# Inverse Design of Composite Metal Oxide Optical Materials based on Deep Transfer Learning and Global Optimization

**Rongzhi Dong · Yabo Dan · Xiang Li ·
Jianjun Hu(✉)**

**Abstract** Optical materials with special optical properties are widely used in a broad span of technologies, from computer displays to solar energy utilization leading to large dataset accumulated from years of extensive materials synthesis and optical characterization. Previously, machine learning models have been developed to predict the optical absorption spectrum from a materials characterization image or vice versa. Herein we propose TLOpt, a transfer learning based inverse optical materials design algorithm for suggesting material compositions with a desired target light absorption spectrum. Our approach is based on the combination of a deep neural network model and global optimization algorithms including a genetic algorithm and Bayesian optimization. A transfer learning strategy is employed to solve the small dataset issue in training the neural network predictor of optical absorption spectrum using the Magpie materials composition descriptor. Our extensive experiments show that our algorithm can inverse design the materials composition with stoichiometry with high accuracy. The source code is freely available at `https://github.com/usccolumbia/TLOpt`

Rongzhi Dong
School of mechanical engineering, Guizhou university, Guiyang 550025, China
E-mail: askemma@cau.edu.cn

Yabo Dan
School of mechanical engineering, Guizhou university, Guiyang 550025, China
E-mail: yabodan152@163.com

Xiang Li
School of mechanical engineering, Guizhou university, Guiyang 550025, China
E-mail: xiangli_0214@163.com

Jianjun Hu(✉)
Department of Computer Science and Engineering, University of South Carolina, Columbia SC 29201
E-mail: jianjunh@csc.sc.edu

## 1 Introduction

Optical materials play a major role in imaging, communications, solar cell and sensor design, and the absorption spectra of meterials are the key research object for these applications. The optical properties of composite metal oxides are an interesting area of optical materials research because these properties may be very different from the properties of individual components. The structure of a composite material depends on its preparation process and the chemical properties of its constituent elements, which further affect its optical properties. Constituent elements of a composite material and the mole ratios among elements also have large impact on its structure and optical properties of the material.

In order to study the optical properties of composite materials, first principles calculations such as Density Functional Theory (DFT) have been widely used [1, 2]. Although first principles calculations are powerful, they are susceptible to the constraints of their excessive calculation cost, which limits the size of the material design space or the number of materials they can screen. To address this problem, machine learning (ML) has been increasingly applied to materials science fields, leading to the emergence of "materials informatics" [3], in which materials learning methods are developed to obtain prior knowledge and predictive models from known material dataset, and then predict complex material properties based on these models. In the past few years, ML has succeeded in predicting new features [4], guiding chemical synthesis and discovering suitable compounds with target properties [5,6,7,8]. While ML based material property prediction models can be used to screen known materials database to find candidates with expected properties, its performance is limited by the available materials, which are not developed for the target properties anyway. When ideal materials are not found in existing databases, discovering and synthesizing new materials with target properties is needed, which is usually based on the experience and knowledge of the researchers and expensive experiments. As materials that can be easily found have been found already and the scope of experimental exploration based on experience is narrow, new methods of material discovery is needed [9], which promoted the development of the inverse material design approaches [10].

Inverse design started in the field of alloy design [11], using genetic algorithm and molecular dynamics simulations to optimize the composition of multi-component alloys. This method received widespread attention in various fields once it was proposed, and now widely used in nanophotonic design [12, 13,14,15,16], surface design [16,17,18,19] catalyst design [20], catalyst design [10,21], drug design [6] and materials design [9]. Inverse design of materials with desirable optical properties is of great significance in many industries such as solar cells, computer monitors, and optical microscopes[22]. Inverse

materials design can be regarded as an optimization problem, in which the desired materials characteristics are used as the optimization objectives. If the material is designed for optimizing only one attribute, it can be formulated as a single-objective optimization problem; if multiple materials attributes need to be optimized at the same time, it can be regarded as a multi-objective optimization problem. There are two major modules in a typical inverse design framework: one is the sampling module to guide the search in the design space in which a variety of optimization algorithms [23] can be used such as genetic algorithm (GA) [24], Bayesian optimization (BO) [25], particle swarm optimization (PSO) [26], and differential evolution (DE) [27, 28]. The other module is the forward property prediction model, which evaluates the performance of each design candidate, for which a set of commonly machine learning algorithms have been used includinig support vector machines (SVM) [29], random forest (RF) [30] and artificial neural network (ANN) [31] etc.

However, few studies have focused on inverse design of optical materials, especially inferring the possible compound formula of the materials only based on their absorption spectrum. The main reason of this phenomenon is lack of large-scale optical characterizations dataset of materials for model training. Previously, a metal oxide optical characterization dataset was published, and autoencoder algorithms for measured optical properties of metal oxides based on this dataset [32, 33] have been developed, which can map composite materials' characterization image patterns to its UV-vis absorption spectrum and vice versa. Comparing the band gap energy from the truth spectra to the predicted spectra, the root mean squared error and mean absolute error are 261 meV and 180 meV respectively, which are very small errors. However, it is not clear how to map the characterization images back to the composite material compositions. Yu et al. [34] proposed a spectroscopic limited maximum efficiency metric, which can be used to guide the search of very thin film photovoltaic devices with high absorption. Inverse design of materials compositions have also been proposed by using generative adversarial networks [35], which is mainly done by screening a large set of generated hypothetical materials.

Our work focuses on the inverse design problem of optical materials compositions with given target UV-vis absorption spectrum. In this study, we propose an algorithm to inverse design composite materials with a given target optical spectrum by combining artificial neural networks and genetic algorithms and Bayesian optimization. First, a large number of known composite material spectra are used to train a neural network model to predict the spectrum from the formula of a given composite material. The transfer learning method is used to train the spectrum prediction model of a specific set constituent elements with limited number of samples of varying stoichiometries. Finally, the metal oxide material is designed inversely based on the given target optical absorption spectrum performance using GA and BO.

Our contributions can be summarized as follows:

· We propose a transfer learning based neural network model for predicting optical absorption spectra from materials compositions. This approach helps to address the small data issue in materials property prediction.
· We develop an approach for inverse design of metal oxide material compositions for achieving a target optical absorption spectrum using both genetic algorithms and Bayesian Optimization.
· We conduct extensive experiments and show that our proposed framework is capable to achieve good performance for target spectra.

The remainder of this paper is organized as follows. Section 2 focuses on the research framework, materials representation, and inverse design models of materials. Section 3 describes our experiments and highlights our inverse design performance. The last section concludes the paper.

## 2 Materials and Methods

2.1 Problem setup and inverse design framework

In our inverse design problem, the goal is to design the materials formula of a potential optical material that can achieve the given target light absorption spectrum. As shown in **Fig.1**, we have 554 different formula groups in the dataset, and each formula group consists of formulas with the same set of elements but different mole ratios. For each independent formula, there is a corresponding absorption spectrum. According to the composition of formulas, we define two versions of the inverse design problem: 1) the elements in the material are given, only the mole ratios need to be determined; 2) the elements are not specified in advance, we need to search both the elements and their mole ratios.

The main components of the framework are shown in **Fig.2**. We use deep neural network models, a type of machine learning model to learn the relationships between material composition and light absorption spectrum. Unlike previous work [36] which explores the relationship between characterization images of different optical material structures and their light absorption properties, we address the real-world need to inversely design suitable materials (in terms of their materials composition) to achieve a specific target light absorption performance. In order to design materials according to performance specification and guide the discovery of new optical materials, we construct two inverse design models through the genetic algorithm and Bayesian optimization respectively to predict the corresponding material compound formula elements and their mole ratios based on the UV-vis absorption spectrum (220 items).

For inverse design with given elements, our framework is shown in **Fig.2** (a), which is composed of a fully connected neural network-based transfer learning model trained with Magpie features and global optimization based search model including a genetic algorithm and a Bayesian optimization. To divide the dataset, we randomly select the target formula from all 100,429
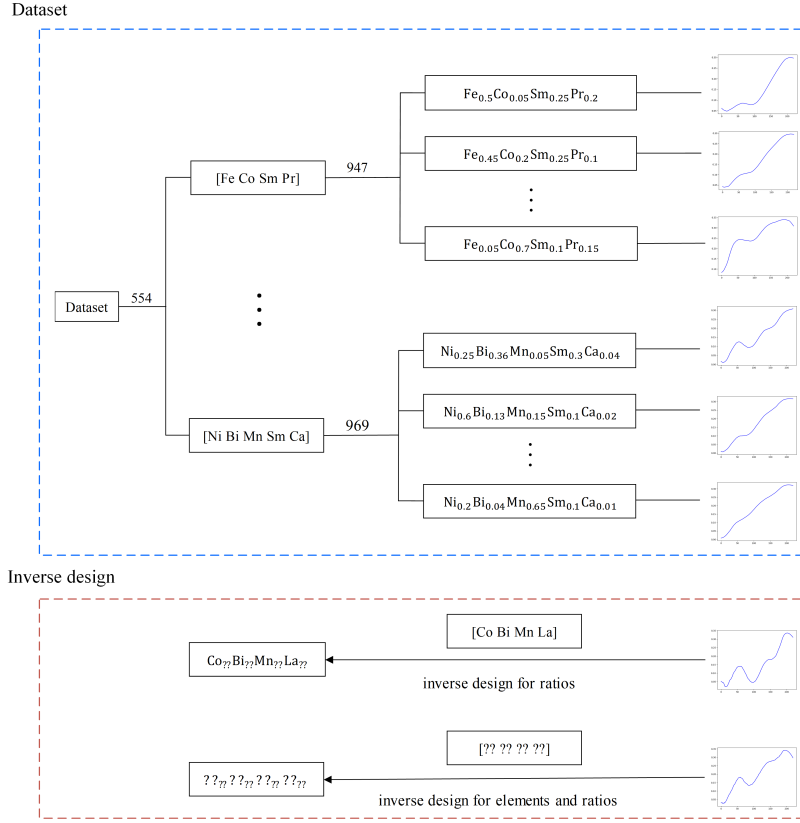
Fig. 1: The inverse composite material design problem for achieving a target optical absorption spectrum

samples, and then set the formulas with same elements but different mole ratios as the target formula and corresponding spectra as Dataset B, while remaining 553 formula-spectrum groups are set to be Dataset A. Firstly, we use large amounts of known data (Dataset A) for initial training of the fully connected neural network model 1. Then we transfer parameters of model 1 to model 2 (with the same type as model 1), and use a small amount of sample data (Dataset B) to fine-tune the model. Finally, a genetic algorithm and Bayesian optimization method are used to inverse design materials that approximate the target optical properties through spectrum fitting. For inverse design without specifying the elements, although we also divide the dataset into A and B according to the known target, the fully connected neural network model is trained only by the Dataset A and without transfer learning( **Fig.2** (b)) to ensure fair comparison with transfer learning.

(a) Transfer learning framework for computational inverse design of optical materials



(b) Framework without transfer learning for computational inverse design of optical materials
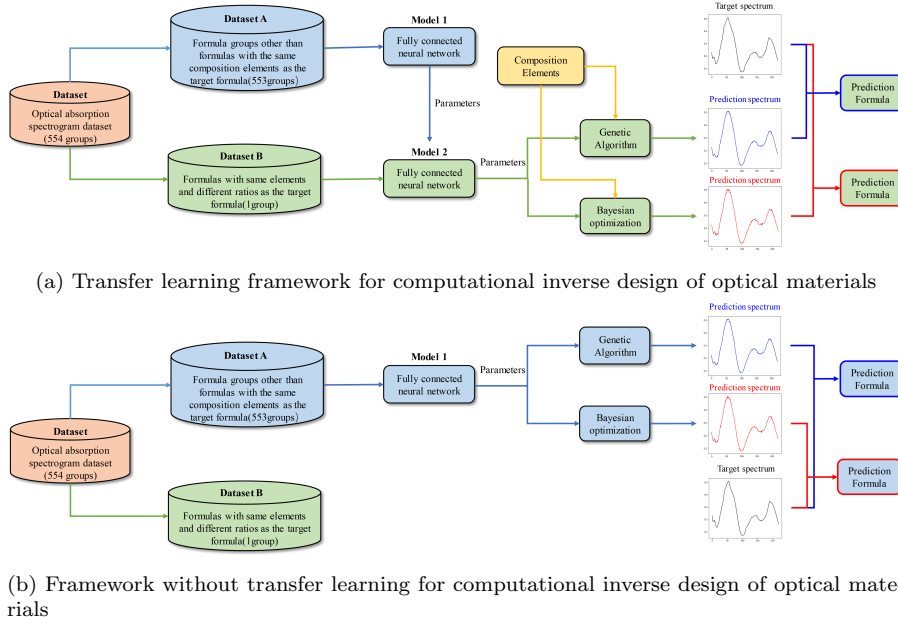
Fig. 2: Frameworks for computational inverse design of optical materials

## 2.2 Materials datasets

The material data used in this study are downloaded from references [32, 33]. This dataset contains sample images, UV-vis spectra and composition of a large set of metal oxide materials selected from the Materials Experiment and Analysis Database (MEAD) of the High-Throughput Experimentation (HTE) group at the Joint Center for Artificial Photosynthesis at Caltech. It is one of the largest publicly available scientific data set of cured metal oxide materials, which are synthesized by metal nitric acid salt with annealing. The absorption spectra are recorded using real-time scanning UV-vis dual-ball spectrometer while a flatbed scanner is used to obtain the sample images. This database contains a total of 178,994 molding material samples and their corresponding optical absorbance values at 220 energies between 1.32 to 3.2eV. The metal oxide samples contain various combinations of 1 to 5 cationic elements, as well as various inkjet printing and heat treatment parameters. These parameters are not used in the model described in this study. Since different preparation processes in the metal nitric acid salt produce many repetitive compound formulas, for materials with a fixed composition, we choose the average photoconductivity under different processes as its reference photoconductivity. After screening, we got a total of 100,429 samples composed of 42 different elements, which are then divided into 554 groups according to their constituent elements. Each group consists of a series of material formulas with the same composition elements and different mole ratios. Our dataset includes 2 groups

of quinary compounds, 114 groups of quaternary compounds, 216 groups of ternary compounds, 181 groups of binary compounds, and 41 groups of simple substances. **Fig.3** shows the distribution of element groups in binary, ternary and quaternary compound materials.



(a) Number of samples for binary element groups

(b) Number of samples for ternary element groups



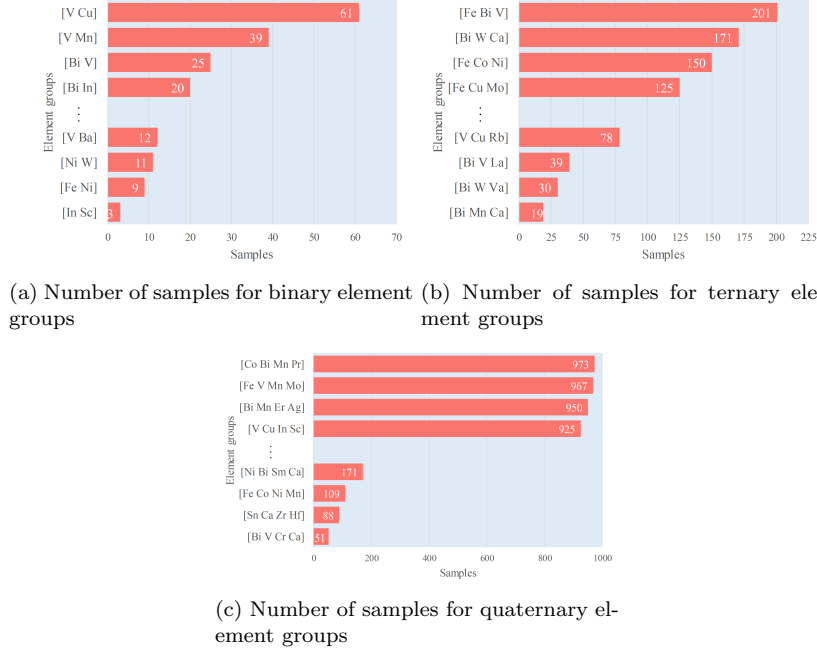(c) Number of samples for quaternary element groups

Fig. 3: Distribution of element groups in binary, ternary and quaternary compound materials

### 2.3 Materials representation

We use the Materials Agnostic Platform for Informatics and Exploration (Magpie) composition descriptor [37] to represent the materials in our dataset, which are calculated from properties of the atom elements in compound formulas to characterize materials. As discussed in the recent review [38], using elemental physical properties as descriptors for structure yields reasonably good performance in predicting various properties. Although composition-based features are unable to distinguish between crystalline polymorphs and molecular isomers/conformers, our data structure type is relatively stable, and there are no crystalline polymorphs and molecular isomers/conformers. For the inverse design process, in order to expand the search space, we should not limit the structure of the material. A Magpie descriptor vector is composed of a set of

Table 1: Model parameters of the fully connected neural network network

| Layer | Input Shape | Output Shape |
|-------|-------------|--------------|
| Fc1 | [batch, input] | [batch, 256] |
| Fc2 | [batch, 256] | [batch, 128] |
| Fc3 | [batch, 128] | [batch, 64] |
| Fc4 | [batch, 64] | [batch, 32] |
| Fc5 | [batch, 32] | [batch, 16] |
| Fc6 | [batch, 16] | [batch, 32] |
| Fc7 | [batch, 32] | [batch, 64] |
| Fc8 | [batch, 64] | [batch, 128] |
| Fc9 | [batch, 128] | [batch, 220] |

statistics of a selected element properties created by Ward [37] and can be used for representing materials with any number of constituent elements. The set of properties is broad enough to capture a widely variety of physical and chemical properties which can be used to create predictive models of many material properties given only composition [39]. The physicochemical properties include stoichiometric properties (depending only on the ratios between elements), element properties (atomic number, atomic radius, melting temperature, etc.), electronic structure properties (valence electron number of s, p, d, and f layers) and ionic compound characteristics. To construct a Magpie feature, 22 weighted element attributes of the compound formula are calculated, and then the minimum, maximum, difference, average, variance and mode characteristics are calculated for each attribute. Finally, the material is characterized as a 132-dimensional data input. Here, the elemental properties are taken from the dataset available in the Wolfram programming language [40].

2.4 Model 1 for spectrum prediction from composition: Initial training

First, we select one target compound formula that needs to be inverse engineered. Except for the formulas in the same group as the target compound formula, we randomly choose 50 compound formulas for each group in Dataset A (select all formulas if less than 50), and then randomly divide them into the training set and the test set according to 70%:30% for performance evaluation. The spectrum prediction model 1 used in this study is a fully connected multi-layer perceptron neural network model. Since our input has 132 dimensions and the expected output has 220 dimensions, the nodes in the first and last layer are determined. From experience, the greater the depth, the better the generalization of tasks, but too many layers often lead to overfitting. Therefore, we chose the number of layers and nodes per layer of our neural network based on experience, and fine-tuned it based on the experimental results. The final neural network parameters are shown in Table 1. The batch size is designed by the user, in this article we selects 64 samples as a training batch.

2.5 Model 2 for spectrum prediction trained by transfer learning

In our inverse design framework, one key issue is that for a given element set, the number of training samples are too few. For example, there are only 51 samples in element set [Bi, V, Cr, Ca]. To address this issue, we use a machine learning strategy called transfer learning. Transfer learning [41] is an algorithm to improve the performance of a target task in a target domain by exploiting some knowledge acquired when solving a source task in a source domain. Transfer learning has been widely used to address small dataset issue in machine learning [42]. Usually, the source domain and the target domain or the source and target tasks are different. Transfer learning can reduce resource consumption and the time required for model training by just fine-tuning a trained model on the target task. In this study, the source domain is the compound formula groups other than the formulas with the same composition elements as the target formula in the previous step (Dataset A); and the target domain is a group of compound formulas with the same constituent elements but different element ratios as the target compound formula (Dataset B). Here both the source task and the target task are predicting light absorption performance. Our transfer learning strategy is to first train the neural network model on Dataset A and then import the parameters of this pretrained model 1 into the identical training model 2 for further training (fine-tuning). We randomly select 30 formulas in Dataset B and divide them into a training set and a validation set according to the ratio of 2:1, as the input to train model 2. The goal of fine-tuning the parameters for model 2 training is to make it more accurate to predict the composition of the target compound formula.

2.6 Genetic Algorithm

A Genetic Algorithm (GA) is a global search method proposed by Holland [43] and Rechenberg [44] inspired by the biological evolution in nature to search for optimal solutions. The algorithm transforms the optimization problem-solving process into a simulated evolution process with inheritance, mutation, selection, and crossover of chromosomal genes in biological evolution through mathematical methods and computer simulation operations. When solving more complex combinatorial optimization problems, compared to conventional optimization algorithms, better optimization results can usually be obtained. Genetic algorithms have been widely used in combinatorial optimization, machine learning, signal processing and adaptive control [45]. In the context of materials science, genetic algorithms have been widely used in crystal structure prediction [46], inverse materials design [47], and materials property prediction [48].

As Gobin and Schuth [49] clarified, the way that materials are presented in the genome may have a significant impact on the performance of evolutionary algorithms. In this study, each material is mathematically represented by a $n \times 7$ bit binary sequence (its chromosome or genotype) representing n mole

ratios of the n elements in the material composition (Figure 4). Each element's mole ratio in the compound formula is represented by a seven-bit binary string. Decoding the binary strings into decimal values, the mole ratios of the elements can be obtained, and the material composition can be determined. The binary sequences of the population can then be mutated and crossovered by specific genetic operators.

Two decoding strategies have been proposed in our study. In the first decoding approach (**Fig.4**(a)), the 7-bit string that encodes the mole ratio for each element is first decoded into a decimal value, and then each of them will be divided by the sum of these decimal values to ensure that the final sum of the decoded mole ratios to be 1. In the second unique decoding approach as shown in **Fig.4**(b), the decoding process is as follows: 1) first all 7-binary strings are decoded into decimal value and then converted into a value between [0,1]; 2) the first ratio $r_1$ will be assigned as the mole ratio of the first element, for remaining ratio $r_i$, we first compare it with $1 - \sum_{i=1}^{i-1} r_i$, then we choose the smaller one as the mole ratio of element $i$. By converting binary encoding to decimal values, the mole ratios of the elements can be obtained, and the material composition can be determined. Compared to the first decoding approach, it has the benefit that each mole ratio vector corresponds to a unique genetic binary string while the first approach allows redundant encodings for the same mole ratio vector.



(a) Redundant decoding approach
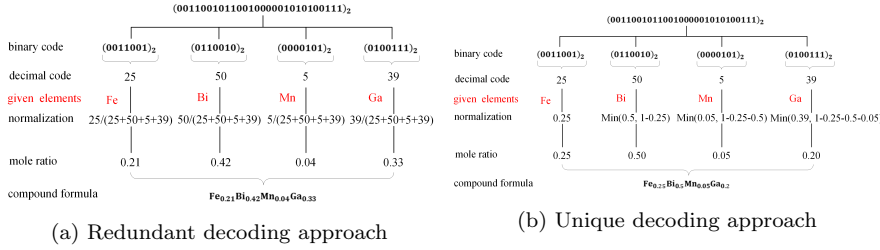
(b) Unique decoding approach

Fig. 4: Genetic encoding approaches for mole ratios with a given element set

For the inverse design problem without specifying the element set, the encoding of the GA is shown in **Fig.5**. The main difference is that the chromosome now has a block for encoding the $n$ elements, where $n$ should be specified by the user.

The initial population of individuals is usually randomly generated, but it can also seeded with suitable known materials. The generations are mutated to generate the next population in an iterative manner. The fitness of each individual in the population is evaluated by a fitness or objective function (here it is the MAE distance between the predicted spectrum for a material composition and the target spectrum). The fitness function may also include undesirable characteristics as constraints that need to be avoided. Then following the idea of survival of the fittest principle, a set of fitter materials will
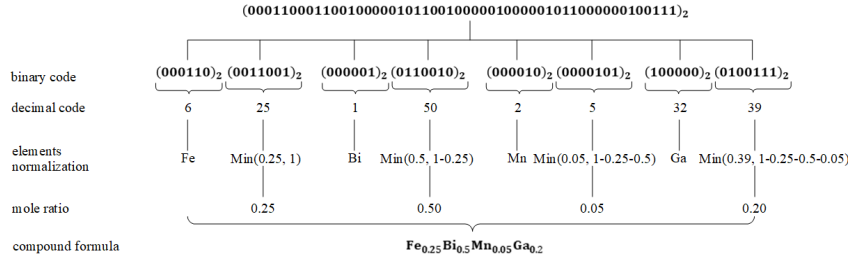
Fig. 5: Genetic encoding for inverse design without specifying the element set

be selected from the current population for breeding via mutation or crossover operations on their genomes to generate a new generation of population. This iterative cycle continues until the maximum number of generations is produced, or some members of the population have characteristics that reach the expected target. **Fig.6** summarizes the basic steps of the evolution process of a GA. The hyper-parameters of a GA include the material genome encoding length, the population size, the mutation and crossover rate, and the number of generations.
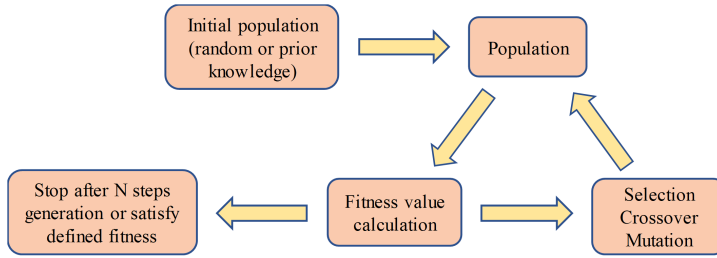


Fig. 6: Basic steps of a genetic algorithm

In this study, the genetic algorithm is used as one of the global search algorithms for inverse material design with a target light absorption performance. The evolution process starts with compound formulas with the same constituent elements but different element ratios as the target compound formula. The parameters of the genetic algorithm in this study are: gene length: 13 (6 for elements and 7 for ratios), initial population size: 500, crossover probability: 0.5, mutation probability: 0.5, generation: 100.

2.7 Bayesian optimization model

The Bayesian optimization (BO) method was proposed by [50]. Jones et al. [51] introduced an effective global optimization (EGO) method and extended

the BO technique. This method has become very popular and well-known in engineering and is now widely used in the design of time-consuming experiments, aimed at reducing experiment costs. Application of BO in machine learning mainly focuses on adjusting the hyperparameters of computationally expensive machine learning models [52]. In this study, we use BO methods to build predictive models for the potential relationships between design variables of the materials and their properties, and then use decision theory to suggest which design is most valuable. BO finds the candidate solutions that minimize the objective function by establishing a substitution function (e.g. Gaussian process model) based on the evaluation results of the objective function. The Bayesian method is different from random or grid search in that it exploits evaluated sampling points to build a surrogate model which not only predict the objective values but also related uncertainty, which allow it to achieve automated balance of exploitation and exploration. The schematic diagram of BO algorithm is shown in **Fig.7**.
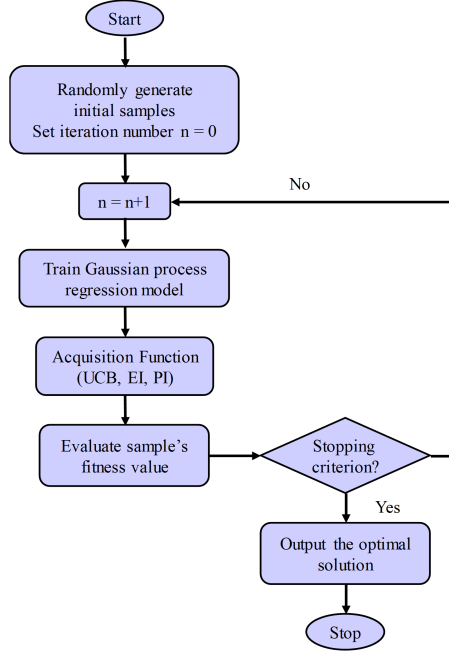


Fig. 7: The flowchart of Bayesian optimization algorithm

BO requires several initial sample points, and through Gaussian process regression (assuming that the optimization variables conform to the joint Gaussian distribution), the posterior probability distribution of first n points is calculated to obtain expectation, mean and variance. The mean represents the final expected effect of this point, the larger the mean, the greater the final

index of the model. The variance represents the uncertainty of the effect of this point, the larger the variance, the more uncertain the value of this point and worth exploring. Therefore, the first step in BO is to implement the Gaussian process regression algorithm.

Another important step is to balance exploitation and exploration. For exploitation, points close to the known points need to be selected as the reference points for the next iteration, that is, excavate points around the known points, the distribution of points will appear in a dense area, which is easy to enter the local maximum. For exploration, points far away from the known points need to be chosen as the reference points for the next iteration, and make the distribution of points as even as possible to explore the unknown area. Sampling points with large mean value can be selected for exploitation and samples with large variance can be selected for exploration. To control the ratio of exploitation and exploration, the acquisition function needs to be defined. The simplest acquisition function is Upper confidence bound algorithm (UCB) which equals the mean plus k times the variance. Where k is the adjustment parameter, which can be intuitively understood as the upper confidence boundary. More complex acquisition functions include expected improvement, entropy search, and so on.

In this study, BO algorithm is used as one of the optimization algorithms for material light absorption performance inverse design. The following configuration parameters are set: the initial population size is 500, the number of generations is 100, and the UCB acquisition function is used to achieve the balance of exploration and exploitation.

2.8 Model evaluation and experiment environment

This study uses fully connected neural networks to construct training spectrum prediction models. We choose average absolute error (MAE), root mean square error (RMSE), and coefficient of determination ($R^2$) as the evaluation criteria of these models. MAE is used to reflect the actual situation of the predicted value error, RMSE is used to measure the difference between the predicted value and the true value, $R^2$ is used to indicate the degree of fit between the predicted value and the true value. The specific calculation formulas are as follows:

$$MAE = \frac{1}{m} \sum_{i=1}^{m} |\mathbf{y}_i - \hat{\mathbf{y}}_i| \tag{1}$$

$$RMSE = \sqrt{\frac{1}{m} \sum_{i=1}^{m} (\mathbf{y}_i - \hat{\mathbf{y}}_i)^2} \tag{2}$$

$$R^2 = 1 - \frac{\sum_{i=1}^{m} (\mathbf{y}_i - \hat{\mathbf{y}}_i)^2}{\sum_{i=1}^{m} (\mathbf{y}_i - \bar{\mathbf{y}})^2} \tag{3}$$

where m is the number of samples, $\mathbf{y}_i$ and $\hat{\mathbf{y}}_i$ are the true and predicted values of the $i$ sample label (the spectrum of formula $i$), $\bar{\mathbf{y}}$ is the average of the $m$ sample real labels. All calculations were performed on a Dell Server workstation equipped with an Intel Xeon W-2123 @3.70 GHz CPU, 16 GB RAM, a Nvidia GTX1080Ti GPU with 12 GB dedicated GPU memory. Software used was Python version 3.6.4, Keras version 2.2.0, and TensorFlow version 1.14.0. The random train-test split is 70% for training and 30% for testing.

## 3 Results and discussion

### 3.1 Prediction performance of the composition descriptor based spectrum predictor

Considering the small sample data for a specific element set, we first train a fully connected neural network Model 1 on Dataset A, and then use the transfer learning method to fine tuning the model parameters of Model 2 whose initial parameters are transferred from Model 1 using Dataset B. **Fig.8** shows how the training and validation performance criteria MAE, RMSE and $R^2$ change during Model 1 training. After 800 epochs the MAE and RMSE are reduced to 0.04eV and 0.004eV and $R^2$ is raised to 0.997. This model will be used as the source model to build target prediction models for different chemical systems of specific element sets.

Since the experimental absorption spectrum usually contains noise, we discuss the performance of Model 1 under noise environment. Signal-to-noise ratio (SNR) is defined as the ratio of signal power to the noise power, often expressed in decibels detailed as follows:

$$\mathbf{SNR}_{dB} = 10\mathbf{log}_{10}\left(\mathbf{P}_{signal}/\mathbf{P}_{noise}\right) \tag{4}$$

where $\mathbf{P}_{signal}$ and $\mathbf{P}_{noise}$ are the power of the signal and the noise, respectively. In this case, the same as [53], Model 1 is tested by adding different SNR white Gaussian. The different SNR ranges from -2 dB to 4 dB. The smaller SNR value is, the stronger power of noise is. Table 2 shows the performance of Model 1 under noise experiment. It is clear that the accuracy increases as the noise gets weak, e.g., the $R^2$ is only 0.9819 when the SNR is -2 dB, while the $R^2$ surges to 0.9929 when SNR is 4 dB. Despite being affected by noise, the prediction accuracy of Model 1 is still of a high point.

Table 2: Performance under noise environment

| errors | SNR(dB) | | | | |
|---|---|---|---|---|---|
| | -2 | 0 | 2 | 4 | None |
| $R^2$ | 0.9819 | 0.9866 | 0.9902 | 0.9929 | 0.9979 |
| MAE (eV) | 0.1154 | 0.0951 | 0.0839 | 0.0714 | 0.0357 |
| RMSE (eV) | 0.1610 | 0.1384 | 0.1180 | 0.0997 | 0.0539 |

(a) MAE changes during training



(b) RMSE changes during training
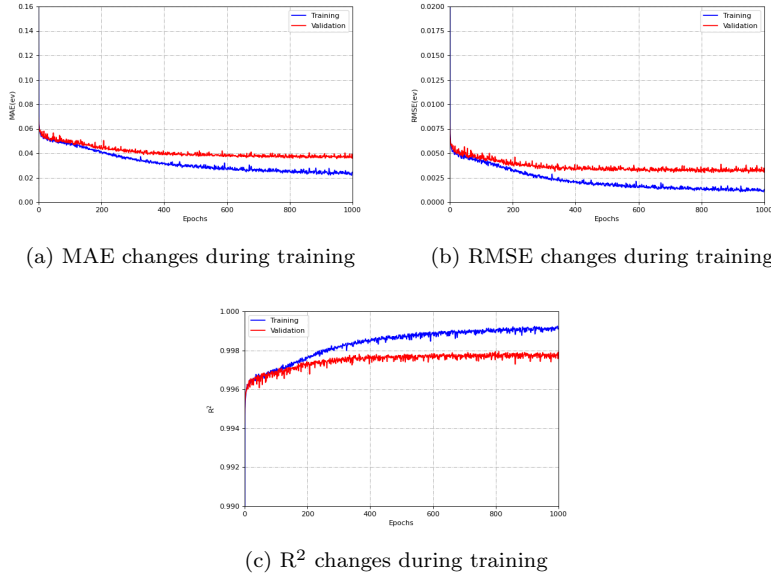


(c) $R^2$ changes during training

Fig. 8: Training and validation errors during training for composition based spectrum prediction

After training the source Model 1, we transfer the model parameters to the target chemical system with a given set of elements and fine-tune its parameters using the samples of the specified chemical system. We calculate the performance of the prediction models for target formula with few samples and with many samples to compare the performances of the models trained with or without using transfer learning strategy. In the process of training, we use randomly selected target formula as the test set. In transfer learning progress, the Dataset B is independently divided into training set, validation set and test set (only the target formula). The results are shown in Table 3. For formula group [Sn, Ca, Zr, Hf] and [Fe, Bi, V, Mn] consisting of 88 and 973 different formula samples respectively, we randomly choose two quaternary compounds $Sn_{0.3}Ca_{0.1}Zr_{0.4}Hf_{0.2}$ and $Fe_{0.25}Bi_{0.15}V_{0.3}Mn_{0.3}$ as representative test formulas with small number of samples and with large large number of samples, respectively.

Table 3: Testing errors for composition based spectrum prediction

|  | $R^2$ | MAE (eV) | RMSE (eV) |
|---|---|---|---|
| small sample set without TL | 0.9925 | 0.0936 | 0.0221 |
| small sample set with TL | 0.9947 | 0.0682 | 0.0127 |
| large sample set without TL | 0.9953 | 0.0416 | 0.0033 |
| large sample set with TL | 0.9982 | 0.0251 | 0.0012 |

3.2 Comparison of two GA encoding-decoding methods

To verify whether our proposed unique decoding approach in the GA for inverse design is better than the standard redundant decoding, we evaluate their performance on two inverse design problems of $Ni_{0.25}Bi_{0.67}Mn_{0.08}$ and $Fe_{0.05}V_{0.05}Cu_{0.25}Ca_{0.65}$ respectively. Both experiments have been run with 30,000 evaluations with pop size of 300 and generation number 100. The MAE distances of our unique decoding (blue color) and standard redundant decoding (red color) are shown in **Fig.9**.

Through **Fig.9**(a) we found that our decoding approach first achieved a lower MAE distance of 0.003488eV in the 27th generation and reached the minimum MAE 0.003484eV in the 68th generation, while MAE of the traditional decoding method converged to 0.003492eV in the 17th generation. As presented in **Fig.9**(b), the minimum MAE of both decoding approaches reached 0.004507eV in 15th generation vs 23rd generation respectively, and our method converged faster.

Both experiments proved that our unique decoding approach performed better than the standard decoding, not only reducing the redundant search space but also achieving the smallest MAE faster.
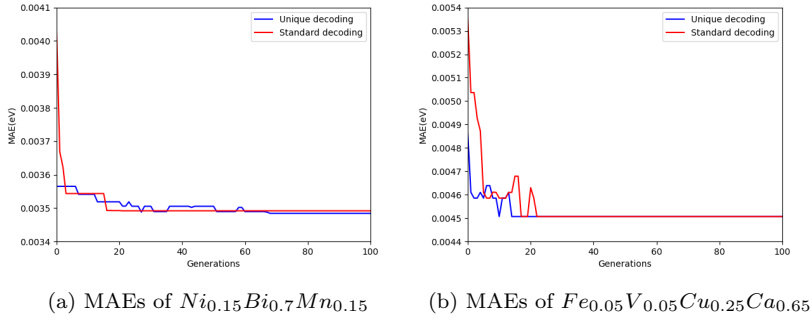


(a) MAEs of $Ni_{0.15}Bi_{0.7}Mn_{0.15}$      (b) MAEs of $Fe_{0.05}V_{0.05}Cu_{0.25}Ca_{0.65}$

Fig. 9: MAE distances of two GA decoding methods

3.3 Inverse design of two compounds using GA and BO

To evaluate the inverse design performance of our TLOpt algorithm, we select two metal oxide materials from the whole dataset as the design target spectra with known composition information. Since most of the materials in the dataset are ternary and quaternary compounds, we randomly selected one for ternary $Ni_{0.15}Bi_{0.7}Mn_{0.15}$ and one for quaternary compounds $Fe_{0.05}V_{0.05}Cu_{0.25}Ca_{0.65}$ respectively.

For each target spectra, we conduct two inverse design tasks: one with specified element set so the inverse design algorithm only needs to determine the mole ratios of the given elements; in the second design task, the elements are not given so the inverse design algorithm also needs to search the elements along with the mole ratios of all elements. For both experiments, we tested our GA and BO-based algorithms based on the same number of evaluations to compare their performance. The final prediction results are shown in **Fig.10** and **Fig.11**. **Fig.10** (a) shows the true spectrum and predicted spectrum by the GA search algorithm with composition elements [Ni, Bi, Mn] specified. After 100 generations, GA identified a formula $Ni_{0.25}Bi_{0.67}Mn_{0.08}$ with very similar spectrum to the target spectrum. The MAE error of the predicted spectrum of GA is only 0.00274eV (See Table 4). **Fig.10** (b) shows the spectrum of inverse designed metal oxide $Ni_{0.04}Bi_{0.71}Mn_{0.25}$, which has a slightly higher MAE error 0.00278eV. Both spectra of the inverse design materials closely approximate the target spectra, which demonstrates the effectiveness of the proposed approach. **Fig.10** (c), (d) compare the prediction spectra and the target spectra of the inverse design metal oxides by GA and BO when the composition elements are not specified. After 100 generations, GA got the most similar spectrum with formula $Mn_{0.22}Fe_{0.45}Pd_{0.33}$, and BO got $Gd_{0.06}Ni_{0.41}Fe_{0.53}$. The MAEs of GA and BO are 0.00273eV and 0.00275eV respectively.

**Fig.11** (a), (b) show the truly spectrum and prediction spectrum of formula $Fe_{0.05}V_{0.05}Cu_{0.25}Ca_{0.65}$ by GA and BO with given composition elements [Fe, V, Cu, Ca]. After 100 generations, GA got the most similar spectrum with formula $Fe_{0.08}V_{0.07}Cu_{0.25}Ca_{0.6}$, and BO got $Fe_{0.02}V_{0.02}Cu_{0.25}Ca_{0.71}$. The MAEs of GA and BO are 0.00280eV and 0.00284eV respectively. **Fig.11** (c),(d) compare the prediction spectra and truly spectra of target formula by GA and BO when the composition elements are not specified. After 100 generations, GA got the most similar spectrum with formula $Al_{0.15}Ti_{0.47}Ba_{0.11}Sc_{0.27}$, and BO got $Mn_{0.36}Yb_{0.34}Mg_{0.22}Pb_{0.08}$. The MAEs of GA and BO are 0.00278eV and 0.00283eV respectively. The results are also shown in Table 4.

From Table 4 we can find that, for the same formula, GA performs better than BO method. When the composition elements are not specified, the prediction spectrum is closer to the true spectrum. For ternary and quaternary compounds, with same generations, the ternary compounds have better prediction performance. In the case of the same number of evaluations, BO requires much less time than GA, the time required for BO is about one-quarter to one-third of GA, but it is easy to fall into local optimization. GA can find the global optimum due to its characteristics of crossover, mutation and elite reservation, but it costs much more time. As the search space becomes larger and more combinations can be made, both GA and BO methods perform better when using random elements. And ternary compounds $Ni_{0.15}Bi_{0.7}Mn_{0.15}$ perform better than quaternary $Fe_{0.05}V_{0.05}Cu_{0.25}Ca_{0.65}$ is determined by the characteristics of the dataset, the element Bi, Mn in [Ni, Bi, Mn] appear more frequently in the dataset, while the Fe, Cu, Ca in [Fe, V, Cu, Ca] appears relatively less frequently. The frequency of all elements in the dataset is shown as **Fig.12**.
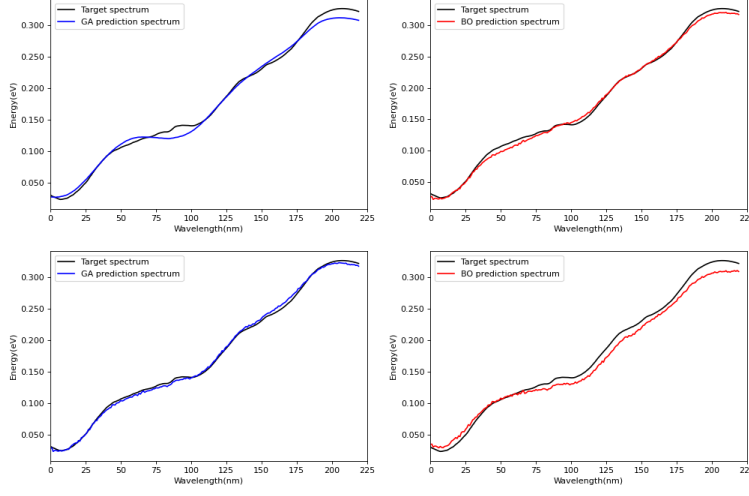
Fig. 10: Inverse designs for target spectrum of $Ni_{0.15}Bi_{0.7}Mn_{0.15}$, (a) Spectra comparison of inverse designs by GA with specified elements, (b) Spectra comparison of inverse designs by BO with specified elements, (c) Spectra comparison of the inverse designs by GA without specifying composition elements, (d) Spectra comparison of inverse designs by BO without specifying composition elements

Table 4: Predicted formulas and MAEs of experiment 1

| Target | $Ni_{0.15}Bi_{0.7}Mn_{0.15}$ | | $Fe_{0.05}V_{0.05}Cu_{0.25}Ca_{0.65}$ | |
|---|---|---|---|---|
| Algorithm | Predicted Formula | MAE(eV) | Predicted Formula | MAE(eV) |
| GA with given | $Ni_{0.25}Bi_{0.67}Mn_{0.08}$ | 0.00274 | $Fe_{0.08}V_{0.07}Cu_{0.25}Ca_{0.6}$ | 0.00280 |
| BO with given | $Ni_{0.04}Bi_{0.71}Mn_{0.25}$ | 0.00278 | $Fe_{0.02}V_{0.02}Cu_{0.25}Ca_{0.71}$ | 0.00284 |
| GA with random | $Mn_{0.22}Fe_{0.45}Pd_{0.33}$ | 0.00273 | $Al_{0.15}Ti_{0.47}Ba_{0.11}Sc_{0.27}$ | 0.00278 |
| BO with random | $Gd_{0.06}Ni_{0.41}Fe_{0.53}$ | 0.00275 | $Mn_{0.36}Yb_{0.34}Mg_{0.22}Pb_{0.08}$ | 0.00283 |

## 3.4 Prediction of the spectrum for a span of representative samples

To further evaluate the performance of our TLOpt inverse design algorithm, we randomly select 25 target compositions and their spectra as design targets from the whole 554 groups. The spectra (red color) of the inverse designed metal oxides predicted by the BO algorithm are shown in **Fig.13** together with the target spectra (blue color).

From **Fig.13**, we can find that most reconstructed spectra patterns (from row 1 to row 4) contain not only the general shape of the truth spectra but also finer details such as the presence of local maxima in absorption. For the
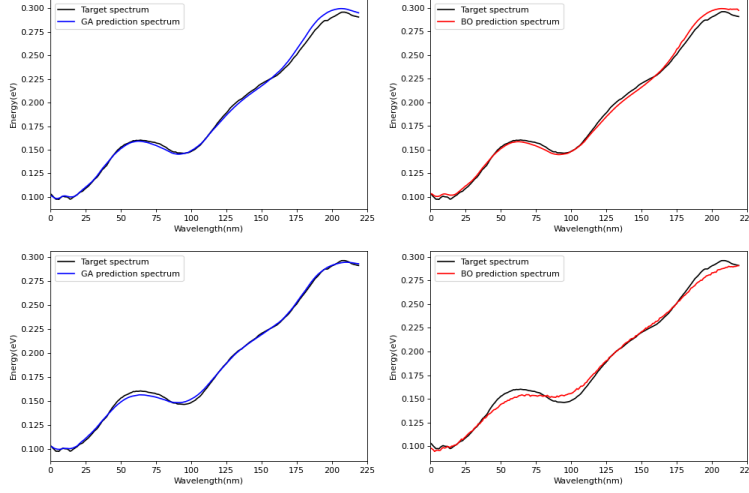
Fig. 11: Inverse designed spectra of $Fe_{0.05}V_{0.05}Cu_{0.25}Ca_{0.65}$, (a) Spectra comparison of inverse designs by GA with specified elements, (b) Spectra comparison of inverse designs by BO with specified elements, (c) Spectra comparison of the inverse designs by GA without specifying composition elements, (d) Spectra comparison of inverse designs by BO without specifying composition elements

targets in the last row in **Fig.13**, the performance of the inverse design is still great, though not as good as those in top 4 rows. After close examination, we find that this is due to the composition elements Lu, Gd, Pd of these target materials rarely appear in the whole dataset while our TLOpt algorithm achieves great performance for the targets in the top 3 rows, which contain common composition elements Bi, Mn, V, Cu, and Ni. The results in **Fig.13** thus validate our inverse design method on the dataset, proving the universality of the model. This inverse design model enables us to exploit hidden relationship between materials composition and their optical absorption properties and discover new materials with desired optical property.

## 4 Conclusion

We propose a transfer learning and global optimization based framework for inverse design of optical materials composition to achieve the target optical absorption spectrum. Our framework is composed of a fully connected neural network-based transfer learning model trained with Magpie features and global optimization based search including genetic algorithm and a Bayesian
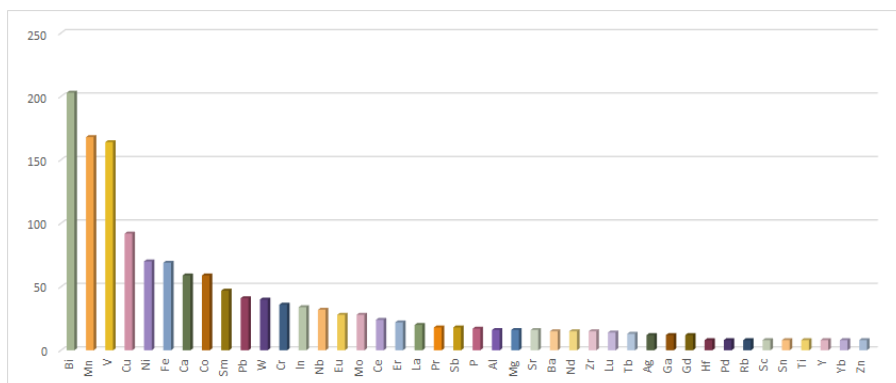
Fig. 12: Frequency of all elements in the dataset

optimization model. Our transfer learning algorithm can be used to address the small dataset problem typical in material informatics, enabling our DNN model for predicting the material's full UV-vis absorption spectrum from only its compound formula. Experiment of our transfer learning shows that after initial training and fine-tuning parameters through transfer learning, our prediction model performs well in spectrum prediction of metal oxide materials with only composition information alone. Extensive experiments show that our frame is able to discover interesting material compositions that approximate the target optical absorption spectrum. Our experiments also present that when running time is not an issue, genetic algorithm methods perform better than the Bayesian optimization in global optimization for our inverse design. Our research proves that machine learning based inverse design method could be used with small dataset when there are additional data of the same type used to initial training prediction model. Our inverse design model can be further improved when combined with materials structure information. Based on the successful design case studies, we believe our inverse design model is of great significance to be used to guide the discovery of new materials with other properties as well.

## 5 Contribution

Conceptualization, J.H.; methodology, J.H. and R.D.; software, R.D. and J.H.; validation, R.D. and J.H.; investigation, R.D., J.H., Y.D.; resources, J.H.; writing–original draft preparation, R.D. and J.H.; writing–review and editing, J.H; visualization, R.D.,Y.D., and X.L.; supervision, J.H.; funding acquisition, J.H.
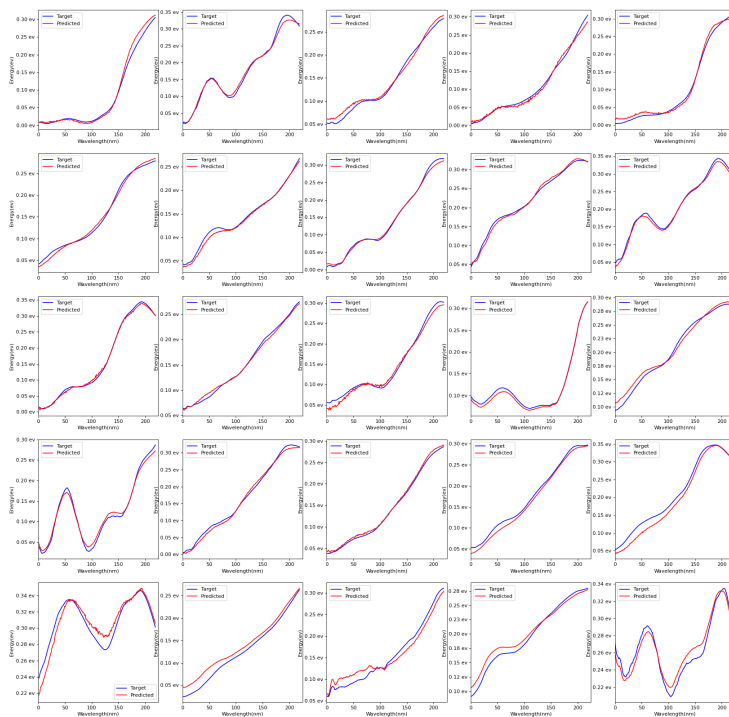
Fig. 13: The inverse design of 25 metal oxides given the target spectra

## 6 Acknowledgement

## References

1. Pornsawan Sikam, Pairot Moontragoon, Zoran Ikonic, Thanayut Kaewmaraya, and Prasit Thongbai. The study of structural, morphological and optical properties of (al, ga)-doped zno: Dft and experimental approaches. *Applied Surface Science*, 480:621–635, 2019.

2. Muhammad Khalid, Malik Aman Ullah, Muhammad Adeel, Muhammad Usman Khan, Muhammad Nawaz Tahir, and Ataualpa Albert Carmo Braga. Synthesis, crystal structure analysis, spectral ir, uv–vis, nmr assessments, electronic and nonlinear optical properties of potent quinoline based derivatives: interplay of experimental and dft study. *Journal of Saudi Chemical Society*, 23(5):546–560, 2019.

3. Krishna Rajan. Materials informatics. *Materials Today*, 8(10):38–45, 2005.

4. Logan Ward and Chris Wolverton. Atomistic calculations and materials informatics: A review. *Current Opinion in Solid State and Materials Science*, 21(3):167–176, 2017.

5. Rafael Gómez-Bombarelli, Jennifer N Wei, David Duvenaud, José Miguel Hernández-Lobato, Benjamín Sánchez-Lengeling, Dennis Sheberla, Jorge Aguilera-Iparraguirre, Timothy D Hirzel, Ryan P Adams, and Alán Aspuru-Guzik. Automatic chemical design using a data-driven continuous representation of molecules. *ACS central science*, 4(2):268–276, 2018.

6. Mariya Popova, Olexandr Isayev, and Alexander Tropsha. Deep reinforcement learning for de novo drug design. *Science advances*, 4(7):eaap7885, 2018.

7. Shuaihua Lu, Qionghua Zhou, Yixin Ouyang, Yilv Guo, Qiang Li, and Jinlan Wang. Accelerated discovery of stable lead-free hybrid organic-inorganic perovskites via machine learning. *Nature communications*, 9(1):1–8, 2018.

8. 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.

9. Alex Zunger. Inverse design in search of materials with target functionalities. *Nature Reviews Chemistry*, 2(4):1–16, 2018.

10. Benjamin Sanchez-Lengeling and Alán Aspuru-Guzik. Inverse molecular design using machine learning: Generative models for matter engineering. *Science*, 361(6400):360–365, 2018.

11. Yuichi Ikeda. A new method of alloy design using a genetic algorithm and molecular dynamics simulation and its application to nickel-based superalloys. *Materials transactions, JIM*, 38(9):771–779, 1997.

12. Sean Molesky, Zin Lin, Alexander Y Piggott, Weiliang Jin, Jelena Vuckovic, and Alejandro W Rodriguez. Inverse design in nanophotonics. *Nature Photonics*, 12(11):659–670, 2018.

13. Alexander Y Piggott, Jan Petykiewicz, Logan Su, and Jelena Vučcković. Fabrication-constrained nanophotonic inverse design. *Scientific reports*, 7(1):1–7, 2017.

14. Dianjing Liu, Yixuan Tan, Erfan Khoram, and Zongfu Yu. Training deep neural networks for the inverse design of nanophotonic structures. *ACS Photonics*, 5(4):1365–1369, 2018.

15. John Peurifoy, Yichen Shen, Li Jing, Yi Yang, Fidel Cano-Renteria, Brendan G DeLacy, John D Joannopoulos, Max Tegmark, and Marin Soljačić. Nanophotonic particle simulation and inverse design using artificial neural networks. *Science advances*, 4(6):eaar4206, 2018.

16. Jiaqi Jiang and Jonathan A Fan. Simulator-based training of generative neural networks for the inverse design of metasurfaces. *Nanophotonics*, 1(ahead-of-print), 2019.

17. Zhaocheng Liu, Dayu Zhu, Sean P Rodrigues, Kyu-Tae Lee, and Wenshan Cai. Generative model for the inverse design of metasurfaces. *Nano letters*, 18(10):6570–6576, 2018.

18. Raphaël Pestourie, Carlos Pérez-Arancibia, Zin Lin, Wonseok Shin, Federico Capasso, and Steven G Johnson. Inverse design of large-area metasurfaces. *Optics express*, 26(26):33732–33747, 2018.

19. Hillel Aharoni, Yu Xia, Xinyue Zhang, Randall D Kamien, and Shu Yang. Universal inverse design of surfaces with thin nematic elastomer sheets. *Proceedings of the National Academy of Sciences*, 115(28):7206–7211, 2018.

20. Jessica G Freeze, H Ray Kelly, and Victor S Batista. Search for catalysts by inverse design: artificial intelligence, mountain climbers, and alchemists. *Chemical reviews*, 119(11):6595–6612, 2019.

21. Benjamin Sanchez-Lengeling, Carlos Outeiral, Gabriel L Guimaraes, and Alán Aspuru-Guzik. Optimizing distributions over molecular space. an objective-reinforced generative adversarial network for inverse-design chemistry (organic). 2017.

22. H Döscher, JF Geisz, TG Deutsch, and JA Turner. Sunlight absorption in water–efficiency and design implications for photoelectrochemical devices. *Energy & Environmental Science*, 7(9):2951–2956, 2014.
23. T Warren Liao and Guoqiang Li. Metaheuristic-based inverse design of materials–a survey. *Journal of Materiomics*, 2020.
24. Longhui Qin, Weicheng Huang, Yayun Du, Luocheng Zheng, and Mohammad Khalid Jawed. Genetic algorithm-based inverse design of elastic gridshells. *Structural and Multidisciplinary Optimization*, pages 1–17, 2020.
25. Xiaoyang Li, Yuqing Hu, Enrico Zio, and Rui Kang. A bayesian optimal design for accelerated degradation testing based on the inverse gaussian process. *IEEE Access*, 5:5690–5701, 2017.
26. Mihir R Khadilkar, Sean Paradiso, Kris T Delaney, and Glenn H Fredrickson. Inverse design of bulk morphologies in multiblock polymers using particle swarm optimization. *Macromolecules*, 50(17):6702–6709, 2017.
27. Yue-Yu Zhang, Weiguo Gao, Shiyou Chen, Hongjun Xiang, and Xin-Gao Gong. Inverse design of materials by multi-objective differential evolution. *Computational Materials Science*, 98:51–55, 2015.
28. Sujin Bureerat and Nantiwat Pholdee. Inverse problem based differential evolution for efficient structural health monitoring of trusses. *Applied Soft Computing*, 66:462–472, 2018.
29. Ziku Wu, Chang Ding, Guofeng Li, Xiaoming Han, and Juan Li. Learning solutions to the source inverse problem of wave equations using ls-svm. *Journal of Inverse and Ill-posed Problems*, 27(5):657–669, 2019.
30. Sebastian J Wirkert, Hannes Kenngott, Benjamin Mayer, Patrick Mietkowski, Martin Wagner, Peter Sauer, Neil T Clancy, Daniel S Elson, and Lena Maier-Hein. Robust near real-time estimation of physiological parameters from megapixel multispectral images with inverse monte carlo and random forest regression. *International journal of computer assisted radiology and surgery*, 11(6):909–917, 2016.
31. Gang Sun, Yanjie Sun, and Shuyue Wang. Artificial neural network based inverse design: Airfoils and wings. *Aerospace Science and Technology*, 42:415–428, 2015.
32. Helge S Stein, Dan Guevarra, Paul F Newhouse, Edwin Soedarmadji, and John M Gregoire. Machine learning of optical properties of materials–predicting spectra from images and images from spectra. *Chemical science*, 10(1):47–55, 2019.
33. Helge S Stein, Edwin Soedarmadji, Paul F Newhouse, Dan Guevarra, and John M Gregoire. Synthesis, optical imaging, and absorption spectroscopy data for 179072 metal oxides. *Scientific data*, 6(1):1–5, 2019.
34. Liping Yu, Robert S Kokenyesi, Douglas A Keszler, and Alex Zunger. Inverse design of high absorption thin-film photovoltaic materials. *Advanced Energy Materials*, 3(1):43–48, 2013.
35. Yabo Dan, Yong Zhao, Xiang Li, Shaobo Li, Ming Hu, and Jianjun Hu. Generative adversarial networks (gan) based efficient sampling of chemical composition space for inverse design of inorganic materials. *npj Computational Materials*, 6(1):1–7, 2020.
36. Helge S. Stein, Guevarra Dan, Paul F. Newhouse, Edwin Soedarmadji, and John M. Gregoire. Machine learning of optical properties of materials – predicting spectra from images and images from spectra. *Chemical Ence*, 10, 2018.
37. Logan Ward, Ankit Agrawal, Alok Choudhary, and Christopher Wolverton. A general-purpose machine learning framework for predicting properties of inorganic materials. *npj Computational Materials*, 2(1):1–7, 2016.
38. Chi Chen, Yunxing Zuo, Weike Ye, Xiangguo Li, Zhi Deng, and Shyue Ping Ong. A critical review of machine learning of energy materials. *Advanced Energy Materials*, 10(8):1903242, 2020.
39. Zhuo Cao, Yabo Dan, Zheng Xiong, Chengcheng Niu, Xiang Li, Songrong Qian, and Jianjun Hu. Convolutional neural networks for crystal material property prediction using hybrid orbital-field matrix and magpie descriptors. *Crystals*, 9(4):191, 2019.
40. https://reference.wolfram.com/language/note/ElementDataSourceInformation.html.
41. Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 2020.

42. Ansi Zhang, Honglei Wang, Shaobo Li, Yuxin Cui, Zhonghao Liu, Guanci Yang, and Jianjun Hu. Transfer learning with deep recurrent neural networks for remaining useful life estimation. *Applied Sciences*, 8(12):2416, 2018.

43. JohnH Holland. adaptation in natural and artificial systems, university of michigan press, ann arbor,". *Cité page*, 100, 1975.

44. Alois Huning. Evolutionsstrategie. optimierung technischer systeme nach prinzipien der biologischen evolution, 1976.

45. Mandavilli Srinivas and Lalit M Patnaik. Genetic algorithms: A survey. *computer*, 27(6):17–26, 1994.

46. Maria Pakhnova, Ivan Kruglov, Alexey Yanilkin, and Artem R Oganov. Search for stable cocrystals of energetic materials using the evolutionary algorithm uspex. *Physical Chemistry Chemical Physics*, 2020.

47. Paul C Jennings, Steen Lysgaard, Jens Strabo Hummelshøj, Tejs Vegge, and Thomas Bligaard. Genetic algorithms for computational materials discovery accelerated by machine learning. *npj Computational Materials*, 5(1):1–6, 2019.

48. Mahdi Shariati, Mohammad Saeed Mafipour, Peyman Mehrabi, Masoud Ahmadi, Karzan Wakil, Nguyen Thoi Trung, and Ali Toghroli. Prediction of concrete strength in presence of furnace slag and fly ash using hybrid ann-ga (artificial neural network-genetic algorithm). *Smart Structures and Systems*, 25(2):183–195, 2020.

49. Oliver C. Gobin and Ferdi Schüth. On the suitability of different representations of solid catalysts for combinatorial library design by genetic algorithms. *Journal of Combinatorial Chemistry*, 10(6):835–846, 2008.

50. Harold J Kushner. A new method of locating the maximum point of an arbitrary multipeak curve in the presence of noise. 1964.

51. Donald R Jones, Matthias Schonlau, and William J Welch. Efficient global optimization of expensive black-box functions. *Journal of Global optimization*, 13(4):455–492, 1998.

52. Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical bayesian optimization of machine learning algorithms. In *Advances in neural information processing systems*, pages 2951–2959, 2012.

53. Ansi Zhang, Shaobo Li, Yuxin Cui, Wanli Yang, Rongzhi Dong, and Jianjun Hu. Limited data rolling bearing fault diagnosis with few-shot learning. *IEEE Access*, 7:110895–110904, 2019.