

Bayesian optimization of nanoporous materials

Aryan Deshwal¹, Cory M. Simon^{‡2}, and Janardhan Rao Doppa^{†1}

¹School of Electrical Engineering and Computer Science, Washington State University. Pullman, WA.

²School of Chemical, Biological, and Environmental Engineering. Oregon State University. Corvallis, OR.

[†]jana.doppa@wsu.edu, [‡]cory.simon@oregonstate.edu

July 10, 2021

Abstract

Nanoporous materials (NPMs) could be used to store, capture, and sense many different gases. Given an adsorption task, we often wish to search a library of NPMs for the one with the optimal adsorption property. The high cost of NPM synthesis and gas adsorption measurements, whether these experiments are in the lab or in a simulation, often precludes exhaustive search.

We explain, demonstrate, and advocate Bayesian optimization (BO) to actively search for the optimal NPM in a library of NPMs— and find it using the fewest experiments. The two ingredients of BO are a surrogate model and an acquisition function. The surrogate model is a probabilistic model reflecting our beliefs about the NPM-structure-property relationship based on observations from past experiments. The acquisition function uses the surrogate model to score each NPM according to the utility of picking it for the next experiment. It balances two competing goals: (a) exploitation of our current approximation of the structure-property relationship to pick the highest-performing NPM, and (b) exploration of blind spots in the NPM space to pick an NPM we are uncertain about, to improve our approximation of the structure-property relationship. We demonstrate BO by searching an open database of $\sim 70\,000$ hypothetical covalent organic frameworks (COFs) for the COF with the highest simulated methane deliverable capacity. BO finds the optimal COF and acquires 30% of the top 100 highest-ranked COFs after evaluating only ~ 120 COFs. More, BO searches more efficiently than evolutionary and one-shot supervised machine learning approaches.

1 Introduction

The selective gas adsorption properties of nanoporous materials (NPMs) endow them with many possible applications in the storage [1, 2], separation [3], and sensing [4] of gases. As examples, promising applications of NPMs include (i, storage) densifying hydrogen (H_2)—a clean fuel—for compact storage onboard vehicles [2, 5], (ii, separation) capturing carbon dioxide from flue gas of coal-fired power plants; subsequently sequester it to prevent global warming [6], and (iii, sensing) detecting toxic compounds and explosives [7, 8].

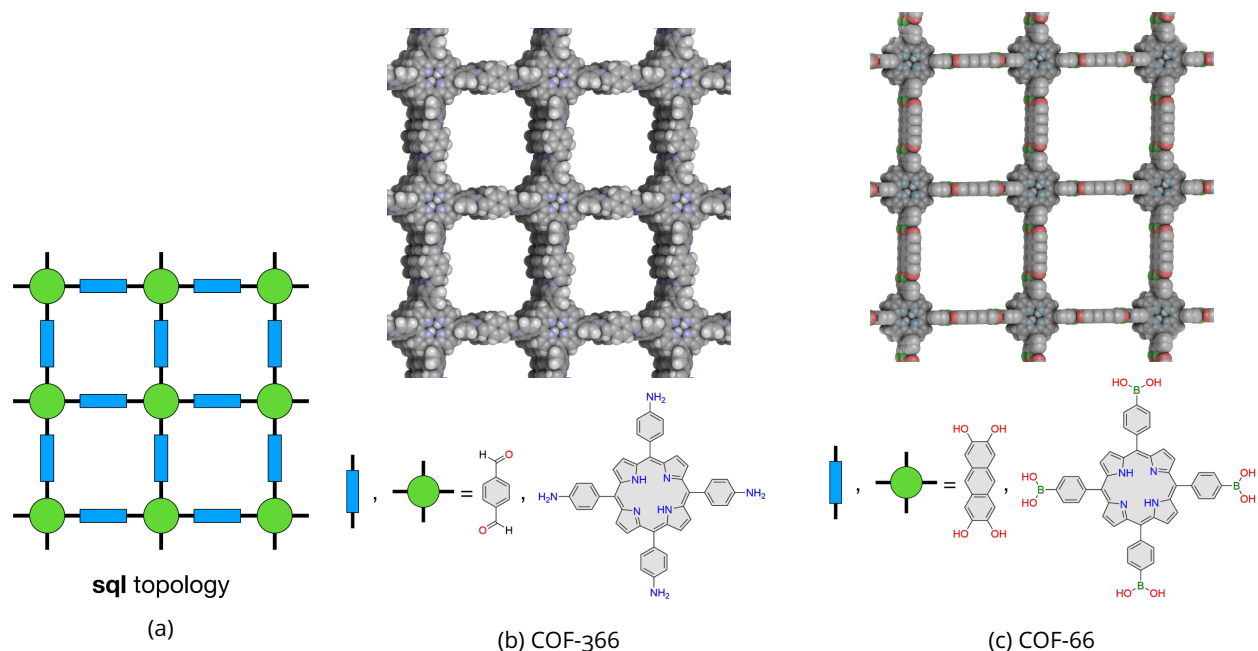


Figure 1: Illustrating the modular synthesis of covalent organic frameworks (COFs) [9–11]. (a) The **sqf** (square lattice) network topology specifies the connectivity of tetratopic (green) and diptopic (blue) building units. (b, c) Two examples of COFs in the **sqf** topology are (b) COF-366 and (c) COF-66 [12]. (top) the crystal structures. (bottom) the building blocks: a planar, tetratopic building unit and a linear, diptopic building unit. The building units are stitched together with covalent bonds, through a condensation reaction, to form 2D COF sheets in the **sqf** topology. The sheets stack into layers to form 3D channels. Owing to the rational design and modular synthesis of COFs, on the order of one hundred COFs have been experimentally synthesized and reported [13].

Several classes of NPMs, such as metal-organic frameworks (MOFs) [14], metal-organic polyhedra (MOPs) [15], covalent organic frameworks (COFs) [9], and porous organic cages (POCs) [16], are synthesized modularly by stitching together molecular building blocks via coordination (MOFs, MOPs) or covalent (COFs, POCs) bonds to form ~ crystalline (MOFs, COFs) or molecular (MOPs, POCs) materials. Fig. 1 illustrates the modular synthesis and rational design of COFs in the **sqf** topology [12]. The many topologies [9, 17, 18], abundance of molecular building blocks, and post-synthetic modifiability [19, 20] permit an unlimited number of possible structures exhibiting diverse adsorption properties.

A common goal is to find, among a large set of candidate NPM structures, the NPM structure(s) with the optimal adsorption property for a given application. As opposed to an exhaustive search, our goal is to search for the optimal NPM *efficiently*, by consuming minimal resources (computational and/or physical) in the process. In the laboratory setting, synthesizing an NPM and measuring its property costs labor and raw materials, and throughput is limited by the capital equipment in the lab. In the computational setting, constructing a high-fidelity computational model of an NPM structure [21–25] and predicting its gas adsorption property through molecular simulations [26, 27] consumes electricity, and throughput is limited by computing resources. Thus, both in the bona fide laboratory

and on the computer, our goal is to find the optimal NPM(s) for a given adsorption task using the fewest experiments (experiment = constructing an NPM and evaluating its adsorption property).

In this article, we explain, demonstrate, and advocate Bayesian optimization (BO) to *actively* and *efficiently* search for NPMs with an optimal property for a given adsorption task. Active, because each BO iteration performs three consecutive steps: (a) conducting an experiment on an NPM, (b) updating our belief about the structure-property relationship, and (c) selecting the NPM for the next experiment. See Fig. 2 for a high-level illustration. Efficient, because BO makes a *data-informed* decision on the NPM to select for the next experiment, while balancing: (a) *exploitation* of our current (uncertain!) data-driven belief about the structure-property relationship to pick an NPM that might have the optimal property and (b) *exploration* of regions of the NPM design space where our belief about the structure-property relationship is weak to pick an NPM we are uncertain about.

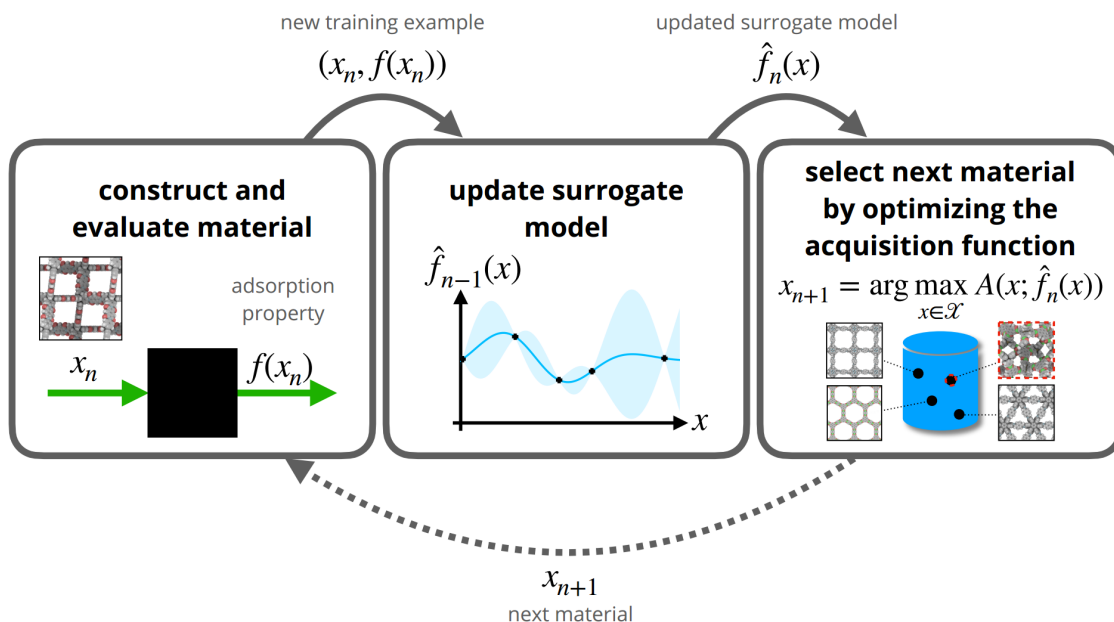


Figure 2: Illustration of a single iteration of active search in Bayesian optimization (BO) of nanoporous materials. First, we conduct an experiment that measures the property $f(x_n)$ of the material represented by x_n . Second, using the new observation, we update our belief about the underlying objective function $f(x)$, reflected by the surrogate model. Third, we use the acquisition function to pick the next material for an experiment. BO is an active search method to find the material x that optimizes the black-box function $f(x)$.

The two key components of BO are a surrogate model and an acquisition function. The surrogate model, with “surrogate” hinting at “substitute for the experiment”- is a probabilistic model for the structure-property relationship. It is trained on all observations from past experiments. The surrogate model cheaply predicts the properties of the unevaluated NPMs and, critically, quantifies the uncertainty in its predictions. The acquisition function is used to make the decision of which NPM to select for the next experiment. It uses the surrogate model to score the utility of selecting each unevaluated NPM for the next experiment by striking a balance in the exploitation-exploration trade-off. The acquisition function is maximal in regions of the NPM design space where (i) we be-

lieve high-performing materials will reside and/or (ii) we are uncertain about the structure-property relationship. The acquisition function relies on the surrogate model to give a resource-efficient, intelligent, active search for optimal NPMs in a wide variety of contexts.

We demonstrate BO of NPMs by efficiently searching an open database [28] of $\sim 70\,000$ hypothetical COFs, represented by vectors of hand-designed features based on domain knowledge, for those with the highest simulated methane deliverable capacity to store natural gas onboard vehicles [2]. BO recovers the optimal (COFs) using fewer experiments than incumbent strategies including random search, evolutionary algorithms, and one-shot supervised machine learning. In the outlook, we discuss several active research areas in BO that are likely to apply to several problems in NPM discovery: (i) batch BO where experiments are parallelized, (ii) multi-fidelity BO, where NPMs can be evaluated using multiple methods which vary in accuracy and resource cost, (iii) multi-objective BO, where we aim to find a Pareto optimal set of NPMs to optimize multiple properties, and (iv) constrained BO, where our goal is to find high-performing NPMs which can be synthesized.

2 Review of previously used NPM search methods

The NPM research community has adopted several approaches to efficiently search a library of NPMs for the optimal NPM(s) [29,30]. We define the efficiency of a search strategy with respect to two naive baselines: (i) exhaustive search, where we conduct the high-fidelity experiment on every NPM in the library, and (ii) random search, where we conduct the high-fidelity experiment on a (uniform) random sample of the NPMs in the library.

Supervised machine learning models. A supervised machine learning model can serve as a cheaper, albeit lower fidelity, surrogate for the high-fidelity experiment [29,31–35], thereby reducing the cost of an exhaustive search. A machine learning approach is predicated upon cheaply computed (relative to the experiment) (i) vector representations of the NPMs—hand-engineered [36,37] or learned [38]—that encode structural features and are correlated to the property or (ii) kernels that capture the tendency for any given pair of NPMs to exhibit similar properties [39]. Training examples are gathered by selecting a small (random or diverse [40,41]) subset of the library of NPMs and labeling them with the property values using the high-fidelity experiment. Using the training examples, the supervised machine learning model learns to predict the property of any given NPM from e.g., its vector representation. The trained model is then used, as a surrogate for the high-fidelity experiment, to cheaply predict, from their vector representations, the properties of the remaining NPMs in the library. Further high-fidelity experiments may be directed on the NPMs predicted to be optimal by the machine learning model. See Refs. [40,42–51] as examples.

Genetic algorithms. Genetic algorithms [52] are iterative, stochastic search methods inspired by Darwinian evolution. Each NPM is represented by a “chromosome”—a vector of categorical variables that uniquely specifies its structure. A small initial generation (set) of property-labeled NPMs is iteratively evolved by applying genetic operations to their chromosomes: mutation, replication, selection, and recombination. At each generation, we conduct experiments on each newly evolved NPM to evaluate its fitness. This guides the genetic operations used to evolve to the next generation of chromosomes (representing NPMs), with the ambition of both exploring NPM space and enriching future generations with high-fitness NPMs. See Refs. [53–55] as examples.

Monte Carlo tree search. When the NPM search space can be framed as a tree, Monte Carlo tree search [56] is more efficient than random search. Starting at the root node, a policy to select a child node is iteratively applied, giving a path through the tree, culminating at a leaf node. The experiment is then conducted on the NPM represented by the leaf node. Its measured property, viewed as a reward, is back-propagated through the tree to update the statistics of each node along the path to it. Both the visit counts and reward allocations of the nodes are used in the tree policy to balance exploration of new branches of the tree that have not been visited often (or at all) and exploitation of current knowledge to follow branches of the tree that appear likely to lead to optimal NPMs. See Refs. [57,58] as examples.

Each of these prior approaches suffer from drawbacks. The supervised machine learning approach selects training data to learn a good predictor of the property from the structure representation, then uses the predictor to greedily acquire the COFs with the highest predicted property. This passive approach can be viewed as one round of exploration and one round of exploitation. Active learning [59] can be used to reduce the number of training examples to learn the structure-property relationship but is not geared towards finding the optimal NPM using the fewest experiments. In BO, we will actively collect training examples according to the goal of finding the optimal NPM with the fewest experiments. Genetic algorithms are sequential, active search procedures aimed at quickly finding the optimal NPM. However, the genetic operations are heuristic and do not balance exploration and exploitation rigorously. As a consequence, genetic algorithms can be difficult to tune and could require many experiments to find the optimal NPMs. MCTS balances exploration and exploitation more rigorously. But, it requires many NPM evaluations to identify promising regions of the tree because it does not explicitly leverage the similarity among structures for principled exploitation. Further, MCTS is limited to NPM design spaces which can be framed as a tree.

3 Problem setup: find the optimal material

Suppose we have a large database of candidate NPM structures, \mathcal{X} , for some adsorption task. Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be a black-box objective function that, given a candidate NPM $\mathbf{x} \in \mathcal{X}$, returns the relevant adsorption property $y = f(\mathbf{x})$. Each evaluation of f corresponds to performing an *expensive* experiment—either in the laboratory or in a molecular simulation—to measure the adsorption property y of NPM \mathbf{x} . Our goal is to find the highest-performing NPM \mathbf{x}^* from \mathcal{X} that maximizes the objective function f ,

$$\mathbf{x}^* = \arg \max_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}), \quad (1)$$

while conducting the fewest number of expensive experiments.

We can interpret $f(\mathbf{x})$ as the [unknown] structure-property relationship [60,61] since \mathbf{x} [abstractly, at this point in our discussion] represents the structure of a unique NPM and evaluating f means conducting an experiment to measure its property, y .

4 Overview of Bayesian Optimization

We explain the key ideas behind the Bayesian Optimization (BO) framework [62, 63] to find the highest-performing NPM—solve the problem in eqn. 1—*efficiently*, by using the fewest (expensive) experiments.

4.1 Defining an NPM feature space

While \mathbf{x} as an abstract representation of an NPM suffices for defining the problem in eqn. 1, for BO we must concretely define a NPM *feature space* or *search space* in which each NPM \mathbf{x} lies.

Without loss of generality, take \mathbf{x} to be a fixed-size (among all NPMs) vector representation of the NPM that lies in a continuous feature space. The NPM feature vectors $\{\mathbf{x}\}$ should be designed to (1) encode the relevant structural and chemical features of the NPMs, (2) be rotation-, translation-, and, if the NPM is a crystal, replication-invariant, (3) be cheap to compute compared to conducting the experiment (evaluating f), and (4) ideally, each correspond to a unique NPM (mapping NPM $\rightarrow \mathbf{x}$ is injective) [64]. As a result, NPMs close in the feature space should exhibit close values of the property y .

The simplest example of a representation \mathbf{x} of an NPM is a list of hand-designed, based on domain expertise, descriptors/features of its structure and composition, such as pore volume, surface area, largest included sphere diameter, density, weight fraction carbon, etc. [65–67]. Alternatively, \mathbf{x} could be learned from a graph [68–70] or 3D image representation [71, 72] of a NPM by a graph neural network [38] or convolutional neural network [73], respectively. See reviews in Refs. [29, 74] for defining feature vectors of NPMs and Refs. [40, 42, 44, 75–78] for different examples of NPM feature spaces. As opposed to dwelling on how to *define* a good feature space of NPMs, we will instead focus on BO, a technique to *search* the defined feature space for the optimal NPM \mathbf{x}^* in an efficient manner.

4.2 Active search: exploitation vs. exploration

Even with an NPM feature space defined, in practice, the structure-property relationship $f(\mathbf{x})$ is a black-box function; analytical expressions for $f(\mathbf{x})$ and/or its gradient $\nabla_{\mathbf{x}}f(\mathbf{x})$ are not known, and it may be multi-modal.

BO is a derivative-free method to *actively* and *efficiently* search the database of NPMs \mathcal{X} for the NPM \mathbf{x}^* that maximizes $f(\mathbf{x})$. Active, because BO sequentially selects NPMs from \mathcal{X} for experimentation (to evaluate with f), iterating between conducting an experiment and making a decision about which experiment to conduct next. Efficient, because BO makes a data-driven decision to select the next NPM for an experiment while taking into account all observed (NPM \mathbf{x} , property $y = f(\mathbf{x})$) pairs from previous experiments. Each decision to select the next unevaluated NPM from \mathcal{X} to evaluate with f must trade off two conflicting goals:

- 1) *Exploitation* suggests to use our current, but uncertain, approximation of the structure-property relationship, based on the past observations, to select the NPM that appears to have the most promise as a high-performing material.
- 2) *Exploration* suggests to select the NPM that we are most uncertain about to improve our approxi-

mation of the structure-property relationship.

So, to balance exploitation and exploration, we must balance visits to regions of NPM feature space that (i) appear to contain high-performing NPMs and (ii) have not been explored well. A colloquial example of the exploitation-exploration dilemma in our lives is, in aiming to maximize our enjoyment of food, whether to choose a restaurant that we have visited and know we like versus a new one [79].

4.3 The ingredients of BO for data-driven decision-making: a surrogate model and an acquisition function

In the BO framework, the two key components used to make each sequential decision of which NPM to conduct an experiment on next are (1) a *surrogate model* that captures our beliefs, based on past observations, about the structure-property relationship and (2) an *acquisition function* that scores each NPM according to the utility of conducting the experiment on it next. The acquisition function uses the surrogate model of the structure property-relationship $f(\mathbf{x})$ to decide which NPM to evaluate next while striking a balance between exploitation and exploration.

The surrogate model. The surrogate model is a probabilistic model of the structure-property relationship $f(\mathbf{x})$ trained on all observations¹ $\{(\mathbf{x}_i, y_i = f(\mathbf{x}_i))\}_{i=1}^n$ from past experiments. Typically, adopting a Bayesian perspective, the surrogate model treats $f(\mathbf{x})$ as a random variable that follows a Gaussian distribution:

$$f(\mathbf{x}) \sim \mathcal{N}(\hat{y}(\mathbf{x}), \sigma^2(\mathbf{x})) \quad (2)$$

with mean $\hat{y} \in \mathbb{R}$ and variance $\sigma^2 \in \mathbb{R}$. The surrogate model reflects our current beliefs about $f(\mathbf{x})$ and serves two purposes in BO. First, to guide exploitation, $\hat{y}(\mathbf{x})$ *cheaply* estimates the properties of the remaining, unevaluated NPMs, i.e., $\hat{y}(\mathbf{x})$ is a cheap-to-evaluate approximation of the expensive-to-evaluate objective function $f(\mathbf{x})$. Second, to guide exploration, $\sigma^2(\mathbf{x})$ quantifies the uncertainties in the predicted properties of the unevaluated NPMs. This makes us aware of “blind spots” of the surrogate model—regions in the NPM feature space we need to explore to improve our approximation $\hat{y}(\mathbf{x})$ and reduce the uncertainty in our beliefs about $f(\mathbf{x})$.

The surrogate model is updated in each iteration of BO, after the new observation $(\mathbf{x}_{n+1}, y_{n+1} = f(\mathbf{x}_{n+1}))$ is gathered, to (i) improve the approximation of $f(\mathbf{x})$ and (ii) account for the reduced uncertainty in the region of the feature space surrounding the newly evaluated NPM \mathbf{x}_{n+1} . Consequently, let us denote the surrogate model after iteration n of BO as $\hat{f}_n : \mathbf{x} \mapsto (\hat{y}, \sigma)$.

The acquisition function. The *acquisition function* $A(\mathbf{x}; \hat{f}(\mathbf{x})) : \mathcal{X} \rightarrow \mathbb{R}$ scores the utility of, next, evaluating NPM $\mathbf{x} \in \mathcal{X}$ with the expensive objective function f . Here, “utility” is defined in terms of our ultimate goal of finding the optimal NPM \mathbf{x}^* in eqn. 1 with the fewest experiments. The acquisition function employs the prediction of the property \hat{y} and the associated uncertainty σ^2 from the surrogate model $\hat{f}(\mathbf{x})$ to assign a utility score to the NPM that balances exploitation and exploration, respectively. Maxima of the acquisition function are located in regions of NPM feature space where the predicted property is large and/or uncertainty is high.

¹Without loss of generality, the observations are assumed noise-less for clarity of presentation.

The decision of which NPM to evaluate next is made by maximizing the acquisition function:

$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x} \in \mathcal{X} \setminus \mathcal{X}_n} A(\mathbf{x}; \hat{f}_n(\mathbf{x})), \quad (3)$$

where $\mathcal{X}_n := \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ is the acquired set of n NPMs that have been evaluated already. Importantly, the acquisition function must be cheap to evaluate and optimize (if \mathcal{X} is not a finite set).

4.4 Summarizing: BO active search iterations

Fig. 2 illustrates an iteration of BO. At the beginning of iteration n , we conduct an experiment on NPM \mathbf{x}_n , i.e., we evaluate NPM \mathbf{x}_n with the objective function f to obtain a new observation $(\mathbf{x}_n, y_n = f(\mathbf{x}_n))$. Next, we update the old surrogate model $\hat{f}_{n-1}(\mathbf{x})$ to account for this new observation, giving the new surrogate model $\hat{f}_n(\mathbf{x})$. We then select the next (unevaluated) NPM to evaluate, \mathbf{x}_{n+1} , as the one that maximizes the acquisition function $A(\mathbf{x}; \hat{f}_n(\mathbf{x}))$.

We terminate BO after we either (i) expend our budget for experiments or (ii) find a material with a satisfactory property value. The BO solution to the problem in eqn.1, \mathbf{x}^* , then follows from the evaluated NPM with the highest observed property, $\arg \max_{i=1}^N y_i$, where N is the number of BO iterations (=experiments) performed. Some theoretical work focuses on characterizing how, under specific assumptions, the quality of the approximate optimum in BO scales with the number of iterations [80].

N.b., typically, the surrogate model is retrained from scratch after each iteration of BO, but some surrogate models can be trained online [81], reducing the computational cost of the search.

4.5 Remark: BO vs. active learning

We remark on a distinction between active learning [82] and Bayesian optimization. Both sequentially collect training examples for a supervised machine learning model. In active learning, the examples are efficiently collected with the goal of reducing the uncertainty in the machine learning model. In Bayesian optimization, the examples are efficiently collected with the goal of, instead, finding the optimal material. BO is more efficient for finding the optimal material than active learning because it avoids collecting examples in regions of feature space that contain poor-performing materials, whereas active learning will do so to reduce the uncertainty of the model in those regions.

5 Surrogate models and acquisition functions

In this section, we explain surrogate models and acquisition functions commonly used in BO.

5.1 Surrogate models: Gaussian processes

Gaussian processes (GPs) [83, 84] are the most commonly used surrogate models in BO owing to their flexibility as function approximators, principled uncertainty quantification, and compatibility with the kernel trick. GPs are non-parametric models that can approximate complicated objective

functions $f(\mathbf{x})$ given labeled training data $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$. Through a Bayesian probabilistic framework, GPs provide uncertainty estimates in their predictions and allow incorporation of prior beliefs. GPs rely on a kernel function $k(\mathbf{x}, \mathbf{x}') : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ [85] to capture the similarity between any two NPMs \mathbf{x} and \mathbf{x}' . This gives GPs the flexibility to approximate arbitrary, complicated (but well-behaved!) functions $f(\mathbf{x})$. Moreover, it gives GPs versatility in how to represent the NPMs, e.g., graph kernels [86] can be used for NPMs represented as crystal graphs (e.g., [39]).

What is a GP? A GP is a stochastic process that treats the value of the objective function at any given point in feature space, $f(\mathbf{x})$, as a random variable. Specifically, GPs assume the joint distribution of any finite collection of function values, say at points $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$ on its domain, follows a multivariate Gaussian distribution

$$\mathbf{f} \equiv [f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_m)]^\top \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}) \quad (4)$$

whose covariance matrix $\mathbf{\Sigma} \in \mathbb{R}^{m \times m}$ is given by the kernel function applied pairwise over the points $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$, $\Sigma_{i,j} := k(\mathbf{x}_i, \mathbf{x}_j)$. The kernel function $k(\mathbf{x}, \mathbf{x}')$ quantifies the similarity of NPMs \mathbf{x} and \mathbf{x}' ; hence, the idea in GPs is that the properties $f(\mathbf{x})$ and $f(\mathbf{x}')$ of similar (dissimilar) NPMs \mathbf{x} and \mathbf{x}' are highly (un)correlated. [87] GPs effectively model the entire function $f(\mathbf{x})$ — in a point-wise manner— by assuming eqn. 4 holds for *any* arbitrary, finite collection of function values on its domain. The mean of zero in eqn. 4 assumes the measurements are centered.

From a Bayesian perspective, eqn. 4 is a prior assumption about the structure-property relationship $f(\mathbf{x})$. When we gather new observations, we will update this prior assumption to arrive at the posterior distribution reflecting our beliefs about the structure-property relationships in light of new data.

Kernel functions. Examples of kernel functions that operate on vector representations of two NPMs \mathbf{x} and \mathbf{x}' include the linear, polynomial, and radial basis function (RBF) kernels:

$$k(\mathbf{x}, \mathbf{x}') = \sigma_f \mathbf{x}^\top \mathbf{x}' \quad \text{linear kernel} \quad (5)$$

$$k(\mathbf{x}, \mathbf{x}') = \sigma_f (\mathbf{x}^\top \mathbf{x}')^d \quad \text{homogeneous polynomial kernel} \quad (6)$$

$$k(\mathbf{x}, \mathbf{x}') = \sigma_f e^{-\|\mathbf{x} - \mathbf{x}'\|_2^2 / (2\gamma^2)}. \quad \text{radial basis function (RBF) kernel} \quad (7)$$

Each kernel possesses the hyperparameter σ_f , the signal variance, which is a scale factor controlling the expected range of the functions represented by the GP. The polynomial kernel has a hyperparameter d that controls the order of the polynomial in the features, and the RBF kernel contains a length-scale hyperparameter γ that controls how close \mathbf{x} and \mathbf{x}' must be in the feature space to be considered “similar” and the expected roughness of the functions represented by the GP. Implicitly, each nonlinear kernel maps the two vectors \mathbf{x} and \mathbf{x}' into a new, higher-dimensional feature space through a mapping τ , then takes the inner product of the vectors in the new feature space:

$$k(\mathbf{x}, \mathbf{x}') = \tau(\mathbf{x})^\top \tau(\mathbf{x}'). \quad (8)$$

Interestingly, the feature map $\tau(\mathbf{x})$ corresponding to the RBF kernel in eqn. 7 maps vectors into an *infinite* dimensional feature space! By implicitly operating in a higher-dimensional feature space,

nonlinear kernels give GPs more flexibility, or expressiveness, for approximating complicated objective functions $f(\mathbf{x})$. Notably, graph kernels [86] and image kernels [88] can define similarities of two NPMs represented as graphs (nodes: atoms, edges: bonds) and images, respectively.

Inference with GPs. In BO, we exploit GPs for regression, with uncertainty quantification, to build a surrogate model for $f(\mathbf{x})$. We have observations $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ from previous experiments (previous iterations of BO), with y_i the measured property of NPM \mathbf{x}_i . Under the Bayesian view, y_i is a noise-free observation² of the random variable $f(\mathbf{x}_i)$. We wish to know the distribution of the random variable $f(\mathbf{x})$ for an *unevaluated* NPM \mathbf{x} , to determine the utility of evaluating it in the next experiment. Imposing the assumption in eqn. 4 for a specific collection of points on the domain of $f(\mathbf{x})$ composed of (i) the n evaluated NPMs from the past experiments $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ and (ii) the unevaluated NPM, \mathbf{x} :

$$\begin{bmatrix} \mathbf{f} \\ f(\mathbf{x}) \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{\Sigma} & \boldsymbol{\sigma} \\ \boldsymbol{\sigma}^\top & k(\mathbf{x}, \mathbf{x}) \end{bmatrix} \right) \quad (9)$$

with $\mathbf{f} = [f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n)]^\top$ the vector of random variables representing the properties of the previously evaluated NPMs, $\boldsymbol{\sigma} = [k(\mathbf{x}, \mathbf{x}_1), k(\mathbf{x}, \mathbf{x}_2), \dots, k(\mathbf{x}, \mathbf{x}_n)]^\top$ the vector of the kernel between the unevaluated NPM and the previously evaluated NPMs, and $\mathbf{\Sigma}$ the kernel matrix for the previously evaluated NPMs with $\Sigma_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j) \quad \forall i, j \in \{1, 2, \dots, n\}$. However, we have observations of the random variables in \mathbf{f} , $\mathbf{y} = [y_1, y_2, \dots, y_n]^\top$. Conditioning the joint distribution in eqn. 9 on the observations \mathbf{y} , we arrive at the posterior distribution for the property $f(\mathbf{x})$ of the *unevaluated* NPM, also Gaussian:

$$f(\mathbf{x})|\mathbf{y} \sim \mathcal{N}(\hat{y}(\mathbf{x}), \sigma^2(\mathbf{x})). \quad (10)$$

with mean and variance:

$$\hat{y}(\mathbf{x}) = \boldsymbol{\sigma}^\top \mathbf{\Sigma}^{-1} \mathbf{y} \quad (11)$$

$$\sigma^2(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}) - \boldsymbol{\sigma}^\top \mathbf{\Sigma}^{-1} \boldsymbol{\sigma}. \quad (12)$$

We can interpret eqns. 11 and 12. The predicted property of the unevaluated NPM, $\hat{y}(\mathbf{x})$, is a linear combination of the observed properties of the evaluated NPMs, \mathbf{y} with weights $\boldsymbol{\sigma}^\top \mathbf{\Sigma}^{-1}$. The weight on each measured property depends on the similarity between that NPM and the unevaluated NPM. The variance $\sigma^2(\mathbf{x})$ describing the uncertainty associated with the prediction of the property of the unevaluated NPM \mathbf{x} is the prior assumption of $k(\mathbf{x}, \mathbf{x})$ reduced by $\boldsymbol{\sigma}^\top \mathbf{\Sigma}^{-1} \boldsymbol{\sigma}$, which captures the similarity of the unevaluated NPM \mathbf{x} with the set of previously evaluated NPMs.

Fig. 3 illustrates a GP model of a toy function $f(x)$, based on an RBF kernel, over a toy one-dimensional NPM feature space $\mathcal{X} = \mathbb{R}$, trained on five observations. The mean in the GP, $\hat{y}(x)$, approximates $f(x)$, and the variance $\sigma^2(x)$ expresses uncertainty in the approximation. Generally, the uncertainty is small close to an observation and large when far from an observation.

²Note that we can pose GPs that relax the assumption that the observations are noise-free [83].

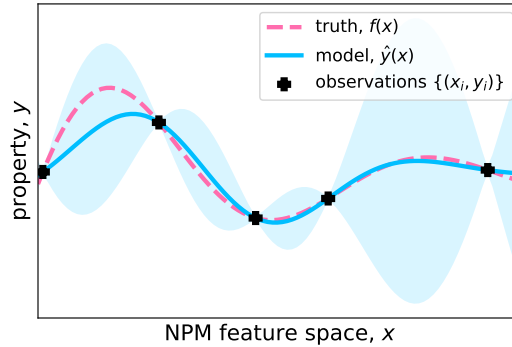


Figure 3: Illustration of a Gaussian process (GP) model with an RBF kernel over a toy one-dimensional NPM feature space. The black points are the observed data from a toy structure-property relationship, $f(x)$. The blue line and shaded region visualize the GP model trained on the observed data: the line is the approximation $\hat{y}(x)$ of the structure-property relationship, while the shaded region illustrates the uncertainty by covering $\hat{y}(x) \pm 2\sigma(x)$. The GP model shows large (small) uncertainty in regions of feature space far from (close to) the observations.

Hyperparameters of GPs. GPs are non-parametric models, but most useful kernel functions used in GPs contain hyperparameters. For example, the RBF kernel in eqn. 7 has the length-scale γ and the signal variance σ_f hyperparameters. To learn the hyperparameters of the kernel that give us the best approximation of $f(\mathbf{x})$, we could perform a grid search over hyperparameter space for those that minimize the validation loss. A more popular and faster way to estimate hyperparameters of the kernel is to maximize the marginal likelihood of the observed data as a function of the kernel hyperparameters. [87] Generally, at each iteration of BO, the hyperparameters of the kernel are updated to account for the newly acquired observation.

Two further interpretations of GP models of functions. A GP model of the objective function $f(\mathbf{x})$ can be interpreted as (i, weight space view) Bayesian linear regression in the implicit feature space of the kernel and (ii, function space view) a distribution over functions [83]. To clarify, the GP model of $f(\mathbf{x})$ in eqn. 4 is equivalent to the parametric model:

$$f(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\tau}(\mathbf{x}), \quad (13)$$

with weights \mathbf{w} on which we place a Gaussian prior and $\boldsymbol{\tau}$ the map associated with the kernel $k(\mathbf{x}, \mathbf{x}')$ used in the GP. In the weight space view, GP inference models the posterior distribution over weights \mathbf{w} in eqn. 13. In the function space view, GP inference models the posterior distribution over the space of functions represented in eqn. 13. Though eqn. 13 is helpful for understanding GPs and sampling functions from the distribution over the function space they describe, we in practice conduct GP inference using the kernel, through eqns. 11 and 12. E.g. $\boldsymbol{\tau}(\mathbf{x})$ is a vector of infinite dimension in the case of the RBF kernel, making the view of GPs in eqn. 13 unfriendly for computations.

5.2 Examples of acquisition functions

We provide three common examples of acquisition functions and explain how they use the surrogate model to trade-off exploration and exploitation to select the NPM to evaluate in the next experiment.

Upper confidence bound (UCB). The UCB acquisition function selects the point that maximizes the upper confidence function:

$$A(\mathbf{x}) := \hat{y}(\mathbf{x}) + \beta\sigma(\mathbf{x}) \quad (14)$$

where $\hat{y}(\mathbf{x})$ and $\sigma(\mathbf{x})$ are the predicted property of NPM \mathbf{x} and its associated uncertainty, respectively, provided by the surrogate model. The parameter β explicitly trades-off exploration and exploitation. If β is large, the UCB tends to select NPMs with the highest uncertainty to explore the feature space. If β is small, the UCB tends to select NPMs with the highest predicted property to exploit our current approximation of the structure-property relationship $f(\mathbf{x})$. To clarify the UCB terminology, the top boundary of the shaded region that bands $\hat{y}(x)$ in Fig. 3 is the UCB for $\beta = 2$.

Expected improvement (EI). Another acquisition strategy is to select the NPM with the highest expected improvement (EI) of the property of the best evaluated NPM so far. Let the random variable $I(\mathbf{x}) = \max(0, f(\mathbf{x}) - \max_i y_i)$ denote the improvement in the property of a NPM \mathbf{x} over the best observed NPM thus far. The EI is then:

$$A(\mathbf{x}) := \int_{-\infty}^{\infty} I(\mathbf{x}) \mathcal{N}(y|\hat{y}(\mathbf{x}), \sigma^2(\mathbf{x})) dy, \quad (15)$$

which can be written in closed form:

$$A(\mathbf{x}) = \begin{cases} (\hat{y}(\mathbf{x}) - \max_i y_i) \Phi\left(\frac{\hat{y}(\mathbf{x}) - \max_i y_i}{\sigma(\mathbf{x})}\right) + \sigma(\mathbf{x}) \phi\left(\frac{\hat{y}(\mathbf{x}) - \max_i y_i}{\sigma(\mathbf{x})}\right) & \sigma(\mathbf{x}) > 0 \\ 0 & \sigma(\mathbf{x}) = 0, \end{cases} \quad (16)$$

with Φ and ϕ the cumulative and probability distribution functions, respectively, of the standard normal distribution. The first and second terms in eqn. 16, respectively, capture the exploitation and exploration component of EI.

Information-theoretic acquisition functions. The principle behind acquisition functions based on information theory is to select the NPM \mathbf{x} that maximizes the mutual information between (i) the observation of its property $y=f(\mathbf{x})$ and (ii) the location of the NPM \mathbf{x}^* in feature space that maximizes $f(\mathbf{x})$. Viewing both $f(\mathbf{x})$ and \mathbf{x}^* as random variables, the following acquisition function describes the mutual information between the observation ($\mathbf{x}, y = f(\mathbf{x})$) of the property of a newly acquired NPM, \mathbf{x} , and the location of the optimal NPM, \mathbf{x}^* .

$$A(\mathbf{x}) = MI[(\mathbf{x}, y), \mathbf{x}^*] \quad (17)$$

$$= H[p(\mathbf{x}^*)] - H[p(\mathbf{x}^*|(\mathbf{x}, y))] \quad (18)$$

where $MI[\cdot, \cdot]$ is the mutual information between two random variables and $H(\cdot)$ is the entropy of a probability distribution $p(\cdot)$. The mutual information is the reduction in the entropy of the probability density function of the location of the optimum NPM, \mathbf{x}^* , as a result of observing the property

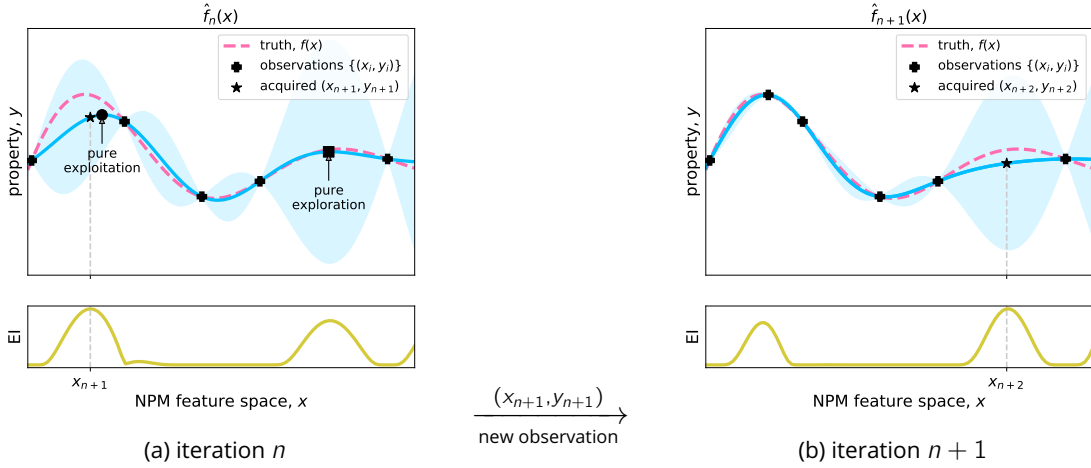


Figure 4: Illustration of one iteration of BO where (i) the acquisition function is used to select the next NPM whilst balancing exploitation and exploration and (ii) the surrogate GP model based on the RBF kernel is updated to account for the newly acquired observation. (a) Iteration n . (b) Iteration $n + 1$, after the surrogate model is updated by the new observation (x_{n+1}, y_{n+1}) acquired to maximize the expected improvement (EI). In both (a) and (b), the top panel shows the surrogate model, and the bottom panel shows the expected improvement (EI) acquisition function. For comparison, in (a) we illustrate the NPM x_{n+1} that would have been selected under a pure exploitation or pure exploration acquisition strategy.

$y = f(\mathbf{x})$ of NPM \mathbf{x} and thus conditioned on knowing (\mathbf{x}, y) . Conceptually, the distribution $p(\mathbf{x}^*)$ could be obtained under a GP surrogate model by sampling functions from eqn. 13 and then optimizing them. In practice, eqn. 18 is difficult to compute, but there are several acquisition functions based on instantiations of this idea [89].

Illustrating BO acquisition and the exploration-exploitation tradeoff. Fig. 4a illustrates the EI acquisition function in a toy one-dimensional NPM space under a GP surrogate model with $n = 5$ observations. The EI acquisition function exhibits two maxima. The first (global) maximum is in a region of the feature space where the predicted property $\hat{y}(x)$ is the largest. The second (local) maximum is where the uncertainty $\sigma(x)$ is largest. We select the NPM for the next experiment, x_{n+1} , as the one that maximizes the EI acquisition function. Fig. 4a shows the acquired NPM assuming the database of NPMs \mathcal{X} covers all points on the domain shown. To illustrate how the EI acquisition strategy balances exploration and exploitation, for comparison, we also show the acquired NPM x_{n+1} if the acquisition strategy were purely exploration and purely exploitation. Pure exploration dictates $x_{n+1} = \arg \max \sigma(x)$, but this NPM has a poor property. Pure exploitation dictates $x_{n+1} = \arg \max \hat{y}(x)$, but this NPM is too close to an existing observation. EI balances the trade-off by picking an NPM with both a high uncertainty and high predicted property.

6 Experiments and Results

We now demonstrate Bayesian optimization by using it to efficiently search for covalent organic frameworks (COFs) for vehicular natural gas storage [2]. Our experiments below use the open data from Mercado et al. [28] and can be fully reproduced on a desktop computer using our computer code at github.com/aryandeshwal/BO_of_COFs.

6.1 Experimental problem setup

Our goal is to efficiently search a database of COFs for the one with the largest methane deliverable capacity.

The database of COFs, \mathcal{X} . The database of COFs contains 69 840 2D and 3D predicted COF structures constructed by Mercado et al. [28].

The COF vector representation, \mathbf{x} . We represent each COF structure with a vector $\mathbf{x} \in \mathbb{R}^{12}$ of structural and chemical features listed in Tab. 1 and computed by Mercado et al. [28]. We Min-Max normalized each feature³ to lie in $[0, 1]$. This defines COF feature space as $[0, 1]^{12}$.

The methane deliverable capacity, y . The COF property we wish to maximize is the simulated deliverable capacity of methane [L STP CH₄/L COF] at 298 K under a 65 bar to 5.8 bar pressure swing. The deliverable capacity of the COF primarily determines the driving range of a vehicle on a “full” adsorbed natural gas fuel tank packed with the COF [2].

The expensive objective function, $f(\mathbf{x})$. Evaluating the objective function f to give $y = f(\mathbf{x})$ involves conducting two grand-canonical Monte Carlo simulations of methane adsorption in the COF structure represented by \mathbf{x} —one at (65 bar, 298 K) and one at (5.8 bar, 298 K). The deliverable capacity y then follows from the difference in the simulated methane adsorption at the two conditions. The function f is expensive to evaluate, as the run time of the molecular simulations is on the order of hours.

Goal: data-efficient search for the optimal COF, \mathbf{x}^* . In an *exhaustive search*, we would conduct expensive molecular simulations to predict the methane deliverable capacity of each candidate COF in the database—i.e., collect $\{(\mathbf{x}, y = f(\mathbf{x})) : \mathbf{x} \in \mathcal{X}\}$ —to find the COF $\mathbf{x}^* \in \mathcal{X}$ with the highest deliverable capacity. In contrast to an exhaustive search, instead, our goal is to find the optimal COF \mathbf{x}^* *efficiently*—while conducting expensive molecular simulations *in only a small proportion of the candidate COFs*.

We hypothesize that BO will provide a simulation-efficient search for the optimal COF, \mathbf{x}^* . In reality, Mercado et. al. [28] already simulated methane adsorption in all of the COFs at (65 bar, 298

³We used the feature vectors of all COFs for the Min-Max normalization (both acquired and non-acquired COFs). This does not constitute data leakage because, in our setting, (i) the features are cheap to compute and (ii) we have a finite library of COFs for which it is feasible to compute all features for all COFs in \mathbf{X} .

K) and (5.8 bar, 298 K) and computed their methane deliverable capacities. Thus, (i) we know the optimal COF \mathbf{x}^* and (ii) as opposed to actually conducting molecular simulations of methane adsorption in a selected COF during the active search, we instead look up the result of the simulations (the deliverable capacity) from the data of Mercado et al. [28]. Each data lookup, conceptually, represents conducting the two expensive molecular simulations of methane adsorption in a COF ourselves. N.b., that we look up data as opposed to conducting a simulation ourselves has no impact on the BO search efficiency when defined in terms of the number of COF evaluations needed to find the optimal COF. The exhaustive search of Mercado et al. [28] allows us to readily evaluate the simulation-efficiency of different search strategies to find the optimal COF(s). We will compare the search efficiency of BO to random search, an evolutionary algorithm, and one-shot supervised learning.

Table 1: Features comprising the vector representation of a COF, \mathbf{x} , broken into those that capture its structure and chemical composition.

| structural (geometrical) | chemical composition |
|----------------------------------|----------------------|
| void fraction | density of carbon |
| density | density of flourine |
| largest included sphere diameter | density of hydrogen |
| largest free sphere diameter | density of nitrogen |
| gravimetric surface area | density of oxygen |
| | density of sulfur |
| | density of silicon |

6.2 Search strategies

We use several different strategies to search for the optimal COF \mathbf{x}^* exhibiting the highest methane deliverable capacity y in the database \mathcal{X} .

Random search. Random search is a naive baseline. At each iteration, we uniform randomly select an unevaluated COF from \mathcal{X} to evaluate.

Random search does not make a *data-informed* selection of a COF for the next evaluation, as it ignores the past observations, $\{(\mathbf{x}_i, y_i = f(\mathbf{x}_i))\}$. Thus, random search is expected to perform poorly.

Bayesian optimization (BO). For BO, we employ (1) a Gaussian process (GP) with the Matérn kernel ($\nu = 2.5$) [83] as our surrogate model and (2) the expected improvement (EI) in eqn. 16 as our acquisition function. To initialize the GP for BO, we first randomly select ten COFs from the database and evaluate their methane deliverable capacity. The GP surrogate model is then trained on $\{(\mathbf{x}_i, y_i)\}_{i=1}^{10}$, which count towards the number of evaluations when we report the search-efficiency of BO. At each iteration of BO, we fit a new GP to all past observations $\{(\mathbf{x}_i, y_i)\}$, which includes choosing the hyperparameters of the Matérn kernel (length-scale and signal variance) by maximizing the marginal likelihood of the data under the GP. We implemented our BO procedure in the BoTorch library [90].

In accordance with the assumption behind GPs, we normalize the deliverable capacities to have mean zero and unit variance (using the training data only). During the acquisition phase, we evaluate the acquisition function for each COF in the database and select the COF with the highest value, in contrast to optimizing the acquisition function over the continuous COF feature space.

Evolutionary search (via CMA-ES). As an evolutionary search method, we use Covariance Matrix Adaptation Evolution Strategy (CMA-ES) [91, 92], a state-of-the-art, stochastic optimizer for rugged, non-convex, black-box objective functions. In CMA-ES, new COFs are stochastically selected from the search space by sampling from a multivariate Gaussian distribution over the feature space. The mean and covariance matrix of the distribution are updated over the search process, as COFs are acquired and evaluated in batches (generations). To update the distribution, with the aim of increasing the likelihood of acquiring and taking search steps towards high-performing COFs, (i) the mean is updated using a weighted average of the most high-performing COFs in the generation (a selection mechanism) and (ii) the covariance matrix is updated using a weighted average of the best search steps (from the mean) towards the high-performing COFs [92].

CMA-ES has two hyperparameters: the initial standard deviation for each feature of the COF representation and the number of new candidate COFs acquired in each iteration (the population size). We initialized the CMA-ES algorithm with a randomly selected COF and set the initial standard deviation to 0.5 to cover our COF feature space $[0, 1]^{12}$. The population size, 11, was determined by a default heuristic in the `cma` library in Python.

CMA-ES operates in a continuous search space. When it selects a point in COF feature space for the next generation, it does not exactly correspond to a feature vector of a COF in the database; to apply CMA-ES, we select the nearest (in feature space), unacquired COF in the database.

One-shot supervised learning (via RF). While one-shot supervised learning does not constitute active search like BO, it is the most popular method for circumventing an exhaustive search for the optimal NPM using the high-fidelity evaluation method. One-shot supervised learning has three stages. (1) COFs are selected and evaluated to serve as training data for the supervised learning model. (2) The trained model is used to predict the deliverable capacity of the remaining COFs. (3) We evaluate the COFs with the highest predicted deliverable capacity. Stages (1) and (3) both incur (costly) COF evaluations. We compare one-shot supervised learning and active search by comparing the deliverable capacities of their acquired set of COFs when given the same budget of COF evaluations. Much like balancing exploration and exploitation, we elect to split the budget of evaluations for one-shot supervised learning among stages (1) and (3) equally. More, for stage (1), we try two training set acquisition strategies: (i) uniform random selection of COFs and (ii) a max-min diversity selection strategy [40, 93], whereby we sequentially acquire COFs for the training set, selecting the COF with the maximum minimum distance (in COF feature space) to a COF already in the training set (starting with an initial, random COF).

As the supervised learning model, we use commonly-used [40–42, 46] random forests (RFs) (100 trees, default parameters in `scikit-learn`) as the supervised learning model to approximate $f(\mathbf{x})$ using the (differently sized) training sets.

6.3 Results

We now execute each strategy to search for the optimal COF \mathbf{x}^* exhibiting the highest methane deliverable capacity y in the database \mathcal{X} .

6.3.1 Search efficiency

The blue curves in Fig. 5 show the search efficiency of BO. Three different search performance metrics are shown as the number of COFs evaluated, n , (= the number of BO iterations = the number of simulations/"experiments") increases. The first metric in (a) is the maximum deliverable capacity among the acquired set of COFs, \mathcal{X}_n . The second and third metric are, among the acquired set of COFs \mathcal{X}_n , (b) the highest deliverable capacity rank, with the rank defined using the entire data set \mathbf{X} , and (c) the fraction of the 100 COFs in \mathcal{X} with the highest deliverable capacity. 95% of the BO searches acquired the optimal COF \mathbf{x}^* after $n = 120$ COF evaluations; all 100 BO searches acquired the optimal COF after $n = 174$ COF evaluations. After $n = 250$ COF evaluations, BO acquires 36% of the top 100 COFs in the data set. This illustrates how BO can provide a simulation-efficient search for the optimal COF as opposed to conducting an exhaustive search; BO picked out the top COF in the database of $\sim 70\,000$ COFs while evaluating less than 200!

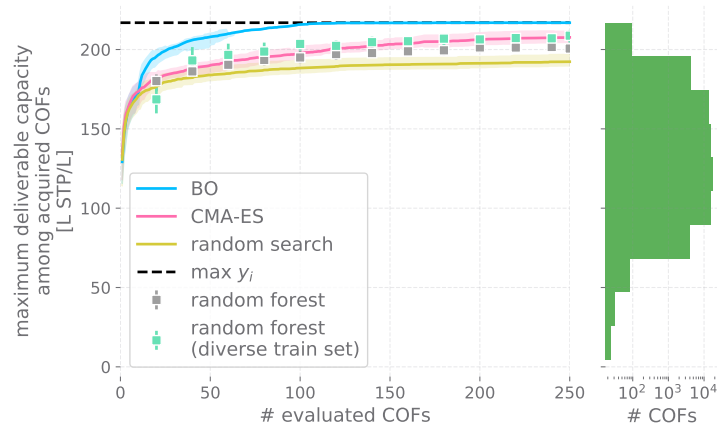
BO also outperforms random search, evolutionary search, and the one-shot supervised learning approach. With the same number of COF evaluations ($n = 250$), a random search (on average) acquires only the 308th ranked COF and 0.25% of the top 100 COFs. The performance of random search is poor because it ignores the past observations when it selects the next COF to evaluate. Evolutionary search and one-shot supervised learning (diverse training set) have a much better search performance than random search, acquiring the 13th/12th ranked COF and the top 11%/15% of the top 100 COFs on a budget of 250 evaluations. Though, evolutionary search nor the one-shot supervised learning strategies recover the optimal COF \mathbf{x}^* after 250 evaluations. Thus, BO outperforms the baseline search methods of evolutionary search and one-shot supervised learning using both metrics of (a) the highest deliverable capacity in the acquired set and (b) the fraction of the top 100 COFs in the acquired set given a budget of 250 evaluations. N.b., BO is designed to optimize the former performance metric, but BO could be tailored to optimize a top- k metric [94, 95].

Regarding the random versus diverse training set acquisition strategies for the one-shot supervised learning approach: Except when the training set is very small (10), the diverse selection of training data gives better search performance than the random selection of training data since it provides better coverage of the feature space.

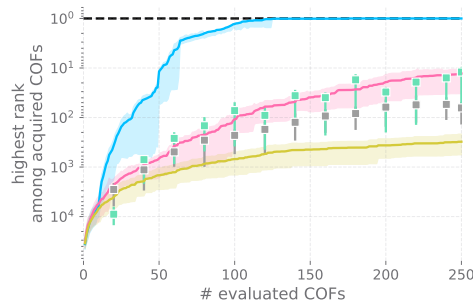
6.3.2 Visualizing the BO acquisition set

To understand the behavior of BO for searching the database of COFs for the optimal COF \mathbf{x}^* , we visualize the acquisition set of COFs in feature space as BO progresses. Given that the feature space is 12-dimensional, we resort to principal component analysis (PCA) to [approximately] reduce the dimension of the feature space to two. I.e., we project each COF feature space onto a 2D reduced feature space through PCA of the data $[\mathbf{x}_1 \mathbf{x}_2 \cdots \mathbf{x}_{|\mathcal{X}|}]$

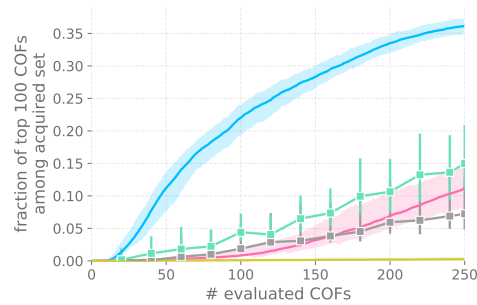
First, we visualize the structure-property relationship $f(\mathbf{x})$. Fig. 6a shows a depiction of $f(\mathbf{x})$ as a



(a)



(b)



(c)

Figure 5: The search efficiency of BO in comparison with random search, evolutionary search, and one-shot supervised learning. The curves show search efficiency in terms of (a) the maximum deliverable capacity, (b) the highest ranking (among the entire data set \mathcal{X}) of the COF deliverable capacity, and (c) the fraction of the top 100 ranked (among \mathcal{X}) COFs- among the acquired COFs, in terms of the number of acquired/evaluated COFs. The shaded region shows the variance over 100 [stochastic] runs. To give (a) context, we show the distribution of the deliverable capacities among the COFs in the entire data set, \mathcal{X} , on the right.

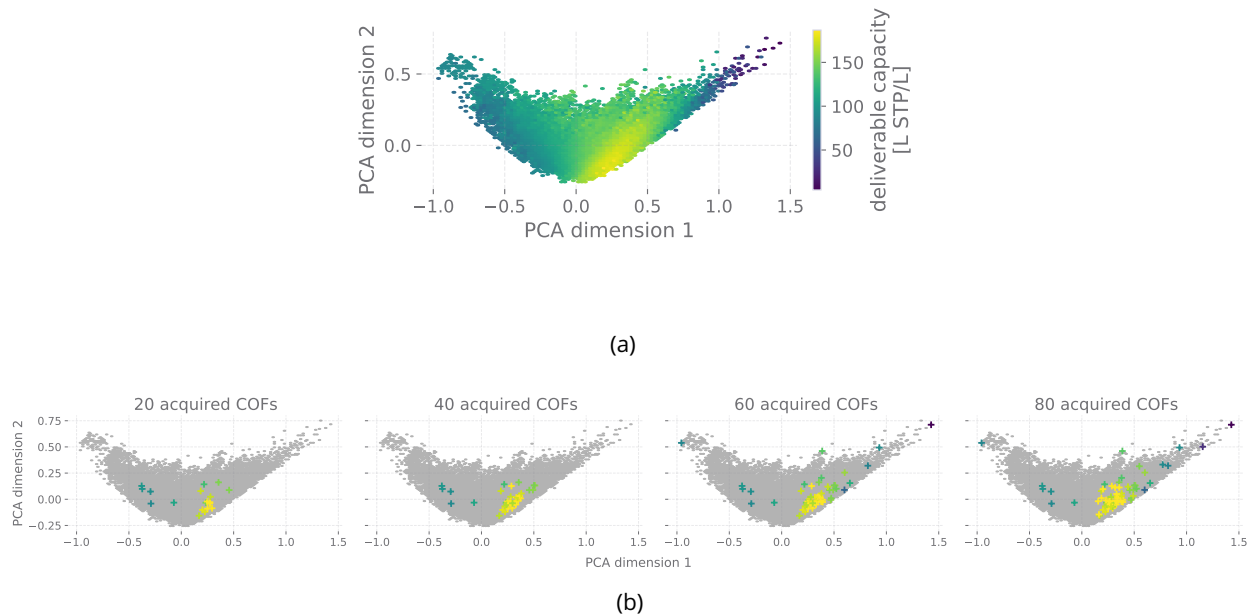


Figure 6: Illustration of the COF acquisition set under BO. Each panel is a visualization of the first two principal components of COF feature space. (a) A heat map where the color of each hexagonal voxel indicates the average deliverable capacity of a COF that falls in that voxel. This is an attempt to visualize $f(\mathbf{x})$. (b) The points represent the acquired COFs at different stages of BO– after 20, 40, 60, and 80 iterations– and are colored according to the deliverable capacity (colorbar in (a) pertains). The gray background shows the region of the feature space covered by the COFs in (a).

2D heatmap over the reduced 2D feature space. The color of each voxel in the reduced COF space indicates the average deliverable capacity of COFs that fall in that voxel of COF space.

Fig. 6b shows the acquired set of COFs, colored by deliverable capacity, at 20, 40, 60, and 80 iterations of BO. For reference, the gray background shows the coverage of COF space by all COFs in the dataset, shown in Fig. 6a. Reflecting the dominance of exploitation, BO concentrates its acquires on the region containing the COFs with highest deliverable capacities and avoids acquiring COFs from regions containing COFs with low deliverable capacities. After 40 iterations, we see the exploratory component of BO acquisition, where it acquires COFs in low-performing regions of the feature space.

6.3.3 Balancing exploration and exploitation

We conceptually illustrated how the expected improvement acquisition function balances exploration and exploitation in Fig. 4. We now show how balancing exploration and exploitation is critical for BO to recover the optimal COF with the fewest experiments. To do so, we compare the search

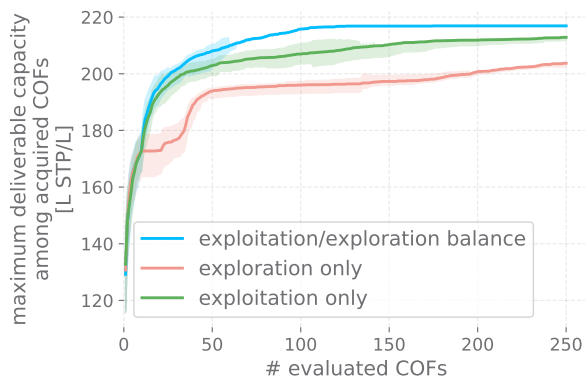


Figure 7: BO search efficiency using three different acquisition functions: expected improvement (balances exploration and exploitation), the predicted deliverable capacity (full exploitation), and the uncertainty in the predicted deliverable capacity (full exploration).

efficiency of BO using three different acquisition functions:

$$A(\mathbf{x}) = EI(\mathbf{x}) \quad \text{exploration-exploitation balance} \quad (19)$$

$$A(\mathbf{x}) = \hat{y}(\mathbf{x}) \quad \text{pure exploitation} \quad (20)$$

$$A(\mathbf{x}) = \sigma(\mathbf{x}) \quad \text{pure exploration} \quad (21)$$

The expected improvement (EI) acquisition function in eqn. 16, used in Fig. 5, balances exploration and exploitation; acquiring the COF with the highest predicted deliverable capacity \hat{y} is pure exploitation; acquiring the COF with the highest uncertainty σ in the predicted deliverable capacity is pure exploration (active learning). Fig. 7 shows the search efficiency of BO under these three different acquisition strategies, in terms of the highest deliverable capacity among the acquired set of COFs. Both the pure exploitation and pure exploration BO acquisition strategies exhibit subpar search performance compared to the the EI acquisition strategy that trades off exploration and exploitation. The pure exploration acquisition strategy (active learning) performs the worst, as it acquires COFs with low deliverable capacities to reduce the uncertainty in the surrogate model's predictions about COFs with low deliverable capacities.

6.4 Conclusions from experiments

BO is an active search method to find the optimal NPM in a data set whilst evaluating, with some expensive method such as a molecular simulation, only a small fraction of the NPMs in the data set. BO achieves this by making acquisition decisions that take into account all past observations and balance exploration and exploitation. The adoption of BO could dramatically impact high-throughput computational screenings of NPMs by reducing the computing cost of finding the optimal NPM, allowing us to screen larger databases of NPMs, and enabling the use of higher-fidelity but more expensive molecular models and simulation methods. Notably, BO applies to NPM search in the experimental domain as well.

7 Outlook

We explained the key ideas behind Bayesian optimization (BO) and advocated for its use to efficiently search databases of NPMs for the one with the optimal property, whilst synthesizing and evaluating the fewest NPMs. The ideas of BO, to sequentially, actively make intelligent decisions on which NPM to synthesize and evaluate based on the past experiments, can be applied to both the laboratory (driven by humans or robots [96–99]) and computational settings. The two core ingredients of BO are (1) a surrogate model that approximates the structure-property relationship and describes our uncertainty in it and (2) an acquisition function that scores the utility of evaluating each NPM next, designed to balance exploration and exploitation. We demonstrated BO of NPMs by using BO to search through a database of ca. 70 000 COFs to find the COF with the highest methane deliverable capacity; all 100 BO searches acquired the optimal COF after evaluating only 174 COFs. While preparing our article, Donval and Hand *et al.* [100] also demonstrated BO of MOFs and COFs for the acceleration of virtual screenings.

There are several extensions to and modifications of Bayesian optimization that are useful for different problem settings in NPM discovery:

- **Batch BO.** In standard BO, we select a single NPM to evaluate in each iteration. However, we may have parallel experimental resources to leverage to further accelerate the search for the optimal NPM. In *batch BO* [95,101–105], we can select multiple NPMs to synthesize and evaluate in parallel during each BO iteration.
- **Multi-fidelity BO.** Often, we have a choice of multiple experimental methods to evaluate the property of an NPM. These methods usually involve a tradeoff in resource cost and the accuracy of the evaluation. For example, a molecular simulation of gas adsorption in an NPM is a low-fidelity experiment (cheap, but inaccurate) while measurement of gas adsorption in an NPM in a physical laboratory is a high-fidelity experiment (costly, but accurate). Intuitively, it is possible to leverage low-fidelity experiments to prune NPMs with low property values and to identify promising NPM candidates that can be searched further using high-fidelity experiments. In *multi-fidelity BO* [106–115], we select both an NPM to evaluate and the fidelity of the experiment in each iteration. This allows optimization of the overall resource cost of experiments—of both low- and high-fidelity—for identifying high-performing NPMs.
- **Multi-objective BO.** We often need to optimize NPMs for multiple property objectives which are conflicting in nature and cannot be optimized simultaneously. For example, for gas separations, we often wish an NPM to have both a high selectivity and a high working capacity for the gas we wish to capture [53]. For multi-objective optimization problems, we need to find the Pareto optimal set of solutions. A solution Pareto optimal if it cannot be improved in any of the objectives without compromising some other objective. The goal of *multi-objective BO* is to find the optimal Pareto set of NPMs using fewest NPM evaluations [116–122]. Similarly, the ϵ -PAL algorithm [123] has recently been used to find the Pareto optimum polymers.
- **Constrained BO.** Possibly, some NPMs in the search space cannot be synthesized. More, often we cannot know if an NPM is synthesizable until we attempt its synthesis, which still incurs a cost. In this context, synthesizability is a black-box constraint over the search space.

In *constrained BO* [124–128], we perform BO where the synthesizability of an NPM cannot be verified without performing an experiment. The typical approach involves learning a statistical model based on the past evaluations of constraint(s) and selecting high-utility NPMs from the predicted feasible region (minimal to no constraint violation).

- **Cost-aware BO.** The evaluation cost can vary from one NPM to another (e.g., cost of synthesizing NPM). We would like to take this cost into account to reduce the overall costs incurred during the search for the optimal NPM. In *cost-aware BO* [129, 130], the acquisition strategy considers not only the information gain of acquiring an NPM but also the cost incurred to synthesize it and measure its property.
- **Robust solutions to BO.** We may be uncertain about the measured/computed features of the NPMs and seek an optimum NPM that is robust to variations in its features. In *robust BO*, we account for the uncertainty in the inputs \mathbf{x} when optimizing $f(\mathbf{x})$ and seek flat as opposed to sharp optima [131, 132].

Some popular software packages for BO include BoTorch [133], BayesOpt [134], and SMAC [135]. COMBO [136] is a BO library tailored to materials science.

In addition to efficiently searching for NPMs with optimal properties, BO is applicable to a wide variety of optimization problems in the chemical and materials sciences [137, 138]. BO has been used to efficiently search for optimal reaction conditions [139, 140], compositions of and processing conditions for materials [98, 141], ligands to dock on proteins [95], crystal structures [142, 143], shape memory alloys [144], and density functional models [145]. For more general overview, see the reviews of Coley [146], Tsuda and coworkers [147], Frazier and Wang [148], and Lookman *et al.* [149].

The effectiveness of BO is predicated upon an accurate surrogate model of the structure-property relationship. In turn, the accuracy of the surrogate model is predicated on (i) an information-rich representation \mathbf{x} of the NPM that encodes the salient features of its structure and chemical composition and (ii) a statistical model that (a) is sufficiently flexible/expressive to approximate the underlying objective function and (b) learns in a data-efficient manner. This gives important and currently active directions for future research. Particularly, engineering useful vector representations \mathbf{x} of NPMs, using domain knowledge, is a very active research area [150, 151]. The representation should be invariant to rotations, translations, replications (if a crystal), and permutations of the list of atoms comprising the structure. Moreover, the mapping from NPM structures to feature vectors should be injective. Graph neural networks can, instead, *learn* vector representations of NPMs from their crystal structures represented as graphs with node labels.

8 Acknowledgements

A.D. and J.D. acknowledge support from NSF grants IIS-1845922, OAC-1910213. C.M.S. acknowledges support from NSF grant 1920945.

References

- [1] Shengqian Ma and Hong-Cai Zhou. Gas storage in porous metal–organic frameworks for clean energy applications. *Chemical Communications*, 46(1):44–53, 2010.
- [2] Alexander Schoedel, Zhe Ji, and Omar M. Yaghi. The role of metal–organic frameworks in a carbon-neutral energy cycle. *Nature Energy*, 1(4), 2016.
- [3] Bin Li, Hui-Min Wen, Wei Zhou, and Banglin Chen. Porous metal–organic frameworks for gas storage and separation: What, how, and why? *The Journal of Physical Chemistry Letters*, 5(20):3468–3479, 2014.
- [4] Lauren E. Kreno, Kirsty Leong, Omar K. Farha, Mark Allendorf, Richard P. Van Duyne, and Joseph T. Hupp. Metal–organic framework materials as chemical sensors. *Chemical Reviews*, 112(2):1105–1125, 2011.
- [5] Leslie J. Murray, Mircea Dincă, and Jeffrey R. Long. Hydrogen storage in metal–organic frameworks. *Chemical Society Reviews*, 38(5):1294, 2009.
- [6] Kenji Sumida, David L Rogow, Jarad A Mason, Thomas M McDonald, Eric D Bloch, Zoey R Herm, Tae-Hyun Bae, and Jeffrey R Long. Carbon dioxide capture in metal–organic frameworks. *Chemical Reviews*, 112(2):724–781, 2012.
- [7] Fei-Yan Yi, Dongxiao Chen, Meng-Ke Wu, Lei Han, and Hai-Long Jiang. Chemical sensors based on metal–organic frameworks. *ChemPlusChem*, 81(8):675–690, 2016.
- [8] Zhichao Hu, Benjamin J Deibert, and Jing Li. Luminescent metal–organic frameworks for chemical sensing and explosive detection. *Chemical Society Reviews*, 43(16):5815–5840, 2014.
- [9] Christian S. Diercks and Omar M. Yaghi. The atom, the molecule, and the covalent organic framework. *Science*, 355(6328):eaal1585, 2017.
- [10] Xiao Feng, Xuesong Ding, and Donglin Jiang. Covalent organic frameworks. *Chemical Society Reviews*, 41(18):6010, 2012.
- [11] Maria S. Lohse and Thomas Bein. Covalent organic frameworks: Structures, synthesis, and applications. *Advanced Functional Materials*, 28(33):1705553, 2018.
- [12] Shun Wan, Felipe Gándara, Atsushi Asano, Hiroyasu Furukawa, Akinori Saeki, Sanjeev K. Dey, Lei Liao, Michael W. Ambrogio, Youssry Y. Botros, Xiangfeng Duan, Shu Seki, J. Fraser Stoddart, and Omar M. Yaghi. Covalent organic frameworks with high charge carrier mobility. *Chemistry of Materials*, 23(18):4094–4097, 2011.
- [13] Daniele Ongari, Aliaksandr V. Yakutovich, Leopold Talirz, and Berend Smit. Building a consistent and reproducible database for adsorption evaluation in covalent–organic frameworks. *ACS Central Science*, 5(10):1663–1675, 2019.
- [14] Hiroyasu Furukawa, Kyle E. Cordova, Michael O’Keeffe, and Omar M. Yaghi. The chemistry and applications of metal-organic frameworks. *Science*, 341(6149):1230444, 2013.

- [15] David J. Tranchemontagne, Zheng Ni, Michael O'Keeffe, and Omar M. Yaghi. Reticular chemistry of metal–organic polyhedra. *Angewandte Chemie International Edition*, 47(28):5136–5147, 2008.
- [16] Tom Hasell and Andrew I Cooper. Porous organic cages: soluble, modular and molecular pores. *Nature Reviews Materials*, 1(9):1–14, 2016.
- [17] Mian Li, Dan Li, Michael O'Keeffe, and Omar M. Yaghi. Topological analysis of metal–organic frameworks with polytopic linkers and/or multiple building units and the minimal transitivity principle. *Chemical Reviews*, 114(2):1343–1370, 2013.
- [18] Valentina Santolini, Marcin Miklitz, Enrico Berardo, and Kim E. Jelfs. Topological landscapes of porous organic cages. *Nanoscale*, 9(16):5280–5298, 2017.
- [19] José L Segura, Sergio Royuela, and M Mar Ramos. Post-synthetic modification of covalent organic frameworks. *Chemical Society Reviews*, 48(14):3903–3945, 2019.
- [20] Sukhendu Mandal, Srinivasan Natarajan, Prabu Mani, and Asha Pankajakshan. Post-synthetic modification of metal–organic frameworks toward applications. *Advanced Functional Materials*, 31(4):2006291, 2020.
- [21] Peter G. Boyd, Yongjin Lee, and Berend Smit. Computational development of the nanoporous materials genome. *Nature Reviews Materials*, 2(8), 2017.
- [22] Yamil J. Colón, Diego A. Gómez-Gualdrón, and Randall Q. Snurr. Topologically guided, automated construction of metal–organic frameworks and their evaluation for energy-related applications. *Crystal Growth & Design*, 17(11):5801–5810, 2017.
- [23] Lukas Turcani, Enrico Berardo, and Kim E. Jelfs. stk : A python toolkit for supramolecular assembly. *Journal of Computational Chemistry*, 39(23):1931–1942, 2018.
- [24] Peter G Boyd and Tom K Woo. A generalized method for constructing hypothetical nanoporous materials of any net topology from graph theory. *CrystEngComm*, 18(21):3777–3792, 2016.
- [25] Richard Luis Martin and Maciej Haranczyk. Construction and characterization of structure models of crystalline porous polymers. *Crystal Growth & Design*, 14(5):2431–2440, 2014.
- [26] Arni Sturluson, Melanie T. Huynh, Alec R. Kaija, Caleb Laird, Sunghyun Yoon, Feier Hou, Zhenxing Feng, Christopher E. Wilmer, Yamil J. Colón, Yongchul G. Chung, Daniel W. Siderius, and Cory M. Simon. The role of molecular modelling and simulation in the discovery and deployment of metal-organic frameworks for gas storage and separation. *Molecular Simulation*, 45(14-15):1082–1121, 2019.
- [27] Hilal Daglar and Seda Keskin. Recent advances, opportunities, and challenges in high-throughput computational screening of MOFs for gas separations. *Coordination Chemistry Reviews*, 422:213470, 2020.
- [28] Rocío Mercado, Rueih-Sheng Fu, Aliaksandr V. Yakutovich, Leopold Talirz, Maciej Haranczyk, and Berend Smit. In silico design of 2d and 3d covalent organic frameworks for methane storage applications. *Chemistry of Materials*, 30(15):5069–5086, 2018.

- [29] Kevin Maik Jablonka, Daniele Ongari, Seyed Mohamad Moosavi, and Berend Smit. Big-data science in porous materials: Materials genomics and machine learning. *Chemical Reviews*, 120(16):8066–8129, 2020.
- [30] Robert Pollice, Gabriel dos Passos Gomes, Matteo Aldeghi, Riley J. Hickman, Mario Krenn, Cyrille Lavigne, Michael Lindner-D’Addario, AkshatKumar Nigam, Cher Tian Ser, Zhenpeng Yao, and Alán Aspuru-Guzik. Data-driven strategies for accelerated materials design. *Accounts of Chemical Research*, 54(4):849–860, 2021.
- [31] Sanggyu Chong, Sangwon Lee, Baekjun Kim, and Jihan Kim. Applications of machine learning in metal-organic frameworks. *Coordination Chemistry Reviews*, 423:213–487, 2020.
- [32] Siwar Chibani and François-Xavier Coudert. Machine learning approaches for the prediction of materials properties. *APL Materials*, 8(8):080701, 2020.
- [33] Seyed Mohamad Moosavi, Kevin Maik Jablonka, and Berend Smit. The role of machine learning in the understanding and design of materials. *Journal of the American Chemical Society*, 142(48):20273–20287, 2020.
- [34] Krishnendu Mukherjee and Yamil J. Colón. Machine learning and descriptor selection for the computational discovery of metal-organic frameworks. *Molecular Simulation*, pages 1–21, 2021.
- [35] Zenan Shi, Wenyuan Yang, Xiaomei Deng, Chengzhi Cai, Yaling Yan, Hong Liang, Zili Liu, and Zhiwei Qiao. Machine-learning-assisted high-throughput computational screening of high performance metal-organic frameworks. *Molecular Systems Design & Engineering*, 5(4):725–742, 2020.
- [36] Marcel F Langer, Alex Goeßmann, and Matthias Rupp. Representations of molecules and materials for interpolation of quantum-mechanical simulations via machine learning. *arXiv preprint arXiv:2003.12081*, 2020.
- [37] Christopher R Collins, Geoffrey J Gordon, O Anatole Von Lilienfeld, and David J Yaron. Constant size descriptors for accurate machine learning models of molecular properties. *The Journal of Chemical Physics*, 148(24):241718, 2018.
- [38] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *Proceedings of International Conference on Machine Learning (ICML)*, pages 1263–1272, 2017.
- [39] Hiroshi Ohno and Yusuke Mukae. Machine learning approach for prediction and search: Application to methane storage in a metal-organic framework. *The Journal of Physical Chemistry C*, 120(42):23963–23968, 2016.
- [40] Cory M. Simon, Rocio Mercado, Sondre K. Schnell, Berend Smit, and Maciej Haranczyk. What are the best materials to separate a xenon/krypton mixture? *Chemistry of Materials*, 27(12):4459–4475, 2015.
- [41] Seyed Mohamad Moosavi, Aditya Nandy, Kevin Maik Jablonka, Daniele Ongari, Jon Paul Janet, Peter G. Boyd, Yongjin Lee, Berend Smit, and Heather J. Kulik. Understanding the diversity of the metal-organic framework ecosystem. *Nature Communications*, 11(1), 2020.

- [42] Eun Hyun Cho, Xuepeng Deng, Changlong Zou, and Li-Chiang Lin. Machine learning-aided computational study of metal-organic frameworks for sour gas sweetening. *The Journal of Physical Chemistry C*, 2020.
- [43] Jake Burner, Ludwig Schwiedrzik, Mykhaylo Krykunov, Jun Luo, Peter G. Boyd, and Tom K. Woo. High-performing deep learning regression models for predicting low-pressure CO₂ adsorption properties of metal-organic frameworks. *The Journal of Physical Chemistry C*, 124(51):27996–28005, 2020.
- [44] Benjamin J. Bucior, N. Scott Bobbitt, Timur Islamoglu, Subhadip Goswami, Arun Gopalan, Taner Yildirim, Omar K. Farha, Neda Bagheri, and Randall Q. Snurr. Energy-based descriptors to rapidly predict hydrogen storage in metal-organic frameworks. *Molecular Systems Design & Engineering*, 4(1):162–174, 2019.
- [45] Ryther Anderson, Jacob Rodgers, Edwin Argueta, Achay Biong, and Diego A. Gómez-Gualdrón. Role of pore chemistry and topology in the CO₂ capture capabilities of MOFs: From molecular simulation to machine learning. *Chemistry of Materials*, 30(18):6325–6337, 2018.
- [46] George S. Fanourgakis, Konstantinos Gkagkas, Emmanuel Tylianakis, and George Froudakis. A generic machine learning algorithm for the prediction of gas adsorption in nanoporous materials. *The Journal of Physical Chemistry C*, 124(13):7117–7126, 2020.
- [47] Maryam Pardakhti, Ehsan Moharrerri, David Wanik, Steven L. Suib, and Ranjan Srivastava. Machine learning using combined structural and chemical descriptors for prediction of methane adsorption performance of metal organic frameworks (MOFs). *ACS Combinatorial Science*, 19(10):640–645, 2017.
- [48] Michael Fernandez and Amanda S. Barnard. Geometrical properties can predict CO₂ and n₂ adsorption performance of metal-organic frameworks (MOFs) at low pressure. *ACS Combinatorial Science*, 18(5):243–252, 2016.
- [49] Hana Dureckova, Mykhaylo Krykunov, Mohammad Zein Aghaji, and Tom K. Woo. Robust machine learning models for predicting high CO₂ working capacity and CO₂/H₂ selectivity of gas adsorption in metal organic frameworks for precombustion carbon capture. *The Journal of Physical Chemistry C*, 123(7):4133–4139, 2019.
- [50] Xiangyu Zhang, Jing Cui, Kexin Zhang, Jiasheng Wu, and Yongjin Lee. Machine learning prediction on properties of nanoporous materials utilizing pore geometry barcodes. *Journal of Chemical Information and Modeling*, 59(11):4636–4644, 2019.
- [51] Zhao Li, Benjamin J. Bucior, Haoyuan Chen, Maciej Haranczyk, J. Ilja Siepmann, and Randall Q. Snurr. Machine learning using host/guest energy histograms to predict adsorption in metal-organic frameworks: Application to short alkanes and Xe/Kr mixtures. *The Journal of Chemical Physics*, 155(1):014701, 2021.
- [52] Tu C. Le and David A. Winkler. Discovery and optimization of materials using evolutionary approaches. *Chemical Reviews*, 116(10):6107–6132, 2016.

- [53] Yongchul G. Chung, Diego A. Gómez-Gualdrón, Peng Li, Karson T. Leperi, Pravas Deria, Hongda Zhang, Nicolaas A. Vermeulen, J. Fraser Stoddart, Fengqi You, Joseph T. Hupp, Omar K. Farha, and Randall Q. Snurr. In silico discovery of metal-organic frameworks for precombustion CO₂ capture using a genetic algorithm. *Science Advances*, 2(10):e1600909, 2016.
- [54] Sean P. Collins, Thomas D. Daff, Sarah S. Piotrkowski, and Tom K. Woo. Materials design by evolutionary optimization of functional groups in metal-organic frameworks. *Science Advances*, 2(11):e1600954, 2016.
- [55] Seyed Mohamad Moosavi, Arunraj Chidambaram, Leopold Talirz, Maciej Haranczyk, Kyriakos C. Stylianou, and Berend Smit. Capturing chemical intuition in synthesis of metal-organic frameworks. *Nature Communications*, 10(1), 2019.
- [56] Cameron B Browne, Edward Powley, Daniel Whitehouse, Simon M Lucas, Peter I Cowling, Philipp Rohlfshagen, Stephen Tavener, Diego Perez, Spyridon Samothrakis, and Simon Colton. A survey of monte carlo tree search methods. *IEEE Transactions on Computational Intelligence and AI in games*, 4(1):1–43, 2012.
- [57] Xiangyu Zhang, Kexin Zhang, and Yongjin Lee. Machine learning enabled tailor-made design of application-specific metal-organic frameworks. *ACS Applied Materials & Interfaces*, 12(1):734–743, 2019.
- [58] Xiangyu Zhang, Kexin Zhang, Hyeonsuk Yoo, and Yongjin Lee. Machine learning-driven discovery of metal-organic frameworks for efficient CO₂ capture in humid condition. *ACS Sustainable Chemistry & Engineering*, 2021.
- [59] David A Cohn, Zoubin Ghahramani, and Michael I Jordan. Active learning with statistical models. *Journal of Artificial Intelligence Research (JAIR)*, 4:129–145, 1996.
- [60] Eugene N. Muratov, Jürgen Bajorath, Robert P. Sheridan, Igor V. Tetko, Dmitry Filimonov, Vladimir Poroikov, Tudor I. Oprea, Igor I. Baskin, Alexandre Varnek, Adrian Roitberg, Olexandr Isayev, Stefano Curtalolo, Denis Fourches, Yoram Cohen, Alan Aspuru-Guzik, David A. Winkler, Dimitris Agrafiotis, Artem Cherkasov, and Alexander Tropsha. QSAR without borders. *Chemical Society Reviews*, 49(11):3525–3564, 2020.
- [61] Tu Le, V. Chandana Epa, Frank R. Burden, and David A. Winkler. Quantitative structure-property relationship modeling of diverse materials properties. *Chemical Reviews*, 112(5):2889–2919, 2012.
- [62] Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P Adams, and Nando De Freitas. Taking the human out of the loop: A review of bayesian optimization. *Proceedings of the IEEE*, 104(1):148–175, 2015.
- [63] Apoorv Agnihotri and Nipun Batra. Exploring bayesian optimization. *Distill*, 2020. <https://distill.pub/2020/bayesian-optimization>.
- [64] Albert P. Bartók, Risi Kondor, and Gábor Csányi. On representing chemical environments. *Physical Review B*, 87(18), 2013.

- [65] Lev Sarkisov, Rocio Bueno-Perez, Mythili Sutharson, and David Fairen-Jimenez. Materials informatics with PoreBlazer v4.0 and the CSD MOF database. *Chemistry of Materials*, 32(23):9849–9867, 2020.
- [66] Thomas F Willems, Chris H Rycroft, Michael Kazi, Juan C Meza, and Maciej Haranczyk. Algorithms and tools for high-throughput geometry-based analysis of crystalline porous materials. *Microporous and Mesoporous Materials*, 149(1):134–141, 2012.
- [67] Lev Sarkisov and Alex Harrison. Computational structure characterisation tools in application to ordered and disordered porous materials. *Molecular Simulation*, 37(15):1248–1257, 2011.
- [68] Andrew S. Rosen, Shaelyn M. Iyer, Debmalya Ray, Zhenpeng Yao, Alán Aspuru-Guzik, Laura Gagliardi, Justin M. Notestein, and Randall Q. Snurr. Machine learning the quantum-chemical properties of metal–organic frameworks for accelerated materials discovery. *Matter*, 2021.
- [69] Tian Xie and Jeffrey C. Grossman. Crystal graph convolutional neural networks for an accurate and interpretable prediction of material properties. *Physical Review Letters*, 120(14), 2018.
- [70] Ali Raza, Faaiz Waqar, Arni Sturluson, Cory Simon, and Xiaoli Fern. Towards explainable message passing networks for predicting carbon dioxide adsorption in metal-organic frameworks. *arXiv preprint arXiv:2012.03723*, 2020.
- [71] Eun Hyun Cho and Li-Chiang Lin. Nanoporous material recognition via 3d convolutional neural networks: Prediction of adsorption properties. *The Journal of Physical Chemistry Letters*, 12(9):2279–2285, 2021.
- [72] Arni Sturluson, Melanie T. Huynh, Arthur H. P. York, and Cory M. Simon. Eigencages: Learning a latent space of porous cage molecules. *ACS Central Science*, 4(12):1663–1676, 2018.
- [73] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.
- [74] Christopher R. Collins, Geoffrey J. Gordon, O. Anatole von Lilienfeld, and David J. Yaron. Constant size descriptors for accurate machine learning models of molecular properties. *The Journal of Chemical Physics*, 148(24):241718, 2018.
- [75] Zhenpeng Yao, Benjamín Sánchez-Lengeling, N. Scott Bobbitt, Benjamin J. Bucior, Sai Govind Hari Kumar, Sean P. Collins, Thomas Burns, Tom K. Woo, Omar K. Farha, Randall Q. Snurr, and Alán Aspuru-Guzik. Inverse design of nanoporous crystalline reticular materials with deep generative models. *Nature Machine Intelligence*, 3(1):76–86, 2021.
- [76] Michael Fernandez, Nicholas R. Trefiak, and Tom K. Woo. Atomic property weighted radial distribution functions descriptors of metal–organic frameworks for the prediction of gas uptake capacity. *The Journal of Physical Chemistry C*, 117(27):14095–14105, 2013.
- [77] Yongjin Lee, Senja D. Barthel, Paweł Dłotko, Seyed Mohamad Moosavi, Kathryn Hess, and Berend Smit. High-throughput screening approach for nanoporous materials genome using topological data analysis: Application to zeolites. *Journal of Chemical Theory and Computation*, 14(8):4427–4437, 2018.

- [78] Felix Faber, Alexander Lindmaa, O. Anatole von Lilienfeld, and Rickard Armiento. Crystal structure representations for machine learning models of formation energies. *International Journal of Quantum Chemistry*, 115(16):1094–1101, 2015.
- [79] UC Berkeley CS188 Intro to AI – Course Materials. http://ai.berkeley.edu/lecture_slides.html.
- [80] Niranjan Srinivas, Andreas Krause, Sham Kakade, and Matthias Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. In *Proceedings of the 27th International Conference on Machine Learning (ICML)*, pages 1015–1022, 2010.
- [81] Thang D Bui, Cuong Nguyen, and Richard E Turner. Streaming sparse Gaussian process approximations. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 30. Curran Associates, Inc., 2017.
- [82] Burr Settles. *Active Learning*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2012.
- [83] Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian processes for machine learning*. Adaptive computation and machine learning. MIT Press, 2006.
- [84] Jochen Görtler, Rebecca Kehlbeck, and Oliver Deussen. A visual exploration of gaussian processes. *Distill*, 4(4):e17, 2019.
- [85] Bernhard Schölkopf and Alexander Johannes Smola. *Learning with Kernels: support vector machines, regularization, optimization, and beyond*. Adaptive computation and machine learning series. MIT Press, 2002.
- [86] Karsten Borgwardt, Elisabetta Ghisu, Felipe Llinares-López, Leslie O’Bray, and Bastian Rieck. Graph kernels: state-of-the-art and future challenges. *arXiv preprint arXiv:2011.03854*, 2020.
- [87] Kevin P Murphy. *Machine learning: a probabilistic perspective*. MIT press, 2012.
- [88] Florent Perronnin, Jorge Sánchez, and Thomas Mensink. Improving the fisher kernel for large-scale image classification. In Kostas Daniilidis, Petros Maragos, and Nikos Paragios, editors, *Computer Vision - ECCV 2010, 11th European Conference on Computer Vision, Heraklion, Crete, Greece, September 5-11, 2010, Proceedings, Part IV*, volume 6314 of *Lecture Notes in Computer Science*, pages 143–156. Springer, 2010.
- [89] Zi Wang and Stefanie Jegelka. Max-value entropy search for efficient bayesian optimization. In *International Conference on Machine Learning*, pages 3627–3635. PMLR, 2017.
- [90] Maximilian Balandat, Brian Karrer, Daniel R. Jiang, Samuel Daulton, Benjamin Letham, Andrew Gordon Wilson, and Eytan Bakshy. BoTorch: A Framework for Efficient Monte-Carlo Bayesian Optimization. In *Advances in Neural Information Processing Systems 33*, 2020.
- [91] Nikolaus Hansen. The cma evolution strategy: a comparing review. *Towards a new evolutionary computation*, pages 75–102, 2006.
- [92] Nikolaus Hansen. The cma evolution strategy: A tutorial. *arXiv preprint arXiv:1604.00772*, 2016.

- [93] Daniel Cosmin Porumbel, Jin-Kao Hao, and Fred Glover. A simple and effective algorithm for the maxmin diversity problem. *Annals of Operations Research*, 186(1):275, 2011.
- [94] Quoc Phong Nguyen, Sebastian Tay, Bryan Kian Hsiang Low, and Patrick Jaillet. Top-k ranking Bayesian optimization. In *Thirty-Fifth AAAI Conference on Artificial Intelligence (AAAI)*, pages 9135–9143. AAAI Press, 2021.
- [95] David E. Graff, Eugene I. Shakhnovich, and Connor W. Coley. Accelerating high-throughput virtual screening through molecular pool-based active learning. *Chemical Science*, 12(22):7866–7881, 2021.
- [96] R. L. Greenaway, V. Santolini, M. J. Bennison, B. M. Alston, C. J. Pugh, M. A. Little, M. Miklitz, E. G. B. Eden-Rump, R. Clowes, A. Shakil, H. J. Cuthbertson, H. Armstrong, M. E. Briggs, K. E. Jelfs, and A. I. Cooper. High-throughput discovery of organic cages and catenanes using computational screening fused with robotic synthesis. *Nature Communications*, 9(1), 2018.
- [97] Connor W. Coley, Dale A. Thomas, Justin A. M. Lummiss, Jonathan N. Jaworski, Christopher P. Breen, Victor Schultz, Travis Hart, Joshua S. Fishman, Luke Rogers, Hanyu Gao, Robert W. Hicklin, Pieter P. Plehiers, Joshua Byington, John S. Piotti, William H. Green, A. John Hart, Timothy F. Jamison, and Klavs F. Jensen. A robotic platform for flow synthesis of organic compounds informed by AI planning. *Science*, 365(6453):eaax1566, 2019.
- [98] B. P. MacLeod, F. G. L. Parlane, T. D. Morrissey, F. Häse, L. M. Roch, K. E. Dettelbach, R. Moreira, L. P. E. Yunker, M. B. Rooney, J. R. Deeth, V. Lai, G. J. Ng, H. Situ, R. H. Zhang, M. S. Elliott, T. H. Haley, D. J. Dvorak, A. Aspuru-Guzik, J. E. Hein, and C. P. Berlinguette. Self-driving laboratory for accelerated discovery of thin-film materials. *Science Advances*, 6(20):eaaz8867, 2020.
- [99] Rebecca L. Greenaway and Kim E. Jelfs. Integrating computational and experimental workflows for accelerated organic materials discovery. *Advanced Materials*, 33(11):2004831, 2021.
- [100] Gael Donval, Calum Hand, James Hook, Emiko Dupont, Malena Sabate Landman, Malina Freitag, Matthew Lennox, and Tina Düren. Autonomous exploration and identification of high performing adsorbents using active learning. 2021.
- [101] Javad Azimi, Alan Fern, and Xiaoli Z. Fern. Batch Bayesian optimization via simulation matching. In John D. Lafferty, Christopher K. I. Williams, John Shawe-Taylor, Richard S. Zemel, and Aron Culotta, editors, *Advances in Neural Information Processing Systems 23: 24th Annual Conference on Neural Information Processing Systems 2010. Proceedings of a meeting held 6-9 December 2010, Vancouver, British Columbia, Canada*, pages 109–117. Curran Associates, Inc., 2010.
- [102] Tarun Kathuria, Amit Deshpande, and Pushmeet Kohli. Batched Gaussian process bandit optimization via determinantal point processes. In Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 4206–4214, 2016.
- [103] Christof Angermüller, David Belanger, Andreea Gane, Zelda Mariet, David Dohan, Kevin Murphy, Lucy Colwell, and D. Sculley. Population-based black-box optimization for biological se-

- quence design. In *Proceedings of the 37th International Conference on Machine Learning ICML*, volume 119 of *Proceedings of Machine Learning Research*, pages 324–334. PMLR, 2020.
- [104] Aryan Deshwal, Syrine Belakaria, and Janardhan Rao Doppa. Mercer features for efficient combinatorial bayesian optimization. In *Thirty-Fifth AAAI Conference on Artificial Intelligence (AAAI)*, pages 7210–7218. AAAI Press, 2021.
- [105] Florian Häse, Loïc M. Roch, Christoph Kreisbeck, and Alán Aspuru-Guzik. Phoenix: A Bayesian optimizer for chemistry. *ACS Central Science*, 4(9):1134–1145, 2018.
- [106] Rémi Lam, Douglas L Allaire, and Karen Willcox. Multifidelity optimization using statistical surrogate modeling for non-hierarchical information sources. In *56th AIAA/ASCE/AHS/ASC Structures, Structural Dynamics, and Materials Conference*, 2015.
- [107] Kirthevasan Kandasamy, Gautam Dasarathy, Junier B Oliva, Jeff Schneider, and Barnabás Poczos. Gaussian process bandit optimisation with multi-fidelity evaluations. In *Conference on Neural Information Processing Systems*, 2016.
- [108] Yehong Zhang, Trong Nghia Hoang, Bryan Kian Hsiang Low, and Mohan Kankanhalli. Information-based multi-fidelity Bayesian optimization. In *Conference on Neural Information Processing Systems Workshop on Bayesian Optimization*, 2017.
- [109] Jialin Song, Yuxin Chen, and Yisong Yue. A general framework for multi-fidelity Bayesian optimization with Gaussian processes. *International Conference on Artificial Intelligence and Statistics*, 2019.
- [110] Shion Takeno, Hitoshi Fukuoka, Yuhki Tsukada, Toshiyuki Koyama, Motoki Shiga, Ichiro Takeuchi, and Masayuki Karasuyama. Multi-fidelity Bayesian optimization with max-value entropy search. *arXiv:1901.08275*, 2019.
- [111] Syrine Belakaria, Aryan Deshwal, and Janardhan Rao Doppa. Multi-fidelity multi-objective Bayesian optimization: An output space entropy search approach. In *AAAI conference on Artificial Intelligence (AAAI)*, 2020.
- [112] Syrine Belakaria, Aryan Deshwal, and Janardhan Rao Doppa. Information-theoretic multi-objective Bayesian optimization with continuous approximations. *CoRR*, abs/2009.05700, 2020.
- [113] H. C. Herbol, P. Clancy, and M. Poloczek. Cost-effective materials discovery: Bayesian optimization across information sources. *Mater. Horiz.*, 7:2113–2123, 2020.
- [114] Olga Egorova, Roohollah Hafizi, David C. Woods, and Graeme M. Day. Multifidelity statistical machine learning for molecular crystal structure prediction. *The Journal of Physical Chemistry A*, 124(39):8065–8078, 2020. PMID: 32881496.
- [115] Anh Tran, Julien Tranchida, Tim Wildey, and Aidan P Thompson. Multi-fidelity machine-learning with uncertainty quantification and bayesian optimization for materials design: Application to ternary random alloys. *The Journal of Chemical Physics*, 153(7):074705, 2020.
- [116] Daniel Hernández-Lobato, Jose Hernandez-Lobato, Amar Shah, and Ryan Adams. Predictive entropy search for multi-objective Bayesian optimization. In *ICML*, pages 1492–1501, 2016.

- [117] Syrine Belakaria, Aryan Deshwal, and Janardhan Rao Doppa. Max-value entropy search for multi-objective Bayesian optimization. In *NeurIPS*, 2019.
- [118] Syrine Belakaria, Aryan Deshwal, Nitthilan Kannappan Jayakodi, and Janardhan Rao Doppa. Uncertainty-aware search framework for multi-objective Bayesian optimization. In *AAAI*, 2020.
- [119] Shinya Suzuki, Shion Takeno, Tomoyuki Tamura, Kazuki Shitara, and Masayuki Karasuyama. Multi-objective Bayesian optimization using pareto-frontier entropy. In *Proceedings of International Conference on Machine Learning (ICML)*, pages 9279–9288, 2020.
- [120] Artur M. Schweidtmann, Adam D. Clayton, Nicholas Holmes, Eric Bradford, Richard A. Bourne, and Alexei A. Lapkin. Machine learning meets continuous flow chemistry: Automated optimization towards the pareto front of multiple objectives. *Chemical Engineering Journal*, 352:277–282, 2018.
- [121] Florian Häse, Loïc M Roch, and Alán Aspuru-Guzik. Chimera: enabling hierarchy based multi-objective optimization for self-driving laboratories. *Chemical science*, 9(39):7642–7655, 2018.
- [122] Syrine Belakaria, Aryan Deshwal, and Janardhan Rao Doppa. Output space entropy search framework for multi-objective Bayesian optimization. *Journal of Artificial Intelligence Research (JAIR)*, 2021.
- [123] Kevin Maik Jablonka, Giriprasad Melpatti Jothiappan, Shefang Wang, Berend Smit, and Brian Yoo. Bias free multiobjective active learning for materials design and discovery. *Nature communications*, 12(1):1–10, 2021.
- [124] Valerio Perrone, Iaroslav Shcherbatyi, Rodolphe Jenatton, Cédric Archambeau, and Matthias W. Seeger. Constrained bayesian optimization with max-value entropy search. *CoRR*, abs/1910.07003, 2019.
- [125] Syrine Belakaria, Aryan Deshwal, and Janardhan Rao Doppa. Max-value entropy search for multi-objective Bayesian optimization with constraints. *CoRR*, abs/2009.01721, 2020.
- [126] Syrine Belakaria, Aryan Deshwal, and Janardhan Rao Doppa. Uncertainty aware search framework for multi-objective Bayesian optimization with constraints. *CoRR*, abs/2008.07029, 2020.
- [127] Zhiyuan Zhou, Syrine Belakaria, Aryan Deshwal, Wookpyo Hong, Janardhan Rao Doppa, Partha Pratim Pande, and Deukhyoun Heo. Design of multi-output switched-capacitor voltage regulator via machine learning. In *DATE*, 2020.
- [128] Syrine Belakaria, Derek Jackson, Yue Cao, Janardhan Rao Doppa, and Xiaonan Lu. Machine learning enabled fast multi-objective optimization for electrified aviation power system design. In *ECCE*, 2020.
- [129] Eric Hans Lee, Valerio Perrone, Cédric Archambeau, and Matthias W. Seeger. Cost-aware Bayesian optimization. *CoRR*, abs/2003.10870, 2020.
- [130] Gauthier Guinet, Valerio Perrone, and Cédric Archambeau. Pareto-efficient acquisition functions for cost-aware Bayesian optimization. *CoRR*, abs/2011.11456, 2020.

- [131] Matteo Aldeghi, Florian Häse, Riley J Hickman, Isaac Tamblyn, and Alán Aspuru-Guzik. Golem: An algorithm for robust experiment and process optimization. *arXiv preprint arXiv:2103.03716*, 2021.
- [132] Lukas Fröhlich, Edgar Klenske, Julia Vinogradska, Christian Daniel, and Melanie Zeilinger. Noisy-input entropy search for efficient robust Bayesian optimization. In *International Conference on Artificial Intelligence and Statistics*, pages 2262–2272. PMLR, 2020.
- [133] Maximilian Balandat, Brian Karrer, Daniel R. Jiang, Samuel Daulton, Benjamin Letham, Andrew Gordon Wilson, and Eytan Bakshy. BoTorch: A Framework for Efficient Monte-Carlo Bayesian Optimization. In *Advances in Neural Information Processing Systems 33*, 2020.
- [134] Ruben Martinez-Cantin. BayesOpt: a Bayesian optimization library for nonlinear optimization, experimental design and bandits. *J. Mach. Learn. Res.*, 15(1):3735–3739, 2014.
- [135] F. Hutter, H. H. Hoos, and K. Leyton-Brown. Sequential model-based optimization for general algorithm configuration. Technical Report TR-2010-10, University of British Columbia, Department of Computer Science, 2010.
- [136] Tsuyoshi Ueno, Trevor David Rhone, Zhufeng Hou, Teruyasu Mizoguchi, and Koji Tsuda. COMBO: An efficient Bayesian optimization library for materials science. *Materials Discovery*, 4:18–21, 2016.
- [137] Qiaohao Liang, Aldair E Gongora, Zekun Ren, Armi Tiihonen, Zhe Liu, Shijing Sun, James R Deaneault, Daniil Bash, Flore Mekki-Berrada, Saif A Khan, Kedar Hippalgaonkar, Benji Maruyama, Keith A. Brown, John Fisher III, and Tonio Buonassisi. Benchmarking the performance of Bayesian optimization across multiple experimental materials science domains. *arXiv preprint arXiv:2106.01309*, 2021.
- [138] Janardhan Rao Doppa. Adaptive experimental design for optimizing combinatorial structures. In *Proceedings of International Joint Conference on Artificial Intelligence (IJCAI)*, 2021.
- [139] Benjamin J. Shields, Jason Stevens, Jun Li, Marvin Parasram, Farhan Damani, Jesus I. Martinez Alvarado, Jacob M. Janey, Ryan P. Adams, and Abigail G. Doyle. Bayesian reaction optimization as a tool for chemical synthesis. *Nature*, 590(7844):89–96, 2021.
- [140] Artur M. Schweidtmann, Adam D. Clayton, Nicholas Holmes, Eric Bradford, Richard A. Bourne, and Alexei A. Lapkin. Machine learning meets continuous flow chemistry: Automated optimization towards the pareto front of multiple objectives. *Chemical Engineering Journal*, 352:277–282, 2018.
- [141] Benjamin Burger, Phillip M. Maffettone, Vladimir V. Gusev, Catherine M. Aitchison, Yang Bai, Xiaoyan Wang, Xiaobo Li, Ben M. Alston, Buyi Li, Rob Clowes, Nicola Rankin, Brandon Harris, Reiner Sebastian Sprick, and Andrew I. Cooper. A mobile robotic chemist. *Nature*, 583(7815):237–241, 2020.
- [142] Tomoki Yamashita, Nobuya Sato, Hiori Kino, Takashi Miyake, Koji Tsuda, and Tamio Oguchi. Crystal structure prediction accelerated by Bayesian optimization. *Physical Review Materials*, 2(1), 2018.

- [143] Edward Pyzer-Knapp, Graeme Day, Linjiang Chen, and Andrew I Cooper. Distributed multi-objective Bayesian optimization for the intelligent navigation of energy structure function maps for efficient property discovery. *ChemRxiv*, 2020.
- [144] Dezhen Xue, Prasanna V. Balachandran, John Hogden, James Theiler, Deqing Xue, and Turab Lookman. Accelerated search for materials with targeted properties by adaptive design. *Nature Communications*, 7(1), 2016.
- [145] R. A. Vargas-Hernández. Bayesian optimization for calibrating and selecting hybrid-density functional models. *The Journal of Physical Chemistry A*, 124(20):4053–4061, 2020.
- [146] Connor W. Coley. Defining and exploring chemical spaces. *Trends in Chemistry*, 3(2):133–145, 2021.
- [147] Kei Terayama, Masato Sumita, Ryo Tamura, and Koji Tsuda. Black-box optimization for automated discovery. *Accounts of Chemical Research*, 54(6):1334–1346, 2021.
- [148] Peter I Frazier and Jialei Wang. Bayesian optimization for materials design. In *Information Science for Materials Discovery and Design*, pages 45–75. Springer, 2016.
- [149] Turab Lookman, Prasanna V Balachandran, Dezhen Xue, and Ruihao Yuan. Active learning in materials science with emphasis on adaptive sampling using uncertainties for targeted design. *npj Computational Materials*, 5(1):1–17, 2019.
- [150] Florian Häse, Loïc M Roch, and Alán Aspuru-Guzik. Gryffin: An algorithm for Bayesian optimization for categorical variables informed by physical intuition with applications to chemistry. *arXiv preprint arXiv:2003.12127*, 2020.
- [151] Aryan Deshwal, Syrine Belakaria, and Janardhan Rao Doppa. Bayesian optimization over hybrid spaces. In *Proceedings of International Conference on Machine Learning (ICML)*, 2021.