Global Optimization via Local Bayesian Optimization. 2020.
- 411 [39] Rommel G. Regis. Trust regions in Kriging-based optimization with expected improvement.
412 *Engineering Optimization*, 48(6):1037–1059, 2015.
- 413 [40] Y Diouane, V Picheny, R Le Riche, A Scotto, and Di Perrotolo. TREGO: a Trust-Region
414 Framework for Efficient Global Optimization. 2021.
- 415 [41] Michalis Titsias. Variational learning of inducing variables in sparse gaussian processes. In
416 David van Dyk and Max Welling, editors, *Proceedings of the Twelfth International Conference
417 on Artificial Intelligence and Statistics*, volume 5 of *Proceedings of Machine Learning Research*,
418 pages 567–574, Hilton Clearwater Beach Resort, Clearwater Beach, Florida USA, 16–18 Apr
419 2009. PMLR.
- 420 [42] Felix Leibfried, Vincent Dutordoir, S T John, and Nicolas Durrande. A Tutorial on Sparse
421 Gaussian Processes and Variational Inference. 2021.
- 422 [43] R. E. Turner and M. Sahani. Two problems with variational expectation maximisation for
423 time-series models. In D. Barber, T. Cemgil, and S. Chiappa, editors, *Bayesian Time series
424 models*, chapter 5, pages 109–130. Cambridge University Press, 2011.
- 425 [44] Bas van Stein, Hao Wang, and Thomas Back. Automatic configuration of deep neural networks
426 with parallel efficient global optimization. In *2019 International Joint Conference on Neural
427 Networks (IJCNN)*, pages 1–7. IEEE, 2019.
- 428 [45] Donald R Jones, Matthias Schonlau, and William J Welch. Efficient Global Optimization of
429 Expensive Black-Box Functions. *Journal of Global Optimization*, 13:455–492, 1998.
- 430 [46] Tinu Theckel Joy, Santu Rana, Sunil Gupta, and Svetha Venkatesh. Fast hyperparameter
431 tuning using Bayesian optimization with directional derivatives. *Knowledge-Based Systems*,
432 205:106247, 2020.
- 433 [47] Aaron Klein, Stefan Falkner, Simon Bartels, Philipp Hennig, and Frank Hutter. Fast Bayesian
434 Optimization of Machine Learning Hyperparameters on Large Datasets. 2017.
- 435 [48] Yuyu Zhang, Mohammad Taha Bahadori, Hang Su, and Jimeng Sun. FLASH: Fast Bayesian
436 Optimization for Data Analytic Pipelines. *Proceedings of the 22nd ACM SIGKDD International
437 Conference on Knowledge Discovery and Data Mining*, 2016.
- 438 [49] Albert Reuther, Jeremy Kepner, Chansup Byun, Siddharth Samsi, William Arcand, David
439 Bestor, Bill Bergeron, Vijay Gadepally, Michael Houle, Matthew Hubbell, Michael Jones, Anna
440 Klein, Lauren Milechin, Julia Mullen, Andrew Prout, Antonio Rosa, Charles Yee, and Peter
441 Michaleas. Interactive supercomputing on 40,000 cores for machine learning and data analysis.
442 In *2018 IEEE High Performance extreme Computing Conference (HPEC)*, pages 1–6. IEEE,
443 2018.
- 444 [50] M. D. McKay, R. J. Beckman, and W. J. Conover. A comparison of three methods for selecting
445 values of input variables in the analysis of output from a computer code. *Technometrics*,
446 42(1):55–61, 2000.
- 447 [51] Elon S Correa and Jonathan L Shapiro. Model Complexity vs. Performance in the Bayesian
448 Optimization Algorithm. In T P Runarsson, H-G Beyer, E Burke, J J Merelo-Guervos, L D
449 Whitley, and X Yao, editors, *Parallel Problem Solving from Nature*, pages 998–1007. Springer,
450 2006.
- 451 [52] Hoss Belyadi, Ebrahim Fathi, and Fatemeh Belyadi. Rock mechanical properties and in situ
452 stresses. *Hydraulic Fracturing in Unconventional Reservoirs*, pages 215–231, 2019.
- 453 [53] Yuriy M. Poplavko. Mechanical properties of solids. *Electronic Materials*, pages 71–93, 2019.

- 454 [54] B. Hinterleitner, I. Knapp, M. Poneder, Yongpeng Shi, H. Müller, G. Eguchi, C. Eisenmenger-
455 Sittner, M. Stöger-Pollach, Y. Kakefuda, N. Kawamoto, Q. Guo, T. Baba, T. Mori, Sami Ullah,
456 Xing Qiu Chen, and E. Bauer. Thermoelectric performance of a metastable thin-film Heusler
457 alloy. *Nature*, 576(7785):85–90, 2019.
- 458 [55] Georg K.H. Madsen and David J. Singh. BoltzTraP. A code for calculating band-structure
459 dependent quantities. *Computer Physics Communications*, 175(1):67–71, 2006.
- 460 [56] Hee Seok Kim, Weishu Liu, Gang Chen, Ching Wu Chu, and Zhifeng Ren. Relationship
461 between thermoelectric figure of merit and energy conversion efficiency. *Proceedings of the*
462 *National Academy of Sciences of the United States of America*, 112(27):8205–8210, 2015.
- 463 [57] Wei Hsin Chen, Po Hua Wu, Xiao Dong Wang, and Yu Li Lin. Power output and efficiency of
464 a thermoelectric generator under temperature control. *Energy Conversion and Management*,
465 127:404–415, 2016.
- 466 [58] H. Julian Goldsmid. Bismuth telluride and its alloys as materials for thermoelectric generation.
467 *Materials*, 7(4):2577–2592, 2014.
- 468 [59] Pedro M. Rodrigo, Alvaro Valera, Eduardo F. Fernandez, and Florencia M. Almonacid. An-
469 nual Energy Harvesting of Passively Cooled Hybrid Thermoelectric Generator-Concentrator
470 Photovoltaic Modules. *IEEE Journal of Photovoltaics*, 9(6):1652–1660, 2019.