




## A perspective on Bayesian methods applied to materials discovery and design

**Raymundo Arróyave** , Department of Materials Science and Engineering, Texas A&M University, College Station, TX 77843, USA; J. Mike Walker '66 Department of Mechanical Engineering, Texas A&M University, College Station, TX 77843, USA; Wm Michael Barnes '64 Department of Industrial and Systems Engineering, Texas A&M University, College Station, TX 77843, USA

**Danial Khatamsaz**, J. Mike Walker '66 Department of Mechanical Engineering, Texas A&M University, College Station, TX 77843, USA

**Brent Vela, Richard Couperthwaite, and Abhilash Molkeri**, Department of Materials Science and Engineering, Texas A&M University, College Station, TX 77843, USA

**Prashant Singh**, Ames Laboratory, U.S. Department of Energy, Iowa State University, Ames, IA 50011, USA

**Duane D. Johnson**, Ames Laboratory, U.S. Department of Energy, Iowa State University, Ames, IA 50011, USA; Department of Materials Science & Engineering, Iowa State University, Ames, IA 50011, USA

**Xiaoning Qian**, Department of Electrical and Computing Engineering, Texas A&M University, College Station, TX 77843, USA

**Ankit Srivastava**, Department of Materials Science and Engineering, Texas A&M University, College Station, TX 77843, USA

**Douglas Allaire**, J. Mike Walker '66 Department of Mechanical Engineering, Texas A&M University, College Station, TX 77843, USA

Address all correspondence to Raymundo Arróyave at [rarrayave@tamu.edu](mailto:rarrayave@tamu.edu)

(Received 20 April 2022; accepted 4 October 2022; published online: 26 October 2022)

### Abstract

For more than two decades, there has been increasing interest in developing frameworks for the accelerated discovery and design of novel materials that could enable promising and transformative technologies. The Integrated Computational Materials Engineering (ICME) program called for integrating computational tools to establish linkages along process–structure–property–performance chains. The Materials Genome Initiative called for integrating experiments and computations within data science frameworks as a strategy to accelerate the materials development cycle. While these frameworks and paradigms have been quite influential, traditional ICME or data science-based approaches tend to have some limitations, mainly when querying the materials space is costly and very little information is available. Bayesian methods are more suitable in this context due to their efficiency gains. To this end, the materials discovery problem is framed as a Bayesian optimization (BO). Different examples in which BO has been applied to solve materials discovery problems are presented. The methods/examples discussed include BO under model uncertainty, multi-information source BO, multi-objective and multi-constraint BO, and batch BO. Bayesian Materials Discovery is a promising area of research that is likely to become more influential as more attention is put on autonomous materials discovery platforms. Therefore, a discussion is provided on the potential development of such methods to increase the ability of existing platforms in materials discovery. The ultimate goal is to pave the way to autonomous materials discovery.

### Introduction

Until well into the twentieth century, the materials discovery/development cycle was mostly carried out using highly heuristic approaches lacking well-defined frameworks. The lack of a systematic approach to materials design meant that very little could be learned/transferred from successful materials development programs. To address this issue, Olson<sup>[1]</sup> extended Cyril Stanley Smith's ideas on materials as hierarchical (multi-scale) systems<sup>[2]</sup> and proposed that materials systems were amenable to design using engineering design principles informed by scientific understanding. In that work, Olson recognized the importance of developing process–structure–property–performance (PSPP) relationships to establish causal connections between processing conditions and materials performance. Such links could then be used to understand how a given material could be modified to optimize its performance.

Olson's idea to use computational tools to establish such linkages was embraced as the underlying framework guiding Integrated Computational Materials Engineering (ICME),<sup>[3]</sup> as

an approach to accelerate the insertion of new materials into technological applications. Provided such linkages could be transformed into quantitative relationships, PSPPs could, in turn, be inverted to optimize the materials to meet specific performance goals. Having established quantitative PSPP relationships, the actual optimization over a given materials design space could be carried out with a number of algorithmic optimization schemes.<sup>[4,5]</sup> Unfortunately, simulation-driven materials design requires solving many significant challenges, such as effectively linking different simulation tools<sup>[6]</sup> and addressing the considerable uncertainties in the models, model parameters, and the experiments used to validate them.<sup>[7]</sup> This problem is compounded by the high computational cost of many of the most precise and predictive tools at the disposal of ICME practitioners.

While the ICME's focus was on linking computational PSPP model chains, the Materials Genome Initiative (MGI)<sup>[8]</sup> put forward an aspirational program to accelerate the materials development cycle through the tight integration of computations and experiments through data analytics techniques. MGI

has been embraced by almost every single community within the materials science discipline, and the field is undergoing an exponential expansion in the number of efforts taking place around the world involving the leveraging of data science, machine learning (ML), and artificial intelligence (AI) to accelerate our understanding of materials spaces. While many types of problems are being addressed using these tools, much work under the MGI paradigm requires a significant amount of data. Consequently, some of the most prominent applications of MGI principles to materials discovery have exploited extensive data, either obtained from experimental combinatorial synthesis of material libraries<sup>[9]</sup> or high-throughput (HTP) computations.<sup>[10]</sup>

ICME- and MGI-inspired approaches have resulted in numerous successful cases of materials discovery and design. However, such approaches tend to be limited, particularly considering that typical materials development efforts are always carried out under stringent resource constraints. ICME methods,<sup>[11,12]</sup> for example, (i) build and exploit process–structure–property–performance (PSPP) relationships<sup>[7]</sup> at a considerable computational expense; (ii) do not readily incorporate data from experiments within their framework; and (iii) are sequential, which means that they tend to be deployed by evaluating materials design choices one at a time. Traditional HTP combinatorial computational<sup>[10]</sup> and experimental<sup>[13,14]</sup> approaches, on the other hand, are (i) incapable of dealing with the high dimensionality and complexity of typical materials design spaces and (ii) are ‘one-shot’ or ‘open-loop’ schemes without a built-in iterative framework and are unable to prescribe future actions given the current state of knowledge.<sup>[15]</sup> *Moreover, traditional ICME and HTP approaches are suboptimal in resource utilization, as they do not guarantee that each proposed choice is optimal.*

The limitations described above can potentially be overcome using inherently resource-aware frameworks that provide a principled way to augment our knowledge of the state of a materials design space iteratively. Because of this, one of the best approaches to accelerate the design of a materials space is to frame this problem as an optimal experimental design task within a Bayesian Optimization (BO) setting.<sup>[16–21]</sup> As a result, the number of works that deploy BO in materials optimization tasks is growing rapidly.<sup>[22–30]</sup> Yet, most of the works published thus far have focused on translating the materials discovery problem into a single black box optimization challenge. Here, we revisit some recent work in which traditional BO frameworks have been modified for them to be better adapted to solving problems in materials science, mainly when one can get access to multiple tools/models/experiments to query the material space and one (potentially) has access to experimental or computational platforms capable of operating in batch/parallel mode.

## A short tutorial on Bayesian optimization

Any Bayesian methods rely on the deployment of the Bayes theorem to update our prior knowledge (beliefs) about a system once new information has been acquired.<sup>[31]</sup> Bayesian

optimization (BO) is a sequential design strategy for increasing (optimizing) the system’s performance or augmenting our knowledge about it. It relies on Bayesian updates of prior beliefs through acquiring new information. BO methods tend to be most effective when confronted with functions that are costly to compute, are not analytical (are ‘black boxes’), cannot be evaluated precisely due to intrinsic noise, and/or provide no facilities for estimating their gradients. We note that many resources explain in great detail what BO is, and the reader is kindly directed to some of these sources<sup>[20,32,33]</sup> if they want to understand these frameworks at a deeper level.

BO methods have many essential features that make them a highly effective optimization scheme. As a consequence of the sequential exploration of a design space, BO methods construct a statistical model of the ground truth being explored that is orders of magnitude cheaper to query than the ground truth itself. The primary ingredients of any BO algorithm are (1) a statistical (surrogate) model used to predict the outcomes of experiments yet to be carried out; (2) a policy that prescribes in a deterministic, principled, and unambiguous manner the best action to take to meet a given (set of) objective(s) given our current state of knowledge of the system.<sup>[34]</sup> Given an initial dataset,  $\mathcal{D}$ , a typical sequential BO scheme would then follow the following general algorithm (see also Fig. 1).

```

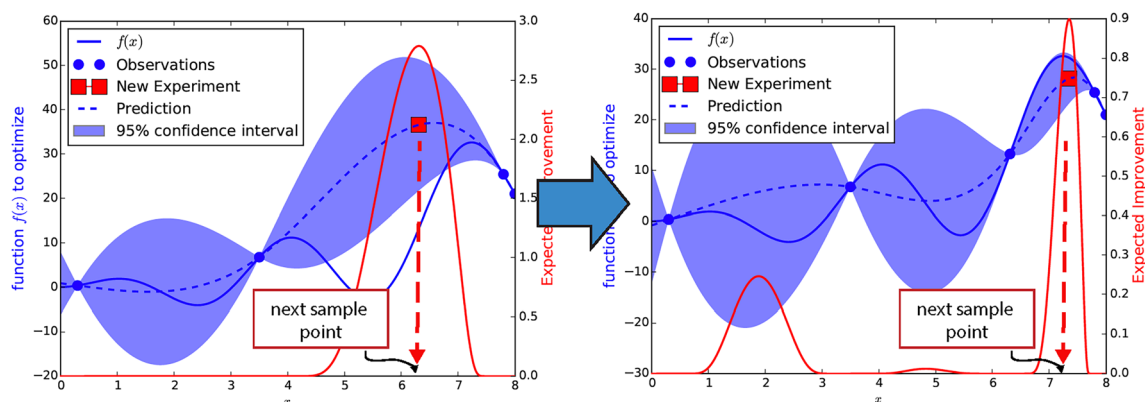
input: initial dataset  $\mathcal{D}$ 
repeat
     $x \leftarrow \text{POLICY}(\mathcal{D})$ 
     $y \leftarrow \text{OBSERVE}(x)$ 
     $\mathcal{D} \leftarrow \mathcal{D} \cup \{(x, y)\}$ 
until termination condition reached
return  $\mathcal{D}$ 

```

**Algorithm 1:** General Bayesian Optimization algorithm [32].

While, in principle, many types of statistical models can be used to simulate a given experimental space, the overwhelming majority of BO approaches use Gaussian processes (GPs).<sup>[35–37]</sup> A GP is a nonparametric statistical model for the objective function,  $f$ , defined as a stochastic process  $p(f) = \text{GP}(f; \mu, K)$  with a mean function,  $\mu$ , and a covariance function (or kernel)  $K$ . The reason why GPs are used in BO is because of their mathematical properties (including smoothness and controllable modeled correlation among observed points). Besides their intrinsic capability to predict uncertainty in our knowledge of the objective function over the design space, the critical characteristic of GPs is the notion of a measure of the degree of correlation between observations: observations that are close tend to yield similar values of the objective function. The degree (and type) of correlation is encoded by the kernel,  $K$ . The latter can take a wide range of forms (as long as it has some appropriate mathematical properties) depending on the nature of the system being optimized.

The second ingredient of BO is the *policy* (or utility function) used to make optimal sequential decisions on the queries. The policy for decision-making is constructed from the posterior distribution of the surrogate model for the design space.



**Figure 1.** Schematic illustration of Bayesian Optimization (BO): from a limited number of observations on a system (blue solid line), a stochastic model (dashed blue line and shaded area) is built. The following observation is determined by accounting for the tradeoff between the exploitation of the current knowledge and the exploration of the unknown regions of the design domain  $x$ . In this case, Expected Improvement (EI) is the metric used and thus the policy falls within the Efficient Global Optimization (EGO) framework.<sup>[34]</sup> Reproduced from Talapatra et al.<sup>[38]</sup>

Over the past decades, dozens of policies or *acquisition functions* have been developed and deployed in BO frameworks. Yet, all acquisition functions address the tension between exploitation (querying the design space in locations predicted to be optimal by the surrogate model) and exploration (querying the design space where there is uncertainty about the surrogate model). These acquisition functions are much cheaper to evaluate than the ground truth/objective function. When GPs are used as the surrogate function, many acquisition functions can be computed as analytical functions of the properties of the GPs, such as the mean response and its variance.<sup>[16,20,32]</sup> Examples of such acquisition functions include expected improvement (EI), probability of improvement (PI), and upper confidence bound (UCB). Unfortunately, the effectiveness of different policies is highly problem dependent, and it is thus impossible to know a priori the optimal policy for a given problem.

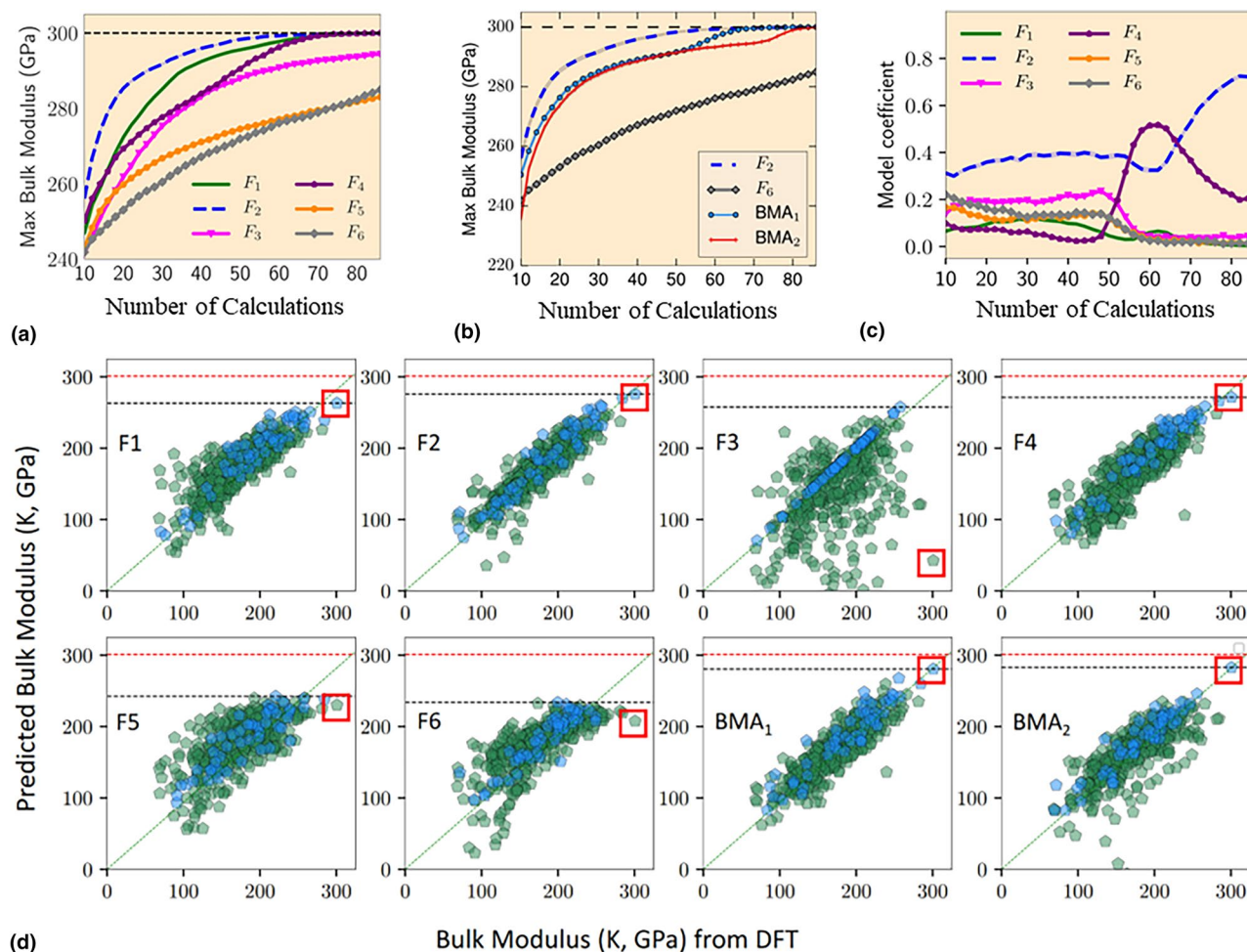
Several approaches have employed BO in a materials science context, for example, in molecular design and discovery and smart and additive manufacturing design optimization of materials, such as piezoelectric materials, catalysts, and structural materials.<sup>[29,30,38–40]</sup> They have shown dramatic improvements in the efficiency with which experimental/computational resources are utilized to find materials with optimal properties/performance, even within exclusively experimentally-based frameworks.<sup>[41]</sup> Yet, most approaches to materials discovery that deploy BO-based frameworks are still limited. BO-based materials discovery has been based on translating existing methods to optimize general ‘black box’ functions. Yet, in materials science, one typically has multiple ways of querying a design space. Most approaches tend to deploy BO to solve single-objective optimization problems when realistic materials development efforts require optimality across many dimensions of the performance space while satisfying several constraints. Moreover, most BO-based materials discovery approaches developed to date do not take advantage of the quickly evolving high-throughput experimental and computational platforms.

As will be seen below, there are many opportunities to tailor BO methods to materials discovery.

## Bayesian optimization under feature importance uncertainty

In typical realistic materials discovery challenges, there is typically little to no prior knowledge of the behavior of the objective function to explore. For example, we may not know the ‘shape’ of the (potentially multi-objective) materials response space. Thus, it may not be possible to develop a realistic model for the correlation (or distance) between points in the design space, encoded in a GP as the covariance function,  $K$ . Furthermore, in many cases, there may not be sufficient information to gauge which features are most correlated with the optimization objective(s). This lack of knowledge stems from cost limitations on doing enough experiments due to the complexity of the objective(s) or the involvement of many features. For example, initially, we may not know which features of a given material are the most correlated with a particular property or performance metric that we want to optimize. In some early examples of BO applied to materials science, feature selection (or feature engineering) was proposed as a preliminary step before initiating a BO optimization loop. Unfortunately, in most practical applications, the feature and/or model selection step is highly uncertain due to the sparseness in the data available. This inability to narrow down the feature space compromises the ability of BO to optimally sample a given design space since the performance of BO methods tends to degrade as the dimensionality of the problem increases.<sup>[42]</sup>

In recent work,<sup>[38,43]</sup> it has been shown that it is possible to carry out feature selection while carrying out an optimal sequential experimental design. In that setting, the feature space was partitioned into feature subspaces that, together with specific covariance functions,  $K$ , constitute different models describing the objective space. To account for this model



**Figure 2.** Bayesian Optimization under Model Uncertainty (BOMU): (a) performance of different models during a sequential BO task. By partitioning features into six different subsets, six different models are created to estimate the same quantity of interest (QoI). (b) Comparison between the best and the worst models and two effective models estimated through Bayesian Model Averaging (BMA). (c) Evolution of the model probability during the sequential experimental campaign. As more data are added to the system, model probability, and importance change. (d) Benchmarking results of using each model to learn and predict the bulk modulus of a set of samples. The red square shows the largest bulk modulus found in DFT calculations of queried samples in the optimization process, and the dashed red line shows the largest value ever seen.<sup>[38]</sup>

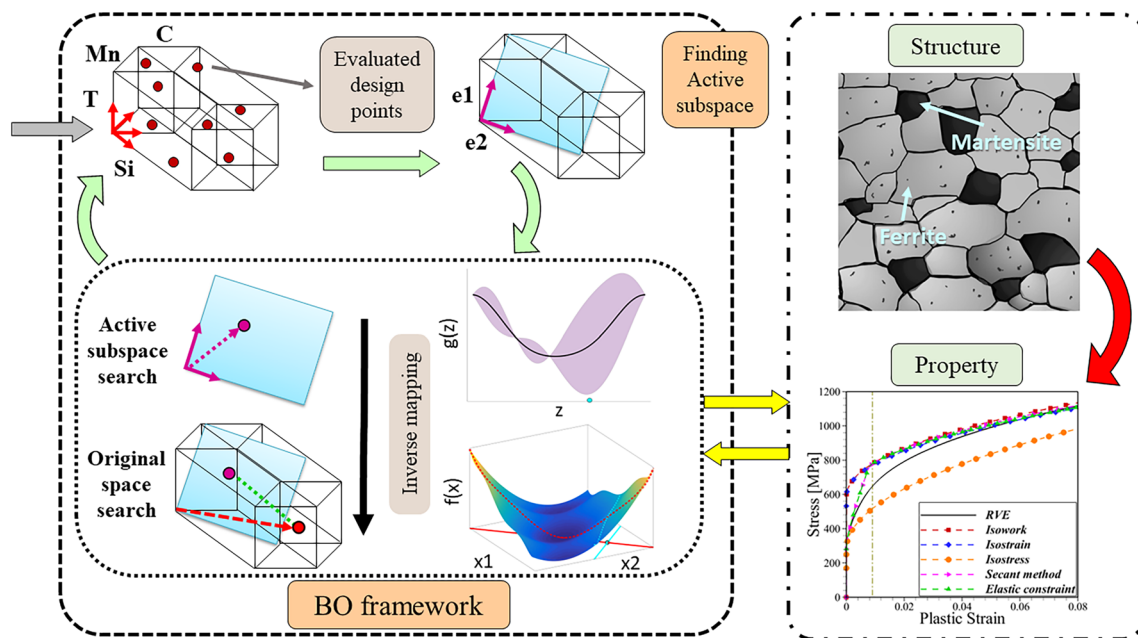
uncertainty, it is possible to weigh all the possible models by their probability of being the true model. When the weighing is based on the Bayesian model evidence, such model assembly is known as Bayesian model averaging (BMA).<sup>[44]</sup> By incorporating BMA within BO, Bayesian optimization under model uncertainty (BOMU) can simultaneously carry out feature selection and BO. In BOMU (see Fig. 2), the feature space available for exploration is partitioned into subsets corresponding to different hypotheses as to what features most correlate with the objective. Each feature set-kernel/covariance function combination is treated as an individual model, and each member of the model set  $\{\mathcal{M}_1 \dots \mathcal{M}_N\}$  would vary in their ability to predict the outcome of the observations that have yet to be made over the design space.

To take action (or query/observation), the policy or utility of observing every yet-to-be-queried point is evaluated as the weighed experiment utility, using the probability of a given

model being the correct model as the weights. After an observation is made, the updated dataset is used to evaluate the probability of the models and a new cycle begins. As a consequence of this model selection step, models would change their probability or weight based on how good they are at predicting the augmented dataset. An interesting aspect of BOMU is that it has been observed that the model probabilities (Fig. 2) tend to evolve in a non-monotonic manner, with the importance of some models decreasing and increasing in importance as a function of the stage of an experimental campaign. BOMU is also capable of immediately detecting uninformative feature sets, reducing their probability to close to zero early in the experimental campaign.

While BOMU carries out feature selection by evaluating the posterior probabilities of competing models, there are other approaches to BO with automated model selection. Having a feature space, the framework described by Malkomes and





**Figure 3.** Implementation of adaptive active *subspace* method within a Bayesian optimization (BO) framework. The basic idea is to find the active subspace, i.e., the directions in the material design space that give the most significant variation in the target property, using the available data at every stage of the optimization task. Next, the design space is mapped to the active subspace, and the first step of the BO framework is applied to find the ‘next best point’ to evaluate within the active subspace. The ‘next best point’ is then mapped back to the original design space by implementing a second BO step. Finally, the design space is evaluated at this best point. This new data is added to the framework for the next iteration. In the example shown inside the BO framework window, the objective function  $f(\mathbf{x})$  is represented as  $g(\mathbf{z})$  in the active subspace. Thus, instead of variables  $x_1$  and  $x_2$ , the objective function is estimated on a 1-dimensional active subspace space of  $z$ . This figure summarizes the implementation of the active subspace method for designing dual-phase steels in the presence of multiple information sources estimating the same property at different levels of fidelity.<sup>[45]</sup>

collaborators<sup>[46–48]</sup> will be capable of dynamically selecting optimal models (with different covariance functions/kernels) through the use of Bayesian optimization over the model space using the concept of ‘kernel of kernels,’ which can be used to measure the distance between the covariance functions (kernels) of Gaussian Process regression models. The framework generates arbitrary candidate kernel functions by utilizing the concept of ‘kernel grammars.’<sup>[49]</sup> We note that BOMU can potentially explore the feature and model space and the policy used. In practice, different feature set-model-policy tuples would be evaluated, and their weight would be modified depending on their effectiveness.

BOMU is a practical approach when there is uncertainty in the feature space. However, a significant limitation is that the partitioning of the feature space is done in an ad hoc manner. Thus there is always a risk of not capturing the truly important features/degrees of freedom. An alternative approach is the *active subspace method*,<sup>[50]</sup> which is a technique to identify the directions in a given design space that have the most significant change relative to a given objective, effectively reducing the dimension of the problem. Once a subspace has been identified, an acquisition function/policy is used to evaluate the best point within this lower-dimensional space. Finally, an observation is prescribed in the actual higher-dimensional design space by inverse mapping from the lower to the higher-dimensional spaces (Fig. 3).

Compared to other approaches for automated feature selection within BO, an advantage of the active subspace method is the ability to dynamically change the effective dimensionality of the problem as the experimental campaign progresses. An extensive evaluation of the active subspace method found that the effective dimension of a given materials design problem did not exceed two for most sequential experimental campaigns.<sup>[45]</sup> Crucially, the effective dimensions of the problem changed continuously through the design sequence, pointing to the fact that attempting to carry out feature selection before the BO stage of a materials design task may not be a robust strategy. The dynamic nature of the feature importance found in this work is in line with what was found in BOMU, where the weight of some models changed dramatically as the BO sequence progressed. These observations imply that in materials design/optimization problems, only a few features/degrees of freedom are active at any given time.

## Multi-information source Bayesian optimization

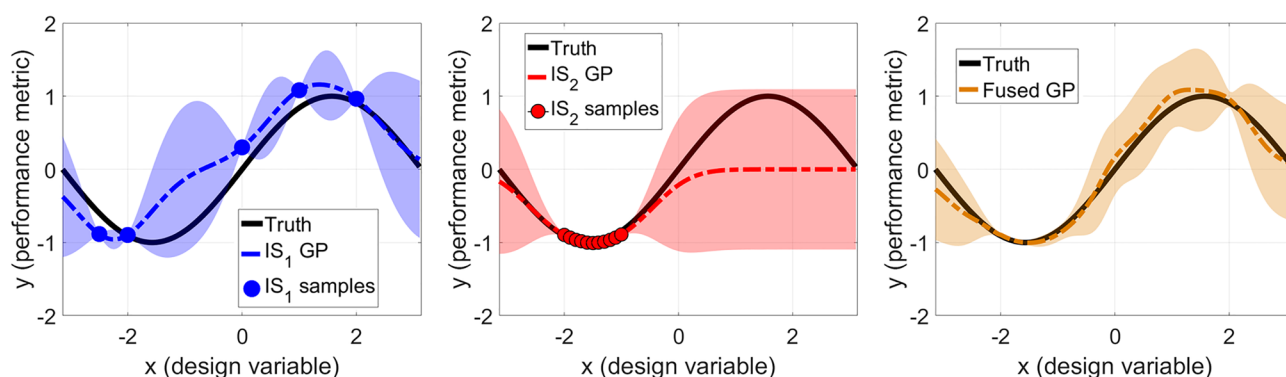
While mapping a materials discovery task to a ‘black box’ optimization can significantly accelerate the materials discovery process, such a simplistic approach does not take

advantage of a crucial aspect of materials research. For example, ‘black box’ BO approaches do not account for the fact that it is always possible to use several methods to query a particular materials space. Such methods include, for example, different models (with different cost and fidelity) connecting different linkages along PSPP relationships or experiments measuring different attributes of the material to be investigated. As such, ‘black box’ BO is not tailored to materials design tasks. On the other hand, traditional ICME-based approaches rely on quantitative PSPP chains, whose linkages can be evaluated using computational models. A significant challenge for these approaches is that they tend to consider a single model per chain—i.e., they ignore the possibility of combining the predictive power of different models simultaneously. A further limitation is that most typical ICME approaches are not capable of combining experiments and simulations since they use experimental information only as a way to verify or validate the models used. Integrating multiple information sources into a materials discovery task could provide significant advantages.

Other fields in science and engineering have already developed sophisticated approaches to combine multiple information sources within optimization schemes. In these approaches, multiple information sources may be used to approximate the behavior of an expensive-to-query ‘ground truth.’ In such multi-fidelity approaches,<sup>[51]</sup> it is possible to combine multiple information sources to have better information about the ‘ground truth,’ at a reduced cost. Figure 4 shows how one can transfer knowledge among information sources to enhance learning through Bayesian information fusion. Practically, such knowledge transfer is constructed by learning and exploiting the statistical correlations among the different information sources. The underlying assumption of this type of information fusion is the existence of correlation among different information sources and between a given information source and the ground truth. In materials science, we can assume that such an assumption is warranted, given that all information sources for a given materials design task predict different aspects of

the material’s behavior and are thus causally correlated. The common goal in any application using multi-information source approaches is reducing associated costs. By collecting information from different sources, which are usually lower in fidelity but cheaper to evaluate, it is possible to obtain information about the ground truth (highest fidelity model) via the correlation between the models and the ground truth. Consequently, it allows for a cheaper learning process and prevents conducting expensive experiments, in contrast to single model approaches which have no choice but to perform costly experiments to learn an objective function.

In many materials design applications, collecting experimental data is considered the most accurate way of estimating a material’s property of interest. However, this comes at the expense of substantial experimental costs (money, time, computational resources, etc.), which makes it almost impossible to complete sufficient numbers of experiments to make any conclusions about the optimal design, particularly in high-dimensional feature spaces. In such cases, using less accurate approaches (exploiting lower fidelity computational models, for example) to estimate a quantity of interest (QoI) and correlating it to the ground truth (highest fidelity model or experimental data) can be beneficial, cost-wise. It is often the case that a set of QoIs can be estimated through different kinds of experiments varying in cost and accuracy (fidelity), yielding different results. It is also usually the case that computational models are exploited to estimate a set of QoIs to avoid the challenges of collecting expensive experimental data. Computational models representing the same system of materials are also different in cost and fidelity since different assumptions and mathematical simplifications are made to create such models. Assuming that any information source contains some useful information about the QoIs, collecting this information from several lower fidelity models helps to reduce the overall design costs. Another point is that the fidelity of information sources may vary through the input space; thus, exploiting multiple information sources is beneficial because more useful sources can be deployed depending on the region of interest.



**Figure 4.** Left: Fitting a GP to information source IS 1. Center: Fitting a GP to information source IS 2. Right: Fused GP from IS 1 and IS 2 GPs. Learning the relationship between IS 1 and IS 2 could improve the fused GP even more in the region where we have no samples of IS 2.

Recently, such a multi-information source framework was used to predict the crystallization tendency of polymers. In work by Venkatram and collaborators,<sup>[52]</sup> it was argued that predicting a new polymer crystallinity is not trivial due to experimental challenges, such as having the necessary infrastructure and sample preparation. On the other hand, group contribution methods are less accurate but much cheaper to obtain. Thus, they fused data from group contribution methods and fewer experimental data available from the literature, obtaining a more accurate model without performing any expensive experiments. Another recent application was the work by Pilania et al. on the bandgap prediction of solids using multi-fidelity approaches.<sup>[53]</sup> In that work, a larger number of samples from a cheap model (PBE) and fewer samples from a more expensive model (HSE06) are used to obtain a fused model and make predictions. The work by Tran and collaborators<sup>[18]</sup> present an application in optimizing the bulk modulus in a ternary composition space using multi-information source BO. A similar work employed multi-information source BO in geometry optimization of CO via employing different DFT functionals as information sources.<sup>[54]</sup>

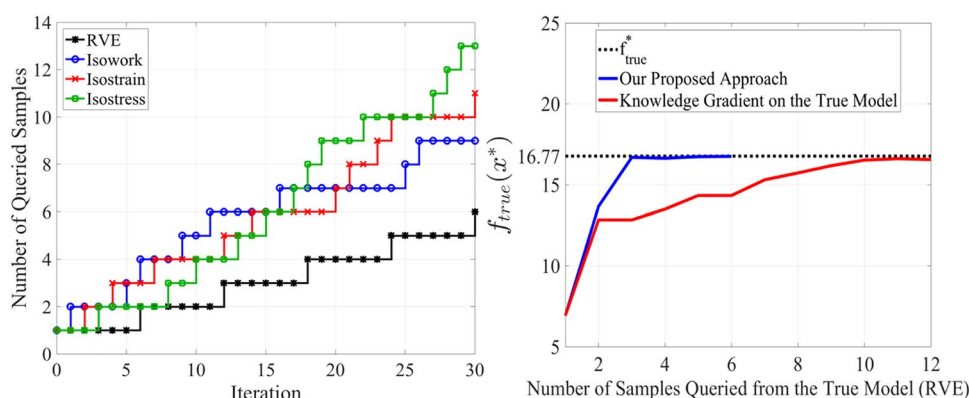
A general multi-information source BO framework has been proposed and deployed to optimize normalized strain hardening rate in designing dual-phase steels—the framework can easily be extended to a wide range of materials discovery/design problems.<sup>[55,56]</sup> In this framework, one first constructs intermediate surrogate models (GPs for convenience due to their mathematical properties) for each available information source. For example, these information sources could be different models connecting materials design inputs and outputs to varying levels of fidelity and physical rigor. Then, using a so-called reification<sup>[57]</sup> procedure, a fused model is constructed from all information sources by learning the statistical correlation among all the sources, including the ‘ground truth.’ The fused model, described as a GP,  $GP(f; \mu, K)$  is then used to evaluate the value of the acquisition function on each candidate point in the experimental space yet to be queried. Having established the best point to query, the multi-information source

Bayesian optimization (MISBO) framework then selects the next information source to query, given an information source fidelity, uncertainty, and cost.

MISBO is much more efficient at solving materials optimization problems<sup>[55]</sup> than even state-of-the-art BO methods that only query the ground truth—see Fig. 5. Since every information source is mapped to a GP, it has been shown that MISBO can be used in problems in which explicit PSPP connections can be exploited by mapping different linkages along the PSPP chain to GPs and then connecting the inputs and outputs of the different GPs along a model chain using simple statistical correlations. As such, MISBO can be mapped directly to ICME frameworks,<sup>[56]</sup> the advantage being a much more efficient utilization of the available information. An additional benefit of MISBO is that it is possible to account explicitly for the cost of a given information source, which is crucial when there are stringent budget constraints.<sup>[56]</sup> Other groups<sup>[58]</sup> have presented similar approaches. Pilania et al.,<sup>[53]</sup> for example, have shown how it is possible to combine low-fidelity approximate predictions for materials properties to supplement high-fidelity, expensive computational simulations, arriving at accurate predictive models at a relatively low cost. MISBO differs from other multi-fidelity approaches in that it does not assume the existence of a hierarchy among the different models as they are all *simultaneously* fused to predict the behavior of the ‘ground truth.’

## Bayesian optimization under multiple objectives and constraints

By far, most approaches to BO have focused on deploying acquisition functions tailored to balance the exploration and exploitation of single-objective design spaces. The reason for this is that in single-objective optimization problems, the objective of the decision-making process is clear to optimize the value of a target system response. However, in materials science, it is often the case that a materials design task has



**Figure 5.** (Left) Number of samples queried from the true model and the information sources in each iteration. (Right) The optimal solution was obtained by the proposed approach and by applying the knowledge gradient on a GP of only the true data for different numbers of samples queried from the true model.<sup>[55]</sup>

multiple objectives to optimize. For example, one would like materials to be strong and ductile or materials with higher energy storage capacity and very short charging times. In many cases, these objectives compete against each other, resulting in performance tradeoffs: strong materials are not necessarily the most ductile, for example.

A multi-objective optimization problem can thus be solved in two different ways:

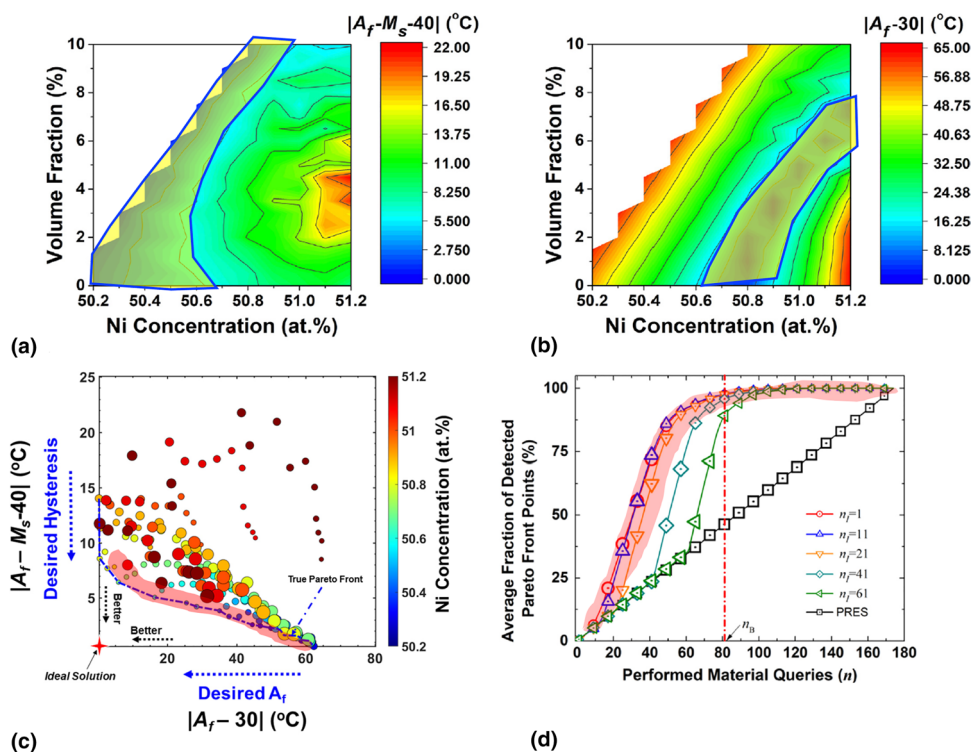
One strategy is to construct a single utility function from different objectives.<sup>[59]</sup> Such a utility function is built based on the preferences of the decision-maker(s) and thus has some degree of subjectivity. Once this single-objective function has been constructed, the optimization is carried out in the usual way, using BO or non-BO methods. Building a single utility function implies that (subjective) preferences do not change as more information is acquired about the system. If the utility function changes because of changes in the preferences for different objectives, the optimization would have to be carried out from scratch.

An alternative approach is to not establish any subjective preference for the different individual objectives and instead try to discover the Pareto front in the multi-objective space.<sup>[60]</sup> The Pareto front is the set of non-dominated solutions being optimal along with one of the directions in the objective space. Once a Pareto front is discovered, a preference model (or utility function) could be constructed from the Pareto set—any utility

function constructed with dominated solutions (i.e., solutions enclosed by the Pareto front) would always be suboptimal. While discovering the Pareto front allows more flexibility in the decision-making process, the discovery of the Pareto set is challenging, particularly in high-dimensional spaces.

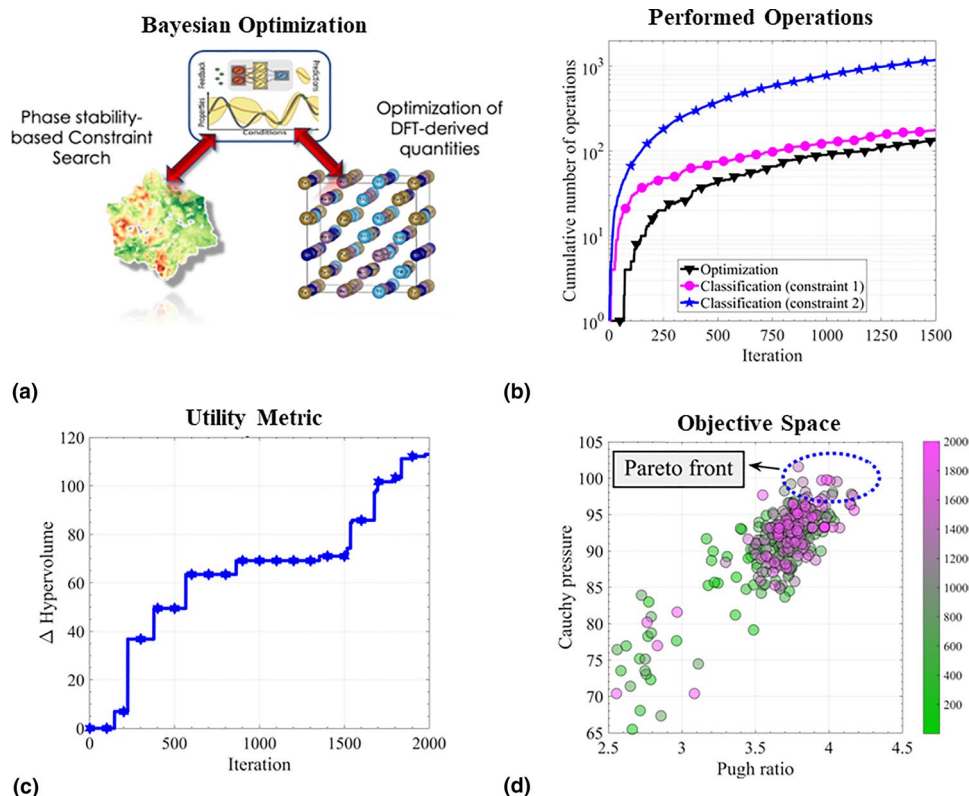
In the context of BO, the problem is paradoxical: on the one hand, one is interested in discovering the Pareto front. On the other hand, however, sequential BO policies always propose a single action to take at any given time. Therefore, one must construct a scalar function that measures the utility of carrying out an experiment at a point in the design space, regardless of the number of objectives. There are not as many acquisition functions in multi-objective BO as in single-objective BO. One of the most widely used metrics in multi-objective BO is the Expected HyperVolume Improvement (EHVI).<sup>[61]</sup> This metric measures the expected change in the Pareto front (toward the direction of non-domination) when a query is made at a given point. Under this policy, the optimizer can identify the Pareto front under resource constraints. This approach has been used before in multi-objective materials optimization—see Fig. 6.<sup>[62]</sup> The efficiency of multi-objective BO methods tend to outperform other optimization approaches, particularly when the number of objectives increases.

While the ability to solve multi-objective problems in a Bayesian optimal setting brings materials discovery to a more realistic setting, one must consider many problems in materials



**Figure 6.** (a) and (b) Multi-objective space in the optimization of chemical composition and precipitate volume fraction to optimize the performance of a precipitation-strengthened shape memory alloy (SMA). (c) Pareto front identified using an optimal Bayesian policy. (d) Optimal BO approaches outperform other policies and can discover a larger fraction of the Pareto front with fewer queries to the system.<sup>[62]</sup>





**Figure 7.** (a) A multi-objective, multi-constraint optimization problem consists of identifying regions in an alloy space that satisfy some minimal performance requirements (i.e., meet performance constraints) while simultaneously optimizing a set of intrinsic properties. (b) Operations performed during the design process are categorized as either classification or optimization depending on querying constraint models or objective functions. (c) As more iterations are completed in the optimization loop, better non-dominated designs are found. Hypervolume is a metric to measure the quality of discovered non-dominated designs and its improvements are directly related to discovering a better Pareto front. (d) Queried samples and non-dominated designs were queried during the process of optimizing the Pugh Ratio and Cauchy pressure as an indication of the ductility of an alloy. The colormap is used to illustrate the order of queries.<sup>[64]</sup>

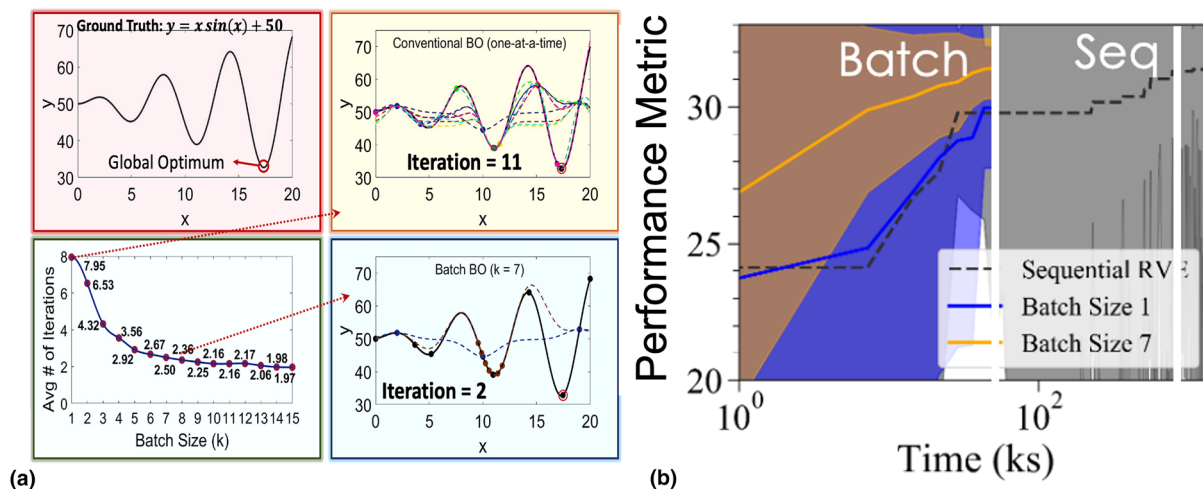
science in which the performance requirements are so stringent that the discovery process is constraint-dominated.<sup>[63]</sup> When there is insufficient information about the materials design space, one needs to allocate resources to identify the feasible space (i.e., the region in the design space with solutions that satisfy all the performance constraints) and find the Pareto set. In some cases, the constraints can be evaluated in terms of easy and fast-to-evaluate analytical models constructed as a function of the design variables. However, in many cases, assessing the constraints constitutes a costly endeavor. For example, evaluating whether an alloy is resistant to oxidation may involve a complex experimental protocol lasting 24–100 h.

The identification of the feasible space can be mapped to a Bayesian Classification (BC) task, in which the boundary separating the feasible and infeasible regions in the materials design space is found in the least number of queries possible. Ideally, such a BC task would be set up using multiple information sources, as in MISBO.<sup>[55,56]</sup> In recent work, it has been shown how it is possible to combine multi-information source BO with multi-information source BC to solve classification/optimization materials design problems—see Fig. 7. While the results of

this approach are promising, there are significant opportunities for further development. For example, it is unclear how to optimally divide resources between optimization and classification in a principled, algorithmic manner. One approach could be first identifying the feasible space and then allocating all the remaining resources to identify the Pareto set. However, this approach may be too risky when resources are limited, as the total (experimental/computational) budget could be spent well before the exploration of the Pareto set begins. Other approaches may switch between classification and optimization as the campaign progresses, resulting in suboptimal performance. Perhaps building criteria for uncertainty measures in the feasible/unfeasible boundary and the Pareto set could be constructed to optimize the process.

### Bayesian optimization over high-throughput parallel computational/experimental settings

One of the significant drawbacks of ICME and modern BO-based materials optimization approaches is their sequential nature and their inability to make use of high-throughput



**Figure 8.** (a) Demonstration of batch Bayesian BO. The sequential approach is much less efficient than batch querying. (b) BBO results in orders of magnitude faster materials optimization than standard sequential BO.<sup>[75]</sup>

workflows pervasive in HPC environments,<sup>[65,66]</sup> as well as modern synthesis<sup>[67,68]</sup> and characterization<sup>[69–71]</sup> facilities. In addition, traditional ICME approaches suffer from their non-optimality already, so in the context of materials design, providing BO with the ability to make *multiple optimal decisions* at once is perhaps the only viable approach toward parallelizing a design framework.<sup>[72]</sup>

The approach for Batch BO by Joy and collaborators<sup>[72]</sup> is among the most promising thus far. Their method is based on the premise that when only a few data points are available, it is impossible to have much confidence in the covariance structure of the GP models describing the correlations between data. A pragmatic solution is to *make no assumption* regarding the shape (e.g., smoothness) of the response surface but instead assume that any possible covariance structure, within bounds, is possible. The hyperparameters describing the model of the observations are sampled as extensively as possible and then multiple recommendations for the next point to explore are elicited. To generate the batch of recommendations, the possible recommendations found are clustered using a k-medoids<sup>[73]</sup> clustering approach. The number of clusters matches the size of the batch of recommendations required. This approach provides excellent flexibility and is relatively easy to integrate with any method that already uses Gaussian Process functions to model the objective function. Figure 8 shows a demonstration of Batch BO and an example in which Batch BO schemes have been integrated with a multi-information source fusion framework developed previously to accelerate the design of dual-phase steels.<sup>[55,74]</sup> The preliminary results show that using batch queries, it is possible to dramatically improve the rate of convergence toward a global optimum and reduce the uncertainty of (and increase confidence in) the optimization framework.

## Conclusions and outlook

When resources are limited, it is clear that Bayesian methods offer the best approach to identifying optimal and feasible regions of a materials design space. While off-the-shelf BO approaches can, in principle, be applied to materials discovery and design, very few materials problems can be mapped to the optimization of ‘black box’ functions. In this short overview, it has been shown that there are modifications to BO that are better suited to materials science problems. Combining multiple information sources with high-throughput (HTP) experimental and/or computational platforms could be deployed to solve materials discovery problems with multiple objectives and constraints. While promising, there is much more opportunity for further development.

The performance of BO depends on the adopted surrogate models approximating underlying objective functions, for which Gaussian Processes (GPs) are typically used. However, in many materials discovery applications, the inherent smooth assumptions of GPs may not always hold. Therefore, more adaptive and flexible Bayesian surrogate models in BO may need to be developed to overcome the weaknesses of widely used GP-based methods when faced with relatively high-dimensional design space or non-smooth patterns of objective functions. Recently, it has been shown that Bayesian Multivariate Adaptive Regression Splines (BMARS)<sup>[76]</sup> and Bayesian Additive Regression Trees (BART)<sup>[77]</sup> are flexible stochastic surrogate ML models that can provide better performance when applying BO to complex materials spaces.<sup>[78]</sup>

In the case of multi-objective BO, it is not uncommon for the different objectives to exhibit some degree of correlation. In the case of materials optimization, this makes sense because other properties or performance metrics of a given material are ultimately connected since they most likely share the same underlying process–structure (PS) relationships. Unfortunately,

in traditional multi-objective BO, each objective is modeled individually and such correlations are not considered when making predictions using GPs. Recently, Deep Gaussian Processes (DGP) have been proposed as ML surrogate models that can be used to build multivariate surrogates to catch this correlation and improve the prediction ability of GPs.<sup>[79–81]</sup> A DGP uses multiple layers of GPs to map from a set of features to latent spaces layer by layer to reach a target objective space. This can potentially be useful in representing different objectives in materials discovery as it is probable that a set of properties in a given material system to be correlated to some degree.

Another challenge in the deployment of Bayesian optimization approaches is the ability to handle qualitative data. In real materials discovery tasks, quantitative and qualitative design variables are required to fully define the materials systems. Existing BO approaches mostly work with quantitative data and cannot involve categorical data or catch the correlation between qualitative and quantitative design variables. Recently, some works have proposed employing Latent variable Gaussian Process (LVGP) models to map categorical variables into continuous variable latent spaces.<sup>[82]</sup> Developing mixed variable BO frameworks has been a subject of recent research<sup>[83–88]</sup> and more platforms capable of handling both types of design variables are expected to be developed in the coming years.

A natural extension of the mathematical frameworks presented here is the deployment of automated, closed-loop materials discovery platforms.<sup>[89]</sup> While somewhat limited, the first examples of automated scientific discovery platforms consist of integrated robotic platforms that execute actions that are, in turn, selected using either BO or Reinforcement Learning (RL) approaches. The field of autonomous experimentation is evolving in a highly accelerated manner and such platforms have already been deployed to solve problems involving a wide range of materials classes and applications.<sup>[90–92]</sup>

Further developments in autonomous materials discovery will likely arise from advances in the decision support tools—BO being one of the dominant paradigms due to its flexibility and applicability in the sparse data regime—used to guide the platform toward a particular objective. Currently, for example, human intuition is seldom incorporated into these platforms. Yet, expert human intuition may be highly valuable in guiding the discovery process. A significant challenge, in this case, would be how to algorithmically introduce human expert opinion within a formal closed-loop BO framework. In other contexts, human-in-the-loop BO (HITL-BO) has been implemented by eliciting human experts to express their preference among BO-generated choices.<sup>[93]</sup> Similar approaches could potentially be used in the context of BO-assisted materials discovery. However, there is no work (to the best of the authors' knowledge) that has exploited such ideas as of this writing.

Another critical issue is that human-driven materials discovery tends to be a very dynamic and flexible process. Yet, the BO policies, actions, and experiments to query the materials design space are hard coded. A further limitation is that traditional BO approaches tend to operate under the assumption that

the problem itself (i.e., the definition of objectives, constraints, degrees of freedom) is fixed. However, it is conceivable that a non-trivial material's discovery loop requires significantly more flexibility when defining the problem space itself. For example, performance metrics that were initially considered to be objectives may be more suitable to be considered as constraints, constraints may need to be softened or hardened, or new degrees of freedom may need to be incorporated into the design space.

While there are many more limitations to current BO-based materials discovery platforms, solutions that provide enhanced flexibility and functionality will likely be gradually deployed as the community pursues further autonomy. It is thus likely that future platforms will be able to gather information, develop new knowledge, and evolve into fully autonomous systems. In the near future, further advancements to these autonomous materials discovery platforms will involve innovations in the BO space, robotics, and data integration. While progress has been slow, advancement is accelerating dramatically, and autonomous experimentation in materials science is likely to become the dominant paradigm for exploring complex, multi-objective, and multi-dimensional materials design spaces in the decades to come.

## Acknowledgments

The authors would like to acknowledge the groups of Profs. Ed Dougherty, Dimitris Lagoudas, Ibrahim Karaman, and Seyede Fetemeh Ghoreishi. The authors also acknowledge Drs. Sahin Boluki, Anjana Talapatra, and Jaylen James. The authors acknowledge the financial support from NSF through Grants No. NSF-CMMI-1663130, NSF-CISE-1835690, NSF-CDSE-2001333, and NSF-DMREF-2119103. BV acknowledges the support of NSF through Grant No. NSF-DGE-1545403. The authors also wish to acknowledge the support from the U.S. Department of Energy (DOE) ARPA-E ULTIMATE Program through Project DE-AR0001427. Calculations used to demonstrate many of the BO methods described in this work were carried out at the High-Performance Research Computing (HPRC) facility at Texas A&M University.

## Data availability

Data shown in this paper can be made available upon reasonable request.

## Declarations

## Conflict of interest

The corresponding author states that there is no conflict of interest.

## References

1. G.B. Olson, *Science* **277**(5330), 1237–1242 (1997)

2. C.S. Smith, *A Search for Structure: Selected Essays on Science, Art and History* (MIT Press, Cambridge, MA, 1983)
3. T.M. Pollock, J.E. Allison, D.G. Backman, M.C. Boyce, M. Gersh, E.A. Holm, R. LeSar, M. Long, A. Powell IV, J.J. Schirra, D. Demania Whitis, C. Woodward, *Integrated Computational Materials Engineering: a Transformational Discipline for Improved Competitiveness and National Security* (National Academies Press, Washington, DC, 2008)
4. E.G. David, *Genetic Algorithms in Search, Optimization and Machine Learning* (Addison-Wesley, Reading, 1989)
5. C.M. Fonseca, P.J. Fleming, Genetic algorithms for multiobjective optimization: formulation, discussion and generalization. in *Fifth International Conference on Genetic Algorithms*, (Morgan Kaufmann, San Mateo, 1993)
6. P. Voorhees, G. Spanos, *Modeling Across Scales: A Roadmapping Study for Connecting Materials Models and Simulations Across Length and Time Scales* (TMS, Warrendale, PA, 2015)
7. J.H. Panchal, S.R. Kalidindi, D.L. McDowell, *Comput. Aided Des.* **45**(1), 4–25 (2013)
8. National Science and Technology Council, *Materials Genome Initiative for Global Competitiveness* (Executive Office of the President, Washington, DC, 2011), pp. 1–18
9. J.R. Engstrom, W.H. Weinberg, *AIChE J.* **46**(1), 2–5 (2000)
10. S. Curtarolo, G.L. Hart, M.B. Nardelli, N. Mingo, S. Sanvito, O. Levy, *Nat. Mater.* **12**(3), 191 (2013)
11. J. Allison, D. Backman, L. Christodoulou, *Jom* **58**(11), 25–27 (2006)
12. J. Allison, *Jom* **63**(4), 15 (2011)
13. H. Koinuma, I. Takeuchi, *Nat. Mater.* **3**(7), 429 (2004)
14. Z. Xiong, Y. He, J.R. Hattrick-Simpers, J. Hu, *ACS Comb. Sci.* **19**(3), 137–144 (2017)
15. A. Talapatra, S. Boluki, P. Honarmandi, A. Solomou, G. Zhao, S.F. Ghoreishi, A. Molkeri, D. Allaire, A. Srivastava, X. Qian, *Front. Mater.* **6**, 82 (2019)
16. E. Brochu, V.M. Cora, N. De Freitas, *arXiv preprint arXiv:1012.2599* (2010)
17. R. Astudillo, P. Frazier, Bayesian optimization of composite functions, in *Proceedings of the 36th International Conference on Machine Learning*, ed. by C. Kamalika, S. Ruslan, (PMLR: Proceedings of Machine Learning Research, 2019), pp. 354–363
18. A. Tran, J. Tranchida, T. Wildey, A.P. Thompson, *J. Chem. Phys.* **153**(7), 074705 (2020)
19. Q. Liang, A.E. Gongora, Z. Ren, A. Tihihonen, Z. Liu, S. Sun, J.R. Deneault, D. Bash, F. Mekki-Berrada, S.A. Khan, K. Hippalgaonkar, B. Maruyama, K.A. Brown, J. Fisher Iii, T. Buonassisi, *Npj Comput. Mater.* **7**(1), 188 (2021)
20. P.I. Frazier, *preprint arXiv:arXiv:1807.02811* (2018)
21. R. Arróyave, D.L. McDowell, *Annu. Rev. Mater. Res.* **49**(1), 103–126 (2019)
22. C. Sharpe, C.C. Seepersad, S. Watts, D. Tortorelli, Design of mechanical metamaterials via constrained Bayesian optimization, in *ASME 2018 International Design Engineering Technical Conferences and Computers and Information in Engineering Conference* (2018)
23. H. Wahab, V. Jain, A.S. Tyrrell, M.A. Seas, L. Kotthoff, P.A. Johnson, *Carbon* **167**, 609–619 (2020)
24. Y. Xie, C. Zhang, H. Deng, B. Zheng, J.-W. Su, K. Shutt, J. Lin, *ACS Appl. Mater. Interfaces* **13**(45), 53485–53491 (2021)
25. J.K. Pedersen, C.M. Clausen, O.A. Krysiak, B. Xiao, T.A.A. Batchelor, T. Löffler, V.A. Mints, L. Banko, M. Arenz, A. Sazan, W. Schuhmann, A. Ludwig, J. Rossmeisl, *Angew. Chem. Int. Ed.* **60**(45), 24144–24152 (2021)
26. S. Greenhill, S. Rana, S. Gupta, P. Vellanki, S. Venkatesh, *IEEE Access* **8**, 13937–13948 (2020)
27. G. Agarwal, H.A. Doan, L.A. Robertson, L. Zhang, R.S. Assary, *Chem. Mater.* **33**(20), 8133–8144 (2021)
28. J. Peng, J.K. Damewood, J. Karaguesian, R. Gómez-Bombarelli, Y. Shao-Horn, *Joule* **5**(12), 3069–3071 (2021)
29. D. Xue, P.V. Balachandran, J. Hogden, J. Theiler, D. Xue, T. Lookman, *Nat. Commun.* **7**, 11241 (2016)
30. D. Xue, P.V. Balachandran, R. Yuan, T. Hu, X. Qian, E.R. Dougherty, T. Lookman, *Proc. Natl. Acad. Sci. USA* **113**(47), 13301–13306 (2016)
31. I.J. Good, *The Estimation Of Probabilities: An Essay on Modern Bayesian Methods* (The MIT Press, Cambridge, MA, 1965)
32. R. Garnett, *Bayesian Optimization* (Cambridge University Press, Cambridge, 2022)
33. B. Shahriari, K. Swersky, Z. Wang, R.P. Adams, N. De Freitas, *Proc. IEEE* **104**(1), 148–175 (2015)
34. D.R. Jones, M. Schonlau, W.J. Welch, *J. Global Optim.* **13**(4), 455–492 (1998)
35. C. Williams, C. Rasmussen, *Adv. Neural Inf. Process. Syst.* **8**, 514–520 (1995)
36. R. Planas, N. Oune, R. Bostanabad, *J. Mech. Des.* **143**(11), 111703 (2021)
37. J.T. Eweis-Labolle, N. Oune, R. Bostanabad, *J. Mech. Des.* **144**(9), 091703 (2022)
38. A. Talapatra, S. Boluki, T. Duong, X. Qian, E. Dougherty, R. Arróyave, *Phys. Rev. Mater.* **2**(11), 113803 (2018)
39. R. Yuan, Z. Liu, P.V. Balachandran, D. Xue, Y. Zhou, X. Ding, J. Sun, D. Xue, T. Lookman, *Adv. Mater.* **30**(7), 1702884 (2018)
40. K. Tran, Z.W. Ulissi, *Nat. Catal.* **1**(9), 696–703 (2018)
41. D. Xue, D. Xue, R. Yuan, Y. Zhou, P.V. Balachandran, X. Ding, J. Sun, T. Lookman, *Acta Mater.* **125**, 532–541 (2017)
42. E. Raponi, H. Wang, M. Bujny, S. Boria, C. Doerr, High dimensional Bayesian optimization assisted by principal component analysis, in *International Conference on Parallel Problem Solving from Nature*, (Springer, Cham, 2020)
43. A. Talapatra, S. Boluki, T. Duong, X. Qian, E. Dougherty, R. Arróyave, *arXiv preprint arXiv:1803.05460*, (2018)
44. L. Wasserman, *J. Math. Psychol.* **44**(1), 92–107 (2000)
45. D. Khatamsaz, A. Molkeri, R. Couperthwaite, J. James, R. Arróyave, A. Srivastava, D. Allaire, *Mater. Des.* **209**, 110001 (2021)
46. G. Malkomes, R. Garnett, Automating Bayesian optimization with Bayesian optimization, in *Advances in Neural Information Processing Systems* (2018)
47. G. Malkomes, C. Schaff, R. Garnett, *Adv. Neural Inf. Process. Syst.* **29**, 1–8 (2016)
48. L. Schlessinger, G. Malkomes, R. Garnett, Automated model search using Bayesian optimization and genetic programming, in *Workshop on Meta-Learning at Advances in Neural Information Processing Systems* (2019)
49. D. Duvenaud, J. Lloyd, R. Grosse, J. Tenenbaum, G. Zoubin, Structure discovery in nonparametric regression through compositional kernel search, in *International Conference on Machine Learning* (PMLR, 2013)
50. S.F. Ghoreishi, S. Friedman, D.L. Allaire, *J. Mech. Des.* **141**(7), 071404 (2019)
51. M.G. Fernández-Godino, C. Park, N.-H. Kim, R.T. Haftka, *arXiv preprint arXiv:1609.07196*, (2016)
52. S. Venkatram, R. Batra, L. Chen, C. Kim, M. Shelton, R. Ramprasad, *J. Phys. Chem. B* **124**(28), 6046–6054 (2020)
53. G. Pilania, J.E. Gubernatis, T. Lookman, *Comput. Mater. Sci.* **129**, 156–163 (2017)
54. H.C. Herbol, M. Poloczec, P. Clancy, *Mater. Horiz.* **7**(8), 2113–2123 (2020)
55. S.F. Ghoreishi, A. Molkeri, A. Srivastava, R. Arroyave, D. Allaire, *J. Mech. Des.* **140**(11), 111409 (2018)
56. D. Khatamsaz, A. Molkeri, R. Couperthwaite, J. James, R. Arróyave, D. Allaire, A. Srivastava, *Acta Mater.* **206**, 116619 (2021)
57. D. Allaire, K. Willcox, Fusing information from multifidelity computer models of physical systems, in *2012 15th International Conference on Information Fusion* (IEEE, 2012)
58. S. Chen, Z. Jiang, S. Yang, W. Chen, *AIAA J.* **55**(1), 241–254 (2016)
59. H. Wang, M. Olhofer, Y. Jin, *Complex Intell. Syst.* **3**(4), 233–245 (2017)
60. P. Ngatchou, A. Zarei, A. El-Sharkawi, Pareto multi objective optimization, in *Proceedings of the 13th International Conference on Intelligent Systems Application to Power Systems* (IEEE, 2005)
61. M.T. Emmerich, A.H. Deutz, J.W. Klinkenberg, Hypervolume-based expected improvement: Monotonicity properties and exact computation, in *2011 IEEE Congress of Evolutionary Computation (CEC)*, (IEEE, 2011)
62. A. Solomou, G. Zhao, S. Boluki, J.K. Joy, X. Qian, I. Karaman, R. Arróyave, D.C. Lagoudas, *Mater. Des.* **160**, 810–827 (2018)
63. R. Arróyave, S. Gibbons, E. Galvan, R. Malak, *JOM* **68**(5), 1385–1395 (2016)
64. D. Khatamsaz, B. Vela, P. Singh, D.D. Johnson, D. Allaire, R. Arroyave, *Acta Mater.* **236**, 118133 (2022)
65. A. Jain, G. Hautier, C.J. Moore, S.P. Ong, C.C. Fischer, T. Mueller, K.A. Persson, G. Ceder, *Comput. Mater. Sci.* **50**(8), 2295–2310 (2011)
66. S. Curtarolo, W. Setyawan, G.L. Hart, M. Jahnatek, R.V. Chepulskii, R.H. Taylor, S. Wang, J. Xue, K. Yang, O. Levy, *Comput. Mater. Sci.* **58**, 218–226 (2012)



67. B.R. Ortiz, J.M. Adamczyk, K. Gordiz, T. Braden, E.S. Toberer, *Mol. Syst. Des. Eng.* **4**(2), 407–420 (2019)
68. M. Moorehead, K. Bertsch, M. Niezgoda, C. Parkin, M. Elbakhshwan, K. Sridharan, C. Zhang, D. Thoma, A. Couet, *Mater. Des.* **187**, 108358 (2020)
69. R.A. Potyrailo, I. Takeuchi, *Meas. Sci. Technol.* **16**(1), 1 (2004)
70. A. Ludwig, *Npj Comput. Mater.* **5**(1), 70 (2019)
71. T. Oellers, V.G. Arigela, C. Kirchlechner, G. Dehm, A. Ludwig, *ACS Comb. Sci.* **22**(3), 142–149 (2020)
72. T.T. Joy, S. Rana, S. Gupta, S. Venkatesh, *Knowl. Based Syst.* **187**, 104818 (2020)
73. L. Kaufman, P. Rousseeuw, *Statistical Data Analysis Based on the L1-Norm and Related Methods* (North Holland, Amsterdam, 1987)
74. S.F. Ghoreishi, A. Molkeri, R. Arróyave, D. Allaire, A. Srivastava, *Acta Mater.* **180**, 260–271 (2019)
75. R. Couperthwaite, A. Molkeri, D. Khatamsaz, A. Srivastava, D. Allaire, R. Arróyave, *JOM* **72**(12), 4431–4443 (2020)
76. J.H. Friedman, *Ann. Stat.* **19**(1), 1–67 (1991)
77. H.A. Chipman, E.I. George, R.E. McCulloch, *Ann. Appl. Stat.* **4**(1), 266–298 (2010)
78. B. Lei, T.Q. Kirk, A. Bhattacharya, D. Pati, X. Qian, R. Arroyave, B.K. Mallick, *Npj Comput. Mater.* **7**(1), 1–12 (2021)
79. A. Hebbal, M. Balesdent, L. Brevault, N. Melab, E.-G. Talbi, *Optim. Eng.* (2022). <https://doi.org/10.1007/s11081-022-09753-0>
80. A. Sauer, R.B. Gramacy, D. Higdon, *Technometrics* (2022). <https://doi.org/10.1080/00401706.2021.2008505>
81. A. Damianou, N.D. Lawrence, Deep Gaussian processes, in *Proceedings of the Sixteenth International Conference on Artificial Intelligence and Statistics*, ed. by M.C. Carlos, R. Pradeep, (PMLR: Proceedings of Machine Learning Research, 2013), pp. 207–215
82. Y. Zhang, S. Tao, W.W. Chen, D.W. Apley, *Technometrics* **62**(3), 291–302 (2020)
83. H. Zhang, W. Chen, A. Iyer, D.W. Apley, W. Chen, arXiv preprint [arXiv:2207.04994](https://arxiv.org/abs/2207.04994), (2022)
84. Y. Zhang, D.W. Apley, W. Chen, *Sci. Rep.* **10**(1), 1–13 (2020)
85. H. Xu, C.-H. Chuang, R.-J. Yang, Mixed-variable metamodeling methods for designing multi-material structures. in *ASME 2016 International Design Engineering Technical Conferences and Computers and Information in Engineering Conference* (2016)
86. J. Pelamatti, L. Brevault, M. Balesdent, E.-G. Talbi, Y. Guerin, *J. Glob. Optim.* **73**(3), 583–613 (2019)
87. J.A. Manson, T.W. Chamberlain, R.A. Bourne, *J. Glob. Optim.* **80**(4), 865–886 (2021)
88. A. Iyer, Y. Zhang, A. Prasad, S. Tao, Y. Wang, L. Schadler, L.C. Brinson, W. Chen, Data-centric mixed-variable bayesian optimization for materials design, in *ASME 2019 International Design Engineering Technical Conferences and Computers and Information in Engineering Conference* (2019)
89. F. Häse, L.M. Roch, A. Aspuru-Guzik, *Trends Chem.* **1**(3), 282–291 (2019)
90. E. Stach, B. DeCost, A.G. Kusne, J. Hattrick-Simpers, K.A. Brown, K.G. Reyes, J. Schrier, S. Billinge, T. Buonassisi, I. Foster, *Matter* **4**(9), 2702–2726 (2021)
91. M.M. Flores-Leonar, L.M. Mejía-Mendoza, A. Aguilar-Granda, B. Sanchez-Lengeling, H. Tribukait, C. Amador-Bedolla, A. Aspuru-Guzik, *Curr. Opin. Green Sustain. Chem.* **25**, 100370 (2020)
92. D.P. Tabor, L.M. Roch, S.K. Saikin, C. Kreisbeck, D. Sheberla, J.H. Montoya, S. Dwaraknath, M. Aykol, C. Ortiz, H. Tribukait, *Nat. Rev. Mater.* **3**(5), 5–20 (2018)
93. Y. Zhou, Y. Koyama, M. Goto, T. Igarashi, arXiv preprint [arXiv:2010.03190](https://arxiv.org/abs/2010.03190), (2020)

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.