

Interpretable Machine Learning for Materials Design

James Dean¹, Matthias Scheffler^{2,3}, Thomas A. R. Purcell³, Sergey V. Barabash⁴, Rahul Bhowmik⁵, and Timur Bazhirov^{1,*}

¹Exabyte Inc, San Francisco, CA, United States

²University of California Santa Barbara, Isla Vista, CA, United States

³The NOMAD Laboratory at the Fritz Haber Institute, Berlin, Germany

⁴Intermolecular Inc, San Jose, CA, United States

⁵Polaron Analytics, Beavercreek, OH, United States

*Corresponding Author Email: timur@exabyte.io

December 2, 2021

Abstract

Fueled by the widespread adoption of Machine Learning and the high-throughput screening of materials, the data-centric approach to materials design has asserted itself as a robust and powerful tool for the *in-silico* prediction of materials properties. When training models to predict material properties, researchers often face a difficult choice between a model's interpretability or its performance. We study this trade-off by leveraging four different state-of-the-art Machine Learning techniques: XGBoost, SISSO, Roost, and TPOT for the prediction of structural and electronic properties of perovskites and 2D materials. We then assess the future outlook of the continued integration of Machine Learning into materials discovery, and identify key problems that will continue to challenge researchers as the size of the literature's datasets and complexity of models increases. Finally, we offer several possible solutions to these challenges with a focus on retaining interpretability, and share our thoughts on magnifying the impact of Machine Learning on materials design.

Keywords: machine learning, materials science, chemistry, interpretability, rational design.

1 Introduction

Today, big data and artificial intelligence revolutionize many areas of our daily life, and materials science is no exception^[1-3]. More scientific data is available now than ever before and the size of the literature is growing at an exponential rate^[4-7]. This has led to multiple efforts in building the digital ecosystem for material discovery, most notably the Materials Genome Initiative (MGI)^[8,9]. The MGI is a multinational effort focused on improving the tools and techniques surrounding materials research, which recently has included suggestions to adopt the set of Findable, Accessible, Interoperable, and Reusable (FAIR) principles when reporting data^[10]. In the years since the creation of the MGI, a number of large materials and chemical datasets have emerged, including the 2D Materials Encyclopedia (2DMatPedia)^[11], Automatic Flow (AFLOW) database^[12,13], Computational 2D Materials Database (C2DB)^[14,15], Computational Materials Repository (CMR)^[16], Joint Automated Repository for Various Integrated Simulations (JARVIS)^[17], Materials Project^[18], Novel Materials Discovery (NOMAD) repository^[19], and the Open Quantum Materials Database (OQMD)^[20]. We note that all of these are primarily computational in nature, and that there is still a scarcity of large databases containing comprehensively-characterized experimental data. Despite this, at least in computational materials discovery, the current availability of data has been a boon for exploration of the materials space, as it allows for highly flexible, data-hungry^[21] models to be trained.

One such approach that has seen widespread popularity in recent years is gradient boosting. Gradient boosting^[22] is an ensemble technique in which a collection of weak learners (typically decision trees) are incrementally trained with respect to the gradient of the loss function^[23]. A well-known implementation

(with over 5,500 citations as of November 2021) is eXtreme Gradient Boosting (XGBoost)^[24], which reformulates the algorithm to provide stronger regularization and improved protection against over-fitting. In chemistry, its applications have been diverse: XGBoost has been used to predict the adsorption energy of noble gases to Metal-Organic Frameworks (MOFs)^[25], biological activity of pharmaceuticals^[26], atmospheric transport^[27], and has even been combined with the representations found in Graph Neural Networks (GNNs) to generate accurate models of various molecular properties^[28].

Neural networks have also seen a lot of interest, owing to their greater flexibility relative to other model types. This has included the influential Behler-Parinello^[29] and Crystal Graph Convolutional Neural Network (CGCNN)^[30] architectures based on chemical structure, the Representation Learning from Stoichiometry (Roost)^[31] architecture based on chemical formula, and many other approaches^[1,2,32-41].

The modern Machine Learning (ML) toolbox is large, although it is still far from complete. As a result model selection techniques are becoming increasingly necessary: this has led to the field of Automated Machine Learning (AutoML). This area of work has seen much progress in recent years^[42,43], and has even been extended to Neural Architecture Search (NAS)^[44], or the automated optimization of neural network architectures. In this work, we leverage the Tree-based Pipeline Optimization Tool (TPOT) approach to AutoML^[45-47], which uses a Genetic Algorithm (GA) to create effective ML pipelines. Although it generally draws from the models of SciKit-Learn^[48], it can also be configured to explore gradient boosting models via XGBoost^[22], and neural network models via PyTorch^[49]. Moreover, TPOT also performs its own hyperparameter optimization, thus providing a more hands-off solution to identifying ML pipelines. The success of GA-based approaches in ML is not isolated to AutoML. Indeed, they are a fundamental part of genetic programming, where they are used to optimize functions for a particular task^[50,51]. Eureka^[52] is a particularly successful example of this^[53], leveraging a GA to generate equations fitting arbitrary functions, and has been used in several areas of chemistry, including the generation of adsorption models to nanoparticles^[54] and metal atoms to oxide surfaces^[55]. This approach of fitting arbitrary functions to a task is also known as "symbolic regression." Recent work surrounding compressed sensing has yielded the Sure Independence Screening and Sparsifying Operator (SISSO) approach^[56]. SISSO also generates equations mapping descriptors to a target property, proceeding by combining descriptors using various building blocks, including trigonometric functions, logarithms, addition, multiplication, exponentiation, and many others. This methodology has been highly successful in a variety of areas including crystal structure classification^[57], as well as the prediction of perovskite properties^[58-60] and 2D topological insulators^[61].

The abundance of scientific data in the literature increasingly makes the use of highly flexible (yet data-hungry) techniques tractable. Although such techniques may deliver models which are highly accurate and generalize well to new data, what is oftentimes lost is the physical interpretation of these models. By physical interpretation, we mean an understanding of the relationship between the chosen descriptors and the target property. Although a black-box model which has a high level of accuracy but little physical interpretation may lend itself well to the Edisonian screening of a wide range of materials, it may be difficult to understand exactly what feature (or combination of features) actually matters to the design of the material. Once the screening is done and the target values are calculated, little may be done to improve performance aside from including new features, adjusting the model's hyperparameters, or increasing the size of the training set. Alternatively, consider a model which has less accuracy, but which has an intuitive explanation, such as an equation describing an approximate relationship between features and target. Although such a model may at first glance seem less useful than a highly-accurate black-box, such a model can help deliver insight into the underlying process that results in the target property. Moreover, by understanding which features are important, the model can give clues into what may be done to further improve it — driving the rational discovery of materials. In addition, interpretability versus accuracy is not a strict trade-off, and it is possible for interpretable and black-box models to deliver similar accuracy^[62]. Therefore, in this work we take steps to compare the performance of all four models for each problem with respect to i) performance and ii) interpretability.

We leverage a diverse selection of techniques in order to draw comparisons of model accuracy and interpretability. Taking advantage of the current abundance of chemical data, we can re-use the Density-Functional Theory (DFT) calculations of others stored on several FAIR chemical datasets. A series of three problems are investigated: 1) the prediction of perovskite volumes, 2) the prediction of 2D material bandgaps, and 3) the prediction of 2D material exfoliation energies. For the perovskite volume problem, we leverage the ABX₃ perovskite dataset (containing 144 examples) published by Körbel, Marques, and Botti^[63]; this dataset is hosted by NOMAD^[19], whose repository has strong focus on enabling researchers to report their data such that it satisfies the FAIR data principles^[64]. For the 2D material problems, we apply the 2DMatPedia published by Zhou et al^[11]; this is a large dataset of 6,351 hypothetical 2D

materials identified via a high-throughput screening of systems on the Materials Project^[18]. We are aware of other 2D material datasets, such as the C2DB^[14] and JARVIS^[17], but for the purposes of simplifying our study we only choose one.

We find that TPOT delivers high-quality models, generally outperforming the other methods in terms of fitness metrics. Despite this, interpretability is not guaranteed, as it can create highly complex pipelines. XGBoost lends itself to interpretation more consistently, as it allows for an importance metric, although it may be harder to understand exactly what the relationship is between the different features (or combinations of different features) and the target variable. We found that Roost performed well on problems that could be approached via compositional descriptors (i.e. without structural descriptors); as a result, it can help us understand when a target property requires more than just the composition. Finally, we achieve the easiest interpretability from SISSO, as it provides access to descriptors which directly capture the relationship between the features and target variable. Using these results, we discuss the advantages and disadvantages of each method, and discuss areas where the digital ecosystem surrounding materials discovery could be improved to improve adherence to FAIR principles. Our work provides a comparison of several common ML techniques on challenging (but relevant) materials property prediction problems.

The manuscript is organized as follows: we begin by training a diverse set of four models, which are XGBoost, TPOT, Roost, and SISSO to investigate each of the three problems, resulting in a total of 24 trained models. Performance metrics (and comparative plots) are presented for each trained model to facilitate comparison, and we discuss the interpretation we can achieve from each of these models. Finally, we provide a discussion of the future outlook of ML in the digital materials science ecosystem and what can be done to further accelerate materials discovery.

2 Methodology

2.1 Data Sources

Crystal structures for the perovskite systems were obtained from the “Stable Inorganic Perovskites” dataset published by Körbel, Marques, and Botti^[63], as hosted by NOMAD^[19]. This dataset contains a total of 144 DFT-relaxed inorganic perovskites identified via a high-throughput screening strategy. Using this dataset, we develop a model of perovskite volume. As we rely on the use of compositional descriptors for these systems, we have scaled the volume of the perovskite unit cell by the number of formula units, such that the volume has units of \AA^3 / formula unit.

Structures for 2D materials were obtained from the 2DMatPedia^[11], a large database containing a mixture of 6,351 real and hypothetical 2D systems. This database was generated via a DFT-based high-throughput screening approach, which investigated bulk structures hosted by the Materials Project database^[18] to find systems which may plausibly form 2D structures. Among other things, the 2DMatPedia provides DFT-calculated exfoliation energies and bandgaps, along with a DFT-optimized structure for each material. We use this dataset to develop models for the bandgap and exfoliation energy of 2D materials. Although the dataset reports exfoliation energies in units of eV, we have converted these into units of J / m². Bandgaps are reported in units of eV.

2.2 Feature Engineering

To facilitate the development of ML algorithms capable of rapidly predicting material properties, we focus primarily on features that do not require further (computationally-intensive) DFT calculations. In the case of the 2D material bandgap, we include the DFT-calculated bandgap of the respective bulk material; we note that these values are tabulated on the Materials Project and can be looked up, thus circumventing the need for further DFT work. A variety of chemical featurization libraries were used to generate compositional and structural descriptors for the systems we investigated, and are listed in sections 2.2.1 and 2.2.2 respectively. Features with values of NaN (which occurred when a feature could not be calculated) were assigned a value of 0.

2.2.1 Compositional Descriptors

Compositional (i.e. chemical formula-based) descriptors were calculated via the open-source XenonPy packaged developed by Yamada et al^[65]. XenonPy uses tabulated elemental data from Mendeleev^[66], Pymatgen^[67], the CRC Handbook of Chemistry and Physics^[68], and Magpie^[69] in order to calculate

compositional features. XenonPy does this by combining the elemental descriptors (e.g. atomic weight, ionization potential, etc.) in various ways to form a single composition-weighted value. For example, three compositional descriptors may be obtained with XenonPy by taking the composition-weighted average, sum, or maximum elemental value of the atomic weight. Leveraging the full list of compositional features implemented in XenonPy results to 290 compositional descriptors, which are explained in greater detail within their publication^[65].

The compositional descriptors were used for the perovskite volume prediction, 2D material bandgap, and 2D material exfoliation energy prediction problems.

2.2.2 Structural Descriptors

Some structural descriptors were calculated using MatMiner^[70], an open-source Python package geared towards data-mining material properties. Leveraging MatMiner, the following descriptors were calculated: Average bond length, average bond angle, Global Instability Index (GII)^[71], Ewald Summation Energy^[72], a Shannon Information Entropy-based Structural Complexity (both per atom and per cell), and the number of symmetry operations available to the system. In the case of the average bond length and average bond angle, bonds were determined using Pymatgen’s implementation of the JMol^[73] AutoBond algorithm. This list of bonds was also used to calculate an average Coordination Number (CN) over all atoms in the unit cell. Finally, we also took the perimeter:area ratio of the 2D material’s repeating unit.

The structural descriptors were used for the 2D material bandgap and 2D material exfoliation energy problems.

2.3 Data Filtering

The choice of data filtering methodology was chosen based on the problem at-hand. The perovskite volume prediction problem did not utilize any data filtering. In the case of the 2D material bandgap and exfoliation energy prediction problems, the data obtained from the 2D MatPedia were required to satisfy all of the following criteria:

1. No elements from the f-block, larger than U, or noble gases were allowed.
2. Decomposition energy must be below 0.5 eV/atom.
3. Exfoliation energy must be strictly positive.

Additionally, in the case of the 2D material bandgap, data were required to have a parent material defined on the Materials Project. This was done because we use the Materials Project’s tabulated DFT bandgap of the bulk system as a descriptor for the bandgap of the corresponding 2D system.

2.4 ML Models

When training models, 10% of randomly selected data was held out as a testing set. The same train/test split was used for all 4 models considered (XGBoost, TPOT, Roost, and SISSO). To facilitate a transparent comparison between models, in all cases we report the Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), Maximum Error, and R^2 score of the test set.

2.4.1 Gradient-Boosting with XGBoost

When training XGBoost models, 20% of randomly selected data was held out as an internal validation set. This was used to adopt an early-stopping strategy, where if the model RMSE did not improve after 50 consecutive rounds, training was halted early. When training, XGBoost was configured to optimize its RMSE.

Hyperparameters were optimized via the open-source Optuna^[74] framework. The hyperparameter space was sampled using the Tree-structured Parzen Estimator (TPE) approach^[75,76]. To accelerate the hyperparameter search, we leveraged the Hyperband^[77] approach for model pruning, using the validation set RMSE to determine whether to prune a model. Hyperband’s budget for the number of trees in the ensemble was set to range between 1 and 256 (corresponding with the maximum number of estimators we allowed an XGBoost model to have). The search space for hyperparameters is found in Table 1.

Table 1: Ranges of hyperparameters screened with Optuna for all XGBoost runs. The search was inclusive of the listed minima and maxima. Hyperparameters use the same variable naming convention as in the XGBoost documentation.

Hyperparameter	Minimum	Maximum
<code>learning_rate</code>	0	2
<code>min_split_loss</code>	0	2
<code>max_depth</code>	0	256
<code>min_child_weight</code>	0	10
<code>reg_lambda</code>	0	2
<code>reg_alpha</code>	0	2

The variable names here (e.g. `learning_rate`) correspond with the variable names listed in the documentation of XGBoost. Additionally, Optuna was used to select a standardization strategy, choosing between Z-score normalization (i.e. subtracting the mean and dividing by the standard deviation) or Min/Max scaling (i.e. scaling the data such that it has minimum 0 and maximum 1). To prevent test-set leakage, the chosen standardizer was fit only with the internal training set, i.e. the portion of the training set that was not held out as an internal validation set. Optuna performed 1000 trials to minimize the validation set RMSE. We report the results of the final optimized model.

2.4.2 AutoML with TPOT

The AutoML tool TPOT was leveraged with a population size of 100 pipelines, with training proceeding for a total of 10 generations. A maximum evaluation time of 10 seconds per model was set. TPOT pipelines were optimized using the default regression configuration. As TPOT is an actively-maintained open-source repository, for the purposes of future replication we enumerate this configuration’s set of allowable components in Table S1. The models listed in this table could be combined in any order any number of times. Models were selected such that their 10-fold cross-validated RMSE was optimized. TPOT also conducts its own internal optimization of model hyperparameters, thus we did not perform our own hyperparameter optimization of the TPOT pipelines.

2.4.3 Neural Networks with Roost

The Roost Neural Network (NN) architecture was leveraged using the “example.py” script provided with its source code. Models are trained for a total of 512 epochs with the default settings. In the case of Roost models, the only feature provided is the composition of the system, given through the chemical formula.

2.4.4 Symbolic Regression with SISSO

The first step of using SISSO is reducing the number of primary features down from a list of hundreds down to the tens. This is done due to the exponential computational cost of SISSO with respect to the number of features and the number of rungs being considered. To perform this down selection we first generate a rung 1 feature space including all of the primary features and operators that are used in the SISSO calculation. We then check how often each of the primary features appear in the ten thousand generated features that are most correlated to the target property. Additionally, we add units to all of the pre-selected primary features to ensure all generated expressions are valid.

In many cases, it was easy to infer what the abstract units are for the XenonPy descriptors. In a few cases where the units weren’t as clear, we compared the reported elemental values of those units to those of known sources (e.g. the NIST WebBook^[78] or the CRC Handbook^[68]) in order to determine the units. Finally, although it was generally easy to determine where the source of a feature was, sometimes we were unable to determine a source. In these cases, we refer to the features as a “XenonPy” feature (for example, “ $r_{XenonPy}$ ”).

The optimal number of terms (up to 3) and rung (up to 2), i.e. the number of times operators are recursively applied to the feature space, are determined using a five-fold cross validation scheme. For all models we allow for an external bias term to be non-zero, and use a SIS selection size of 500. The resulting descriptors were then evaluated using the same external test set for each of the other methods. To take advantage of SISSO’s ability to generate new composite descriptors and operate in large feature spaces, additional features were included in the SISSO calculations. A full list of features used in the SISSO work can be found in the linked GitHub repository.

3 Results

3.1 Perovskite Volume Prediction

XGBoost, TPOT, and SISSO were applied to investigate the volume of perovskites as a function of the compositional features described in section 2.2.1. Additionally, we trained a Roost model on the chemical formula of the perovskites to predict the volume. The train/test split resulted in a total of 129 entries in the training set, and 15 in the test set. We find generally good performance on the perovskite volume problem across all 4 models, although the TPOT and SISSO model display the best performance by all metrics investigated (see Table 2), including respective test-set R^2 of 0.979 and 0.990. The Roost model also performs well with a test-set R^2 of 0.935, but it also has a non-normal error, as can be seen in Figure 1. Finally, we find that while XGBoost is the worst performing method, it still has a relatively good test-set R^2 of 0.866.

Table 2: Performance metrics for the XGBoost, TPOT, Roost, and SISSO models on the perovskite volume prediction problem. The parity plots for these models are depicted in 1.

Error Metric	Partition	XGBoost	TPOT	Roost	SISSO
MAE	Train	8.15	0.860	9.11	7.27
RMSE	Train	11.88	1.14	11.21	9.08
Max Error	Train	43.29	4.96	22.94	26.74
R^2	Train	0.950	1.000	0.955	0.971
MAE	Test	12.89	4.02	8.83	4.05
RMSE	Test	17.17	6.78	12.00	4.71
Max Error	Test	37.10	21.75	31.69	10.27
R^2	Test	0.866	0.979	0.935	0.990

The performance of all 4 models is summarized in Figure 1. Visually, we find a very tight fit by the TPOT model in both the training and test sets, with good correlation from the XGBoost and SISSO models. We also find a systematic under-prediction of perovskite volumes in the roost model in both the training and test set, with the under-prediction beginning at approximately 75 \AA^3 / formula unit, achieving a maximum deviation at approximately 130 \AA^3 / formula unit, and returning to parity at approximately 200 \AA^3 / formula unit.

The good performance of the TPOT model results from a generated pipeline with four stages. The first two stages are `MinMaxScaler` units; although the second one is redundant, as the data is already scaled to be between 0 and 1 by the first scaler. The third stage is a `OneHotEncoder` unit, which leverages one-hot encoding for categorical features (the TPOT implementation defines a categorical feature as one with fewer than 10 unique values). Finally, the perovskite volume is predicted using support vector regression.

In the case of the XGBoost model, we can extract feature importances. Although various different feature importance metrics can be derived from XGBoost, in this case we use the “gain” metric, which describes how the model’s loss function improves when a feature is chosen for a split while constructing the trees. A large number of features (290) were input into this model, so we display only the 10 most-important features identified by XGBoost in Supporting Information Figure S1. Here, we find that the average Rahm atomic radii^[79,80] have the highest importance score, followed by the average Van der Waals radius used by the Universal Force Field (UFF)^[81]. The remaining 288 features fall off as a long tail of low importance scores, indicating that they did little to improve the model’s performance in predicting the perovskite volume.

For SISSO, we used reduced the feature space as outlined in Section 2.4.4, with the pre-screened features listed in the Supporting Information along with the assumption we made about the units of the descriptor when fed into SISSO. Generally, we find that the main descriptors selected by the procedure are related to volume and atomic radius. Some other descriptors with less interpretability are found, such as the C6 dispersion coefficients, polarizability, melting points, and Herfindahl-Hirschman Index (HHI)^[82] production and reserve values. Although typically used to help indicate the size of a company within a particular sector of the economy, the XenonPy definition of HHI appears to come from the work of Gaultois et al^[82]. In the referenced work, the HHI production value refers to the geographic distribution of elemental production (e.g. answering the question of concentrated the industry that produces pure Fluorine is), and HHI reserve value describes the geographic distribution of known deposits of these materials (e.g. whether they are spread out over a wide area, or concentrated in a small area).

We report the best descriptor found in Equation 1. In this equation, the variables c_0, a_0, a_1 are the regression coefficients determined by SISSO.

$$V_{Perovskite} \approx c_0 + a_0 \cdot \frac{Z^{ave}}{C^{ave} \cdot (r_{Slater}^{ave} - r_{pyyikko,triple}^{ave})} + a_1 \cdot (V_{gs}^{ave} - V_{gs}^{min}) \cdot \frac{r_{pyyikko,triple}^{ave}}{r_{pyyikko}^{ave}} \quad (1)$$

where $c_0 = -10.547$, $a_0 = 4.556$, $a_1 = 3.050$, Z^{ave} is the average atomic number, C^{ave} is the average mass specific heat capacity of the elemental solid, r_{Slater}^{ave} is the average atomic covalent radius predicted by Slater, $r_{pyyikko,triple}^{ave}$ is the average triple bond covalent radius predicted by Pyyko, $r_{pyyikko}^{ave}$ is the average single bond covalent radius predicted by Pyyko, and V_{gs}^{ave} and V_{gs}^{min} are the average and minimum ground state volume per atom as calculated by DFT. Unsurprisingly the ground state atomic volumes and covalent radii play an important role in determining the final volume of the perovskite structures. Interestingly, both the atomic number and specific heat capacity of the material appear in the final descriptor; however, these likely only act as minor corrections to the final results.

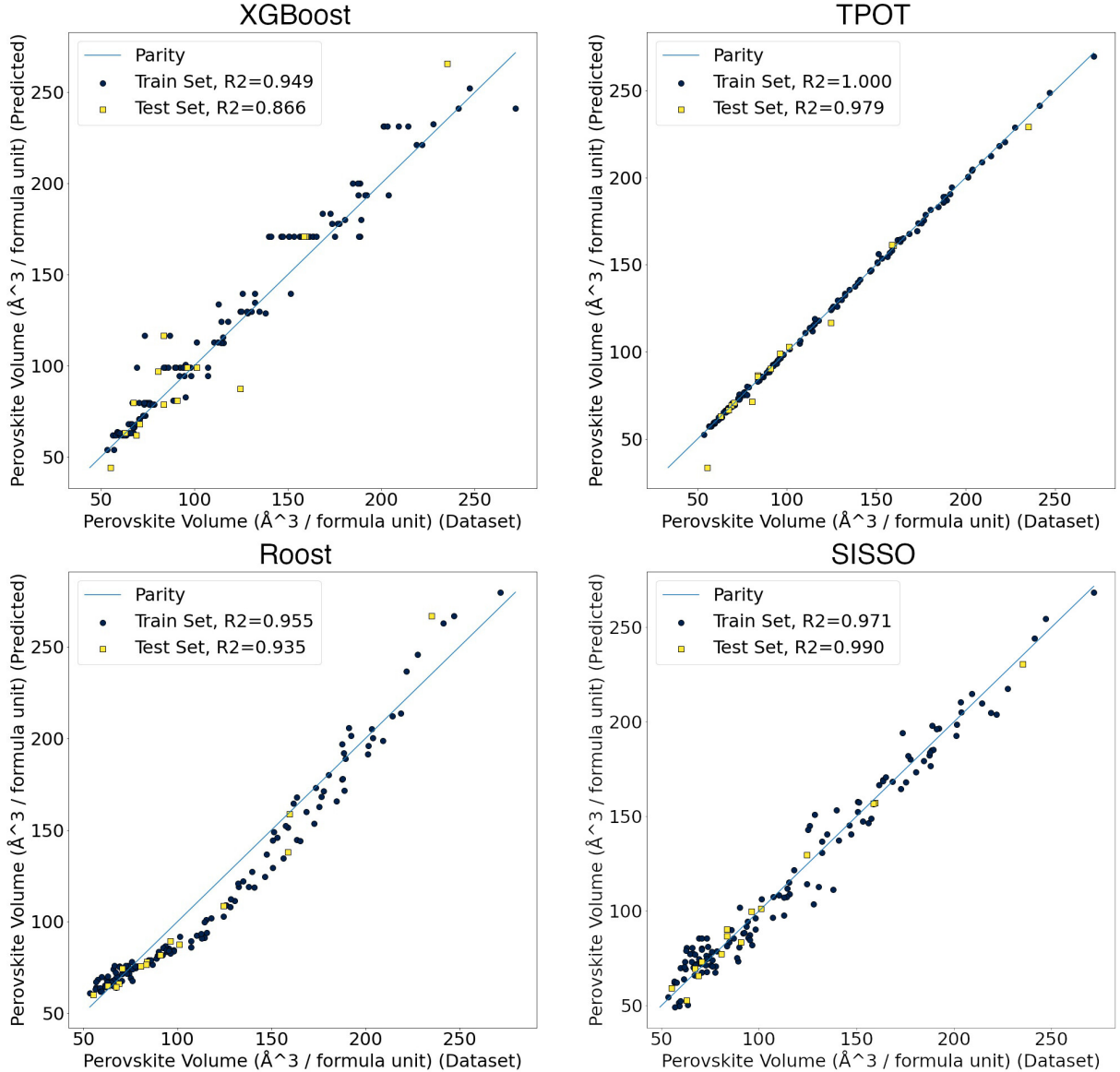


Figure 1: Parity plots for the XGBoost, TPOT, Roost, and SISSO models on the perovskite unit cell volume problem. Included are the training and testing sets. A diagonal line indicating parity is drawn as a guide to the eye.

3.2 2D Material Bandgaps

The bandgap predictions leveraged a data filtering strategy (described in Section 2.3). As a result of our data filtering approach, the 6,351 entries in the dataset were reduced to 1,412 entries. The train/test

split divided the data into a training set of 1,270 rows, and a test set containing 142 entries. The performance metrics of the XGBoost, TPOT, Roost, and SISSO models of 2D Material Bandgap can be found in Table 3. Performance is generally worse on this problem when compared to the perovskite volume predictions. As a result, in addition to the compositional features of XenonPy (Section 2.2.1) we also used several structural features (section 2.2.2). We also leveraged the bulk bandgap of the parent-3D material for each of the 2D materials, as we observed the performance of the TPOT, SISSO and XGBoost models increased when this value was included.

Table 3: Performance metrics for the XGBoost, TPOT, Roost, and SISSO models on the 2D material bandgap problem. The rung-2, 4-term SISSO model is reported. The parity plots for these models are depicted in 2.

Error Metric	Partition	XGBoost	TPOT	Roost	SISSO
MAE	Train	0.12	0.27	0.11	0.29
RMSE	Train	0.25	0.41	0.27	0.47
Max Error	Train	2.75	3.66	2.83	4.35
R ²	Train	0.973	0.928	0.968	0.907
MAE	Test	0.31	0.29	0.65	0.29
RMSE	Test	0.50	0.48	1.07	0.47
Max Error	Test	2.04	2.75	4.80	3.27
R ²	Test	0.892	0.900	0.507	0.919

Although test-set model performance was worse compared to the perovskite problem, the TPOT and SISSO models retained their status as having the best performance metrics for the test-set R², MAE, and RMSE. The XGBoost model additionally performed well, with the best performance in terms of maximum error, and nearly as good performance on the other metrics when compared to TPOT and SISSO. We find the Roost model overfit the data to some extent on the data, as the test-set error metrics are considerably worse than their training-set counterparts.

A parity plot summarizing these results can be found in Figure 2. In all cases, we can see a spike of misprediction for systems with a DFT bandgap of 0. We note here that a large portion of these entries had DFT bandgaps of 0: of the 382 of the 1,412 entries in the dataset, a total of 27% of all training data.

The pipeline generated by TPOT is less complex than that of the perovskite volume problem. The first stage of the pipeline is a `MinMaxScaler` unit, scaling each feature such that it is between 0 and 1. The second stage is then an `ElasticNetCV` unit, which uses 5-fold cross-validation to optimize the alpha and L1/L2 ratio of the Elastic Net model. The converged alpha value was 1.011×10^{-3} , and the converged L1/L2 ratio was 0.95, which strongly leans towards the L1 (Least Absolute Shrinkage and Selection Operator (LASSO)) regularization penalty.

We can also extract feature importances from the XGBoost model, and we report the 10 highest-ranked features in Supporting Information Figure S2. In contrast with the perovskite results, the features are generally ranked similarly; although there is a clear ranking, we do not see 1-2 features dominating followed by a long tail of unimportant features. The selected features are also less interpretable; we find the minimum Mendeleev number as the most important feature, an alternative numbering of the periodic table in which numbers are ordered by their group rather than period. Specifically, we use the variant proposed by Villars et al^[83]. For example, in the referenced system of numbering, Li, Na, K, and Rb are elements 1, 2, 3, and 4. This numbering is generally consistent, with the exception is H, which is assigned a value of 92 to place it above F, which has value 93. By taking the minimum value of the composition, this metric can be intuited as identifying the element in the earliest group and earliest period of the periodic table. Other descriptors we find are the average atomic weight, the variance of the number of unfilled s-orbitals in the system, and DFT-calculated ground-state atomic volumes.

The results of the prescreening procedure for SISSO are presented in the Supporting Information Table S3. These features are similar to the previous results, with the addition of the parent-3D material bandgap and electronegativity information playing an important role as well.

The selected SISSO model is

$$E_{Bandgap}^{2D} \approx c_0 + a_0 \cdot \frac{r_{vdw}^{min} \cdot E_{Bandgap}^{3D,parent}}{\sqrt[3]{Period}^{min}} \quad (2)$$

where $c_0 = 0.143$, $a_0 = 8.054 \times 10^{-3}$, r_{vdw}^{min} is the minimum Van der Waals radius of the atoms in the material, $E_{Bandgap}^{3D,parent}$ is the bandgap of the 3D-parent material, and $Period^{min}$ is the minimum period of the elements in the material. This descriptor represents a simple rescaling of the bandgap of the 3D-parent material, which makes sense as both bandgaps are highly correlated to each other. Furthermore

this results is consistent with the TPOT model which is also primarily controlled by the bandgap of the 3D-parent material.

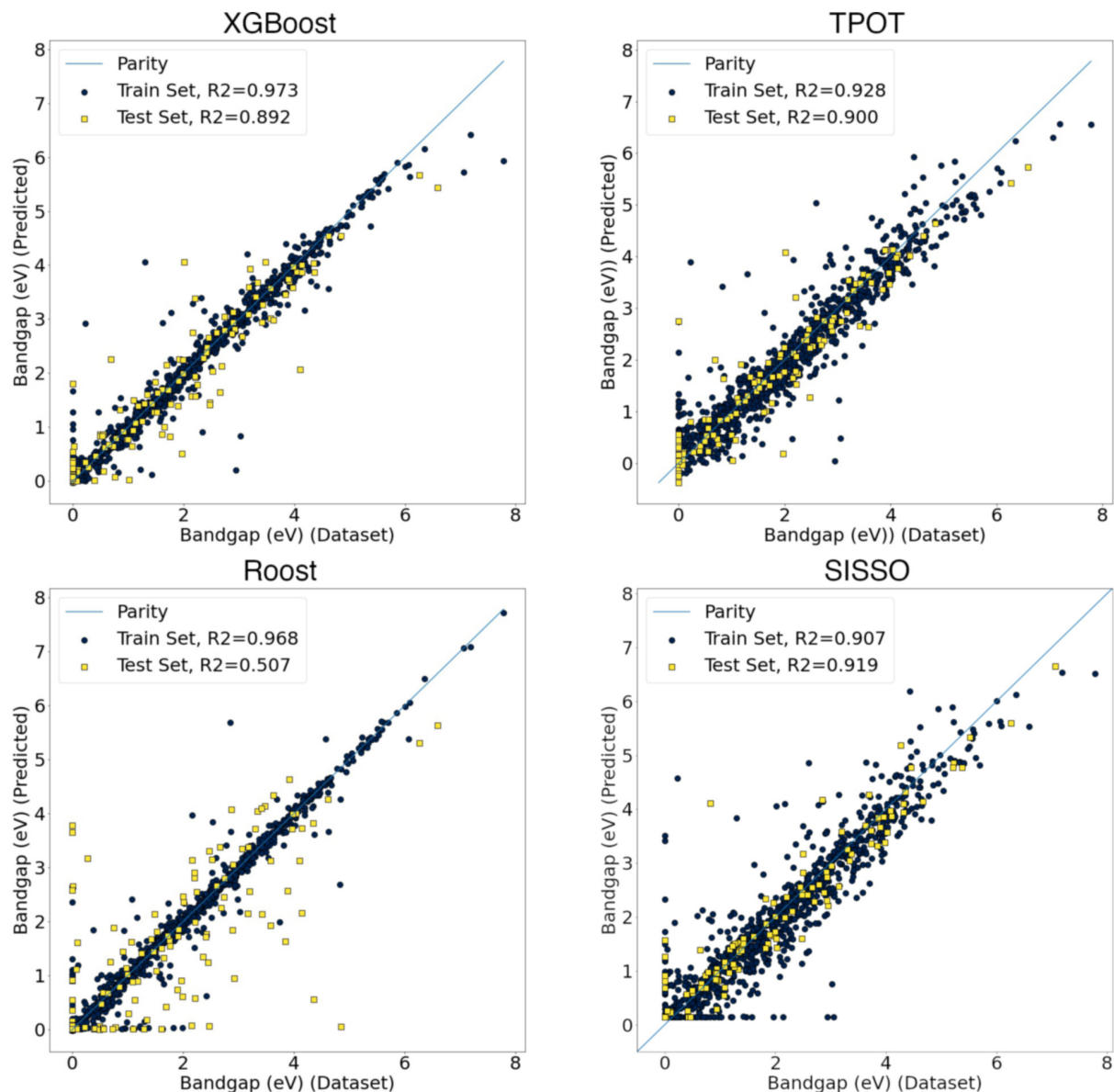


Figure 2: Parity plots for the XGBoost, TPOT, Roost, and SISSO models on the 2D material bandgap problem. Included are the training and testing sets. A diagonal line indicating parity is drawn as a guide to the eye. Regression statistics for the models shown on this plot can be found in Table 3.

3.3 2D Material Exfoliation Energy

In the case of the 2D material exfoliation energy problem, the training and test-set statistics for the XGBoost, TPOT, Roost, and SISSO models can be found in Table 4. In this case, our feature selection methodology down-selected the 6,351 rows of our dataset into 3,388 rows. The train/test split further divided this into a training-set of 3,049 entries, and a test set of 339 entries. Generally, we see the worst performance of the models in this problem, compared to the perovskite volume and 2D material bandgap problems.

Table 4: Performance metrics for the XGBoost, TPOT, Roost, and SISSO models on the 2D material exfoliation energy problem. The parity plots for these models are depicted in 3.

Error Metric	Partition	XGBoost	TPOT	Roost	SISSO
MAE	Train	0.20	0.14	0.06	0.26
RMSE	Train	0.35	0.31	0.24	0.45
Max Error	Train	7.11	8.34	9.63	9.58
R ²	Train	0.624	0.702	0.827	0.365
MAE	Test	0.23	0.21	0.19	0.27
RMSE	Test	0.35	0.33	0.34	0.42
Max Error	Test	1.64	1.85	1.96	2.48
R ²	Test	0.476	0.543	0.499	0.2412

A set of parity plots for all four models is presented in 2. To facilitate easier comparison at experimentally-relevant energy ranges, we have zoomed the plot in such that the highest exfoliation energy is 2 eV. Plots showing the entire energy range explored can be found in Supporting Information Figure S6. Here, we find that all models perform generally poorly, with the largest errors occurring at higher exfoliation energies in the case of XGBoost, TPOT, and SISSO. (see Figure 3). The best test-set R² and RMSE this time is only TPOT, although they are still relatively poor, with a test-set R² of only 0.543. Roost displays the best test-set MAE, although the model seems to have overfit, as it displays drastically better performance on the training set than it does on the test set. The XGBoost model performs slightly worse than either TPOT or Roost, and the SISSO approach did not perform well for this problem.

The TPOT algorithm again results in a relatively simple model pipeline. The first stage is a `SelectFwe` unit, which down-selects the features according to the Family-wise Error (FWE)^[23]. An alpha value of 0.047 is selected for this purpose, with results in a down-selection of the 299 input features into 125 features. This is then fed into an `ExtraTreesRegressor` unit, which is an implementation of the Extremely Randomized Trees method proposed by Geurts, Ernst, and Wehenkel^[84].

We again extract features from the XGBoost model (Supporting Information Figure S3), and find the Mendeleev Number again appears as an important feature, albeit as the maximum instead of the minimum. Additionally, we see descriptors related to bond strengths in the corresponding elemental systems: average melting points, and average heats of evaporation.

The list of preselected features can be found in the Supporting Information Table S4. Overall, we see that this set of features is the least similar out of all three problems with the bulk modulus, thermal conductivity, and decomposition energy being the most prevalent. Additionally, the bandgap of the material is also selected, suggesting the possible hierarchical learning. For this calculation additional features such as the decomposition energy, ratio between the perimeter and area of the surface, and the electronic bandgap of the material are also included.

The best SISSO model found for this problem is

$$E_{Exfoliation} \approx c_0 + a_0 \cdot \frac{P}{A} \cdot q_{evaporation}^{min} \cdot \sqrt{\sum_i \kappa_i} + a_1 \cdot \frac{E_{decomposition}}{V_{ICSD}^{ave} \cdot r_{vdw,uff}^{min}{}^3} \quad (3)$$

where $c_0 = 0.327$, $a_0 = 1.24 \times 10^{-4}$, $a_1 = 9.26 \times 10^8$, $\frac{P}{A}$ is the perimeter to area ratio of the surface, $q_{evaporation}^{min}$ is the minimum atomic evaporation heat of each element in the material, $E_{decomposition}$ is the decomposition energy, κ_i is the thermal conductivity of the elements at 298 K, V_{ICSD}^{ave} is the average atomic volume in the ICSD database, and $r_{vdw,uff}^{min}$ is the minimum atomic Van der Waals radius from UFF for each element in the material.

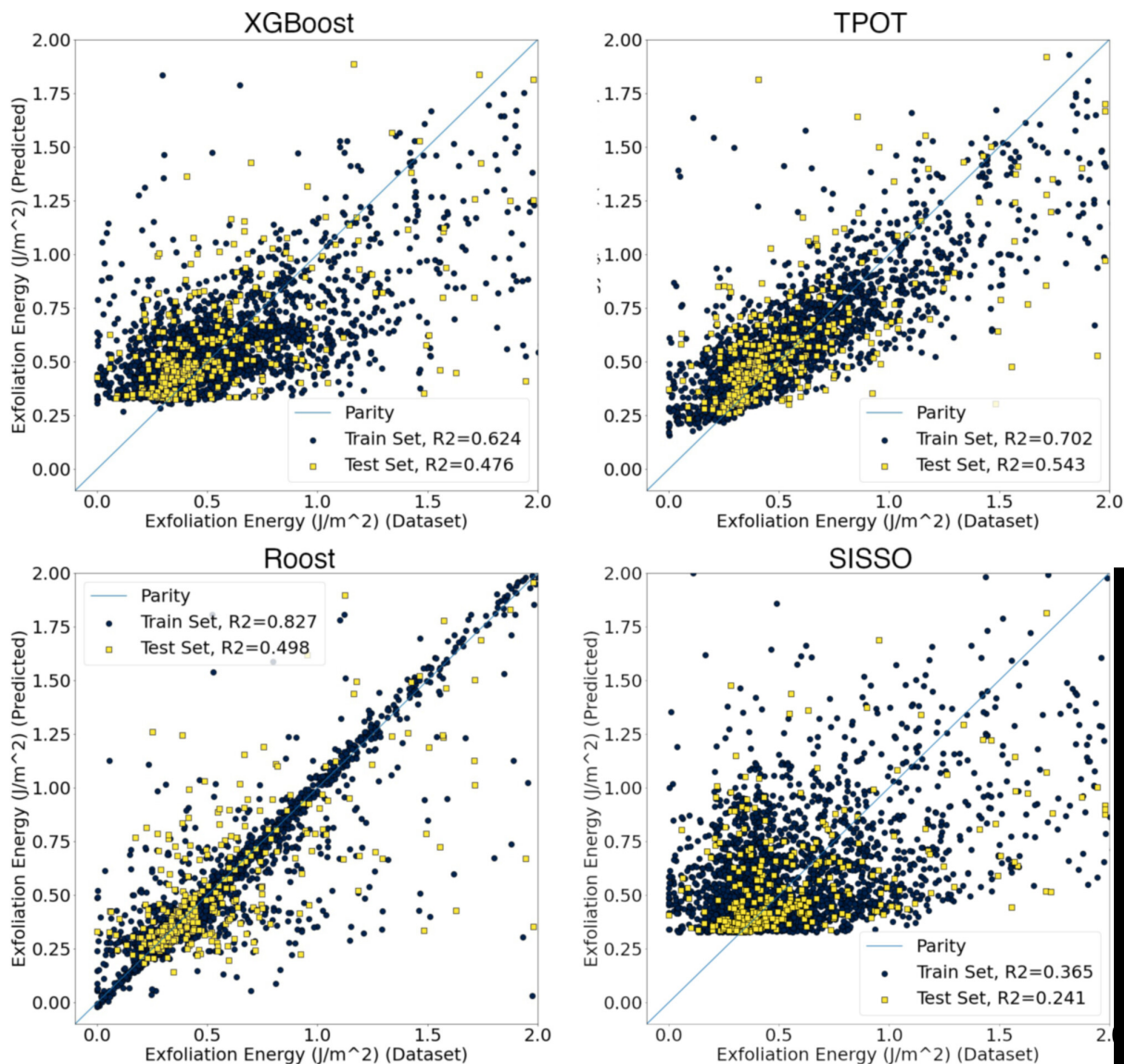


Figure 3: Parity plots for the XGBoost, TPOT, Roost, and SISSO models on the 2D material exfoliation energy problem. Included are the training and testing sets. A diagonal line indicating parity is drawn as a guide to the eye. Regression statistics for the models shown on this plot can be found in Table 4. To facilitate comparison at energy ranges that are more experimentally-relevant, we have zoomed in the plot to study energies no higher than 2 eV. The full data range is plot in Supporting Information Figure S6.

4 Discussion

We have developed a series of models which are capable of generating predictions for (1) the volume per formula unit of a series of ABX_3 perovskites, (2) the DFT-calculated bandgap of several 2D materials, and (3) several 2D material exfoliation energies. These problems encompass a variety of outcomes that one may find when training models of predictive properties.

4.1 Perovskite Volume per Formula Unit

In the case of the volume per formula unit of ABX_3 perovskites, we observe all four model types perform well. Overall, we find that the volume per formula unit for ABX_3 perovskites can be predicted using only compositional descriptors (i.e. with no structural descriptors). The likely reason all four models perform well despite having no structural information is the general similarity in crystal structure between these systems — they are all perovskites, and therefore all possess very similar crystal structures. Supporting this is that the Roost model, which only leverages the chemical formula as an input, and which we did not optimize the hyperparameters or architecture for, performed just as well on this problem

— albeit with some systematic deviation from parity at intermediate volumes. Although interpretability is reduced by virtue of being a neural network, we can still achieve an important insight from this model — just by knowing the chemical formula of the system, we can achieve accurate predictions of perovskite volumes, which further justifies our use of compositional descriptors (see Section 2.2.1) on this problem as we move to the SISO, XGBoost, and TPOT models. Additionally, we note this performance was achieved with a training set containing only 129 entries in the dataset - compared to the original Roost paper^[31] that leveraged approximately 275,000 entries from the OQMD dataset^[85].

Like the Roost model, we have difficulty in interpreting the pipeline generated by TPOT. The TPOT model delivers the best performance — which is clearly visible from the parity plot in Figure 1. This performance came at a price, however, and the rather complex pipeline containing scaling, one-hot encoding, and linear support vector regression does not lend itself well to interpretation.

Entering into the realm of interpretability, although the XGBoost model does not produce a direct formula for perovskite volumes, we can still gain some insight using it. It is still, however, relatively accurate — and allows us access to a feature importance metric (see Supporting Information Figure S2). In this case, we see the five most-important features are the average Rahm^[79,80] atomic radius, average UFF^[81] atomic radius, sum of elemental velocities of sound in the material, average Ghosh^[86] electronegativity, and the sum of the Pyykko^[87] triple-bond covalent radii. Overall, we see a strong reliance on descriptors of atomic radius — which as we noted in the TPOT discussion makes intuitive sense.

Finally, the SISO model (Equation 1) offers the most direct interpretation, as it is simply an equation. Immediately, we can see that ground state DFT atomic volumes are important. This result is highly intuitive and is not surprising when we consider that i) we are predicting volume, therefore using an average atomic volume makes sense and ii) the TPOT and XGBoost models extensively leveraged atomic radius descriptors that are also related to the volume of the atoms.

The overall good performance of SISO for this application is promising, as it is one of the most accurate models, while being by far the most interpretable. This represents a key advantage to symbolic regression, as if you can find an accurate model, then it will be easy to understand and analyze the results. Moreover, we note that are not alone in the literature when it comes to leveraging SISO to generate models of perovskite properties — the last several years have seen success in the creation of models of perovskite properties with this tool. The work of Xie et al.^[60] achieved good success in predicting the octahedral tilt in ABO_3 perovskites, the work of Bartel et al.^[59] resulted in the creation of a new tolerance factor for ABX_3 perovskite formation, and Ihalage and Hao^[58] leveraged descriptors generated by SISO to predict the formation of quaternary perovskites with formula $(A_{1-x}A'_x)BO_3$ and $A(B_{1-x}B'_x)O_3$.

4.2 2D Material Bandgap

The 2D material bandgap models did not achieve the same performance as for the perovskite systems (see Figure 2). Even in the case of TPOT and SISO, which still had the best test-set performance by most metrics, there were a few outliers. Specifically, we find that the test-set MAE for the models ranged between 0.29 eV (TPOT and SISO) and 0.65 eV (Roost) relative to the PBE DFT calculations reported by the 2DMatPedia. Putting this number in perspective, we note the recent work of Tran et al.^[88], which benchmarked the bandgap predictions of several popular DFT functionals for many of the systems in the C2DB; the work identified that the PBE functional exhibited a MAE of 1.50 eV relative to the G_0W_0 method. Other investigators have studied the prediction of 2D material bandgaps: Rajan et al.^[89] also achieved a test-set MAE of 0.11 eV on a dataset of 23,870 MXene systems (which, as far as we are aware, has not been made publicly available) using a Gaussian Process regression approach, with DFT-calculated properties including the average M-X bond length, volume per atom, MXene phase, fand heat of formation, and compositional properties including the mean Van der Waals radius, standard deviation of periodic table group number, standard deviation of the ionization energy, and standard deviation of the meting temperature. Zhang et al.^[90] improved on this error slightly, achieving a test-set MAE of 0.10 eV on the C2DB dataset (around 4000 entries)^[14] with both Support-Vector Regression and Random Forest approaches, albeit using descriptors such as the fermi-energy density of states and total energy of the system (requiring further DFT work for additional prediction). In contrast to both approaches, which used DFT-calculated values that would need to be obtained for new systems to be predicted, the only DFT-calculated value we leverage in our feature set is a bulk bandgap tabulated on the Materials Project^[18]. Thus, although our TPOT model had a slightly higher MAE of 0.29 eV, we note that this would not require further DFT work to generate new predictions. As 2D systems

are still relatively new, we note that much more work has been performed in the 3D materials space, particularly in the leveraging of neural networks to predict bandgaps. The recent Atomistic Line Graph Neural Network (ALIGNN)^[91] reported a test-set MAE of 0.218 eV for the prediction of bulk materials hosted by Materials Project^[18] (which as of October 2021 has over 144,000 inorganic systems). The Materials Graph Network (MEGNet) architecture^[92] achieved a test-set MAE of 0.32 eV on the bulk systems of the Materials Project. Although these neural network models are on 3D systems, we note that they do not leverage DFT properties (which we re-iterate would cause any resulting model to require a DFT calculation for future prediction) and had access to much larger datasets than the training set we obtained after filtering the 2D MatPedia entries (see Section 2.3). Overall, although the systems we investigate are not 3D bulk systems, we believe this puts the TPOT MAE of 0.29 eV for the bandgap of 2D systems in perspective.

In all 4 models we trained, many of the incorrect predictions occur where the DFT bandgap is 0 eV (which represented 27% of the training set values). Because of this, we tried simplifying the bandgap problem, by training an XGBoost model to predict whether the system was a metal (see Supporting Information section S.6.3), and showed that we could achieve good results — although we incorporated a purely structural fingerprint, the Sine Matrix Eigenspectrum (see Supporting Information Section S.6.1). As this descriptor resulted in some rather large vectors (of length 40, the maximum number of atoms in any system) with little direct physical intuition, we do not directly include it for the purposes of this section. Instead, what we can derive from it is the knowledge that structural features can provide information to predict the bandgap.

If we take a closer look at the Roost model, which takes a purely compositional approach to materials property prediction, we can see a poor generalization to the test set (see Table 3). This indicates that we have likely caused it to over-fit (which could have been improved for example through the use of early stopping). Overall, we may also use this result to further show the need for structural descriptors when predicting bandgaps of these systems.

The TPOT, SISO, and XGBoost models for 2D material bandgap achieved similar performance to one-another (see Table 2). Moreover, we again have the opportunity to extract meaning from the TPOT model here, as it leveraged Elastic Net to perform its prediction — a model which is achieved by mixing together Ridge Regression and the LASSO. Additionally, as the L1/L2 ratio is 0.95, this is nearly entirely L1 regularization (see section 3.2). As a result of this, and that the data were all scaled such that they ranged between 0 and 1, we can view the the feature coefficients approximately (but not entirely) through the lens of LASSO, which does perform feature selection. Doing this, we see that the corresponding bulk 3D-parent material’s bandgap (as found on Materials Project) is the most strongly-weighted feature by the elastic net model. This is extremely intuitive as well — it is reasonable for the bulk bandgap to be correlated with the 2D bandgap.

The XGBoost model, however, suggests an entirely different set of descriptors as being “important.” Instead, we find that the five most-important features are, in order from greatest to least, the minimum Mendeleev number, average atomic weight, variance of the number of unfilled s-orbitals, average ground-state DFT volume, and maximum Cordero^[93] covalent radius. None of these descriptors have a particularly intuitive relationship with the bandgap either. This is a similar result to that found by Rajan et al.^[89], who leveraged several models to predict the GW bandgap of MXene systems; in this work, they found that several of the most-important features exhibited only a statistical (i.e. non-intuitive) relationship with the bandgap.

Finally, the SISO model (Equation 2) shows an equivalent performance to TPOT and XGBoost, while using only three primary features. Similar to the TPOT model the bandgap of the bulk 3D-parent material is the most important feature, with only minor non-linear contributions from the minimum atomic Van der Waals radius and period. Interestingly, the SISO model not only misses some of the metallic materials, but also incorrectly predicts some materials to metallic in the training set. This suggests the increase simplicity of the SISO models may slightly reduce their reliability.

Future work on this problem may achieve better performance on the bandgap problem by investigating the bond strengths of the different elements in the system. We also note the very good performance that recent neural network approaches have had on the 3D bandgap problem^[91,92], likely due to their representation of the structure of the 3D systems. Similar to how the Roost model achieves good success when compositional descriptors are appropriate, we may find good success in leveraging neural network approaches when structural features are required. We note here that Deng et al.^[28] achieved good results on a variety of molecule properties by incorporating various graph representations from different neural network architectures. Hence, future work in this domain may benefit from the incorporation of the information-dense structural fingerprints that may be obtained from neural network-based approaches.

4.3 2D Material Exfoliation Energy

We observed some of the worst model performance (across all models) in the case of the 2D material exfoliation energy. Despite being a larger dataset than either the perovskite (144 total, 129 in the training set) or bandgap (1,412 total, 1,270 in the training set), the 3,049 entries in the training set (out of 3,388 total) for the exfoliation energy proved insufficient to achieve good results for any of the models. Moreover, the compositional and structural features were not sufficient to adequately describe the system.

Some interpretation can at least be derived from the XGBoost and TPOT models. The five most-important features according to XGBoost (see Supporting Information Figure S3) are, in order, the maximum Mendeleev number, average melting point, average evaporation heat, maximum number of unfilled orbitals, and the sum of the melting points. In the case of the TPOT model, we arrived at an extremely randomized tree approach, which also has a feature importance metric. Here, we find that the average Van der Waals radius, maximum dipole polarizability, minimum atomic weight, maximum atomic weight, and maximum elemental (and tabulated) DFT bandgap are weighted as important. Between the two models, we see much difference in which features are deemed important. In the XGBoost model, the average melting point, sum of melting points, and average evaporation heat of the elemental systems can be rationalized if we realize that these are all driven by the strength of the interaction between atoms in the material, thus these descriptors may provide information relating to the forces that must be overcome when exfoliation is performed. Many of the other features in the two models correlate with size: the maximum Mendeleev number, average Van der Waals radius, and minimum / maximum atomic weights all provide a description of the size of the atoms in the system.

Finally, the SISSO model (Equation 3) echoes these findings. Although it is less performant, it provides intuitive descriptors that tell a similar story. One term takes the ratio of the decomposition energy and approximations to the atomic volume. This captures a description of atomic size as well as the strength of atomic bonds. Additionally, the second term in the model incorporates information about the surface, thermal conductivity and the heat of evaporation. Although the second term is less descriptive, it still captures terms relating to the size of the atoms and bond strengths of that atoms involved in the 2D material. The relatively poor performance of the SISSO models, also highlights the need for better input features to describe the exfoliation energy of a material. In cases where only a loose correlation between a target property and the inputs exist, the relative simplicity of symbolic regression models can be detrimental. While TPOT and XGBoost can utilize information from all features in their final predictions, symbolic regression in general and SISSO in particular only acts on the order of tens of features. Because of this, it is likely that more structural information is needed to get accurate models with SISSO.

In effect, the models can be interpreted as collectively describing the exfoliation energy as a function of i) the size of the atoms in the 2D material and ii) the strength of the intermolecular forces between those atoms. When we predict exfoliation energies, we’re predicting the interaction between layers in an exfoliable material. Overall, finding better methods of cheaply approximating these weak interactions may provide better results in the prediction of exfoliation. Additionally, as the number of datasets which contain exfoliation energies increases (such as the 2DMatPedia^[11], C2DB^[14,15] and JARVIS^[17]), further insight into this problem will be possible, and more-complex (albeit less-interpretible) models will become feasible.

Additionally, in order to obtain more-accurate predictions of exfoliation energy, data generated via a more thorough computational treatment may be required. We illustrate this by examining an outlier in the training set at 9.9 J / m², which all four models heavily under-predicted (see Supporting Information figure S6) (7-8 eV in the case of XGBoost and TPOT, and over 9 eV in the case of Roost and SISSO). Upon closer examination of this system, we find that it is actually a pair of layers containing N atoms (Figure 4 A). The 2DMatPedia^[11] reports that this system (2dm-id 5985) was not directly sourced by a simulated exfoliation from a bulk structure, but instead was obtained by substituting the atoms in a hypothetical 2D Sb structure (Figure 4 B). The Sb structure (2dm-id 4275) was obtained by a simulated exfoliation from a structure obtained from materials project (Figure 4 C). The parent bulk material (mp-567409) is reported by the Materials Project^[18] to be a monoclinic crystal which undergoes a favorable decomposition (energy above hull is reported as 0.121 eV / atom) to a triclinic system. That being said, as this is a hypothetical 2D system, comparison with the hypothetical 3D bulk system was necessary for the calculation of exfoliation energy. As the prediction of crystal structure is a very challenging field with few easy approximations^[94], this may have contributed further to the extreme value of the exfoliation energy. Indeed, as Zhou et al report^[11], the decomposition energy lends itself better to assessing whether a material is truly stable. Indeed, despite the extremely high exfoliation energy of this hypothetical 2D

N system, it is reported by the 2DMatPedia to have a decomposition energy of 0 eV/atom. This too seems somewhat high, as systems containing N-N bonds tend to be high-energy materials, typically undergoing strongly exothermic decomposition to inert, gaseous N_2 ^[95]. With this in conjunction with the observation that our models all predict exfoliation energies significantly lower than the tabulated values, we have reason to believe that this system would be far easier to exfoliate than the 10 eV exfoliation energy implies. Moreover, this system may have a strong energetic preference to decompose further into N_2 , which additional DFT work could reveal. Overall, this underscores the importance of obtaining high-quality data, and filtering that high-quality data, for the training of interpretable models.

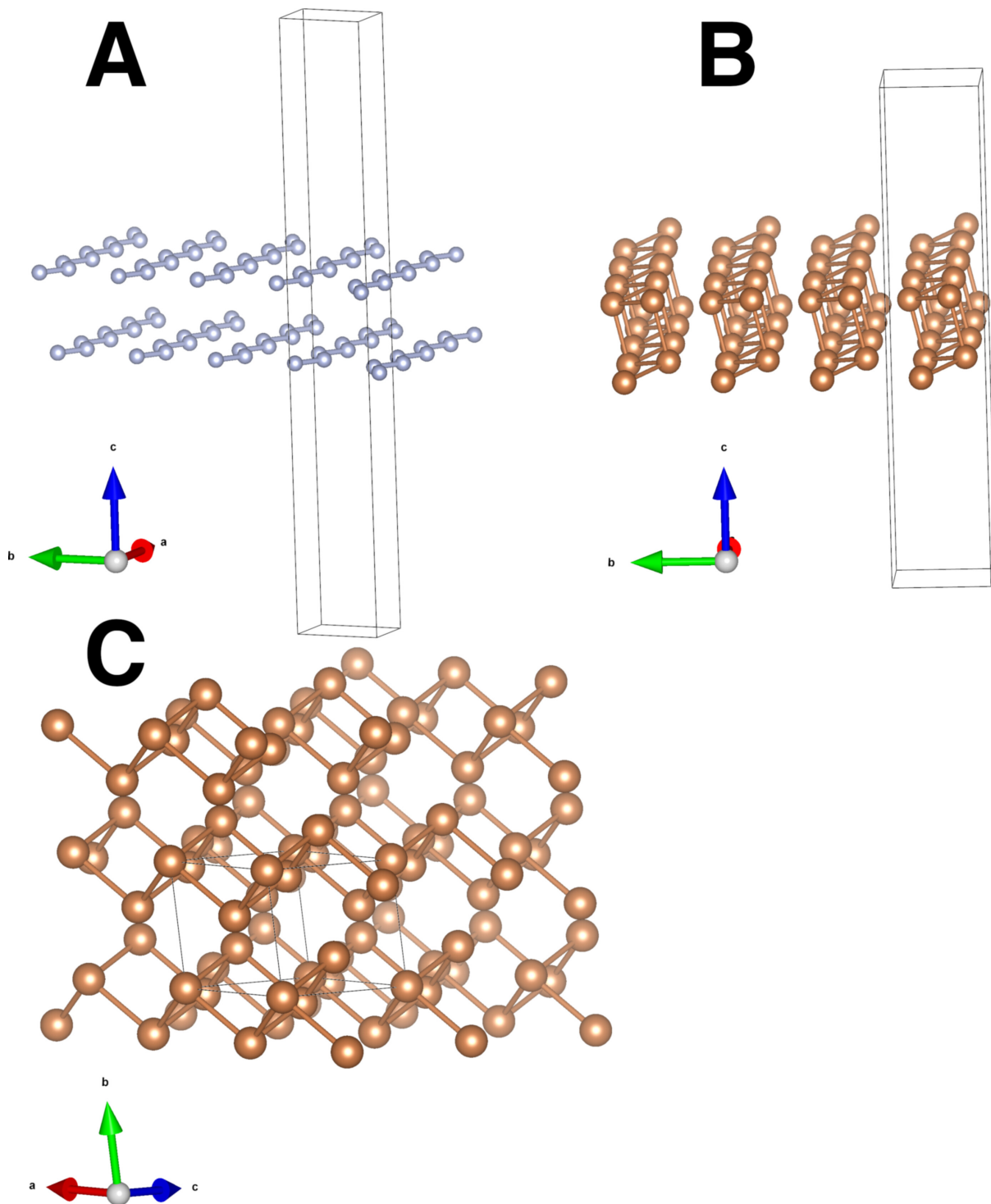


Figure 4: Illustrations of of A) a N-containing system (2dm-id 5985) which persisted as a large outlier across all exfoliation models in the training set, B) the Sb structure (2dm-id 4275) the N-containing system was derived from, and C) the bulk structure from Materials Project (mp-567409) from-which the exfoliation of the Sb system was simulated.

4.4 Future Outlook

As ML is further integrated into materials discovery workflows, we anticipate that the numerous successes neural networks have presented^[1–3,29–41] will continue to propel them onto the cutting edge of chemical property prediction. This comes with the challenge of honing our techniques for their interpretation, an area which has seen much interest in recent years, and where there is still plenty of opportunity for further development^[96,97]. We also expect AutoML techniques such as TPOT will continue gaining traction in materials discovery, due to the amount of success and attention they have recently had^[42–47]. This too presents the challenge of interpretability if highly complex pipelines are generated (see Sections 3.1 and 4.1). We note here that part of the value that AutoML techniques bring is the ability to make advanced techniques accessible to a wider audience of researchers by lowering the barrier of entry. Hence, we expect that the problem of interpretation may be compounded for AutoML (and especially NAS) systems: the ability to automatically extract some level of interpretation from the generated pipelines is important for automation to make ML truly accessible to non-experts. Overall, we expect that as neural network models and AutoML algorithms continue to grow in capability and complexity, work in developing the tools and techniques needed to interpret them will see a greater attention.

In contrast with challenge of interpreting neural networks or the pipelines found by AutoML systems, symbolic regression tools like Eureka and SISSO yield an exact equation describing the model, and are thus easier to interpret. This makes it easier to achieve key insights with physical interpretations — such as the very intuitive way in which SISSO is able to describe the systems. Overall, despite its reduced ability to predict the exfoliation energy of a material when compared to the models of TPOT, XGBoost, and Roost, we note the mathematical equations returned by SISSO provide a direct relationship between the target properties and model predictions. Additionally, in the case of the exfoliation energy, we believe that we may see further improvements by including richer structural information. We base this on the observation that the Roost model performed poorly on both of this problems – recalling that Roost is only provided the chemical formula of the system, this could indicate that compositional descriptors alone are insufficient to describe these properties. Indeed, it is well-known that structure and energy are intimately related (the fundamental assumption of geometry optimization techniques is that energy is a function of atomic position), hence it can be inferred that exfoliation energy and structure are similarly related. In the case of bandgaps, we note that there is also a strong dependence on structure; Chaves et al^[98] notes that the number of layers in a 2D material can strongly influence the band gap, reporting differences of up to several eV can occur between the bulk and monolayer form of a material.

Interoperability is still a challenge in the materials discovery ecosystem. Although it is possible to easily convert between different chemical file formats (e.g. by OpenBabel^[99]), and packages such as Pymatgen^[67], Atomic Simulation Environment (ASE)^[100], and RDKit^[101] can easily convert to each others’ format, we note that there is a challenge of calculating features using a variety of different packages. Some tools expect Pymatgen objects (e.g. XenonPy), others expect ASE objects, whereas others require RDKit objects (e.g. all of the descriptors in the RDKit library) to perform a calculation of features, thus creating some standard for the interoperability of these packages would be beneficial. Additionally, further efforts should be made to report the sources of data used by featurization packages. We note that MatMiner^[70] is exemplary in this regard: each of the featurization classes it defines has a “citation” method returning the appropriate source to credit. Mendeleev^[66] is another good example of this; within its documentation, a table lists citations for many (though not all) of the elemental properties it can return. Overall, by placing a stronger focus on i) interoperability and ii) data provenance, the Python materials modeling ecosystem can be made stronger — and therefore help accelerate materials discovery.

All of the models we have investigated in this work required sufficient training data to avoid over-fitting. Although techniques such as cross-validation, early-stopping (in the case of neural networks and XGBoost), and train/test splitting can help guard against (and detect) over-fitting, having a sufficiently-large dataset is of the utmost importance to achieve truly generalizable models. As a result, there is a critical need for data management approaches that satisfy the set of FAIR principles. This crucial need for effective data management has led to the incorporation of data storage tooling in popular chemistry packages including Pymatgen^[67], ASE^[100], and RDKit^[101]. Moreover, advances in both computational capacity and techniques has given rise to studies performing the high-throughput screening of chemical systems^[102–104]. This has resulted in the development of tools focusing on the provenance of data, such as the Automated Interactive Infrastructure and Database for Computational Science (AiiDA) system^[105,106].

Overall, we have identified a series of key issues should see more attention as the digital ecosystem

surrounding materials modeling continues to develop. First, interpretability of models allows us to derive physical understanding from the available data. This is a key benefit of symbolic regression tools like SISSO, which result in the creation of human-readable equations describing the model. Additionally, increasing the accessibility of ML techniques through automation (such as in the field of AutoML) will allow a wider range of researchers the ability to benefit from advances in modeling techniques. Data management and data provenance is another major issue, which allows us to better understand which datasets can be combined (e.g. when combining DFT datasets, the methodologies should be consistent between them), and to help us understand if something intrinsic to the training data is affecting model performance. These data management goals are core focus of platforms such as Exabyte^[107], which provides an all-in-one solution for i) storing material data and metadata, ii) storing the methodology required to derive a property from a material, and iii) providing the means to automatically perform calculations, and iv) automatically extracting calculation results and storing them for the user. This focus on providing a tool that manages materials, workflows, and calculations has allowed Exabyte to be a highly successful platform, which has led to studies involving automated phonon calculations^[108], high-throughput screening of materials for their band-structure^[109,110]. Future capabilities of the platform are slated to include a categorization scheme for computational models to provide even more metadata to track the provenance of calculated material properties^[111].

5 Conclusion

In this work, we have performed a series of benchmarks on a diverse set of ML algorithms: gradient boosting (XGBoost), AutoML (TPOT), deep learning (Roost), and symbolic regression (SISSO). These models were used to predict i) the volume of perovskites, ii) the DFT bandgap of 2D materials, and iii) the exfoliation energy of 2D materials. We identify that TPOT, SISSO and XGBoost tend to produce more-accurate models than Roost, but Roost works well in systems where compositional descriptors are enough to predict the target property. Finally, although SISSO was unable to find an accurate model for the exfoliation energy, it provides a human-readable equation describing the model, facilitating an easier interpretation compared to the other algorithms. We believe that interpretability will remain a key challenge to address as complex techniques (i.e. neural networks and AutoML) become more mainstream within the digital materials modeling ecosystem. Overall, as tools improving the accessibility of machine-learning continue to be developed, data provenance and model interpretability will become even more important, as it is a critical part of ensuring the accessibility of these techniques. By working to ensure that a wider audience of researchers can achieve insight from the rich digital ecosystem of materials design, materials discovery can be accelerated.

Acknowledgements

This research was supported by the US Department of Energy (DoE) Small Business Innovation Research (SBIR) program (grant no. DE-SC0021514). Computational resources were provided by Exabyte Inc. The authors also wish to acknowledge fruitful discussions with Rhys Goodall (University of Cambridge) regarding the Roost framework.

Data Access

Jupyter (Python) notebooks are available on Exabyte’s GitHub (https://github.com/Exabyte-io/Scientific-Projects/tree/arXiv_interpretableML_nov2021), which contains code to reproduce our results and figures.

Disclosures

James Dean and Timur Bazhirov are employed by Exabyte Inc.

References

- [1] Claudia Draxl and Matthias Scheffler. Big Data-Driven Materials Science and Its FAIR Data Infrastructure. In Wanda Andreoni and Sidney Yip, editors, *Handbook of Materials Modeling: Methods: Theory and Modeling*, pages 49–73. Springer International Publishing, Cham, 2020. ISBN 978-3-319-44677-6. doi: 10.1007/978-3-319-44677-6_104.
- [2] Adam C. Mater and Michelle L. Coote. Deep Learning in Chemistry. *Journal of Chemical Information and Modeling*, 59(6):2545–2559, June 2019. ISSN 1549-9596. doi: 10.1021/acs.jcim.9b00266.
- [3] Keith T. Butler, Daniel W. Davies, Hugh Cartwright, Olexandr Isayev, and Aron Walsh. Machine learning for molecular and materials science. *Nature*, 559(7715):547–555, July 2018. ISSN 1476-4687. doi: 10.1038/s41586-018-0337-2.
- [4] Lutz Bornmann and Rüdiger Mutz. Growth rates of modern science: A bibliometric analysis based on the number of publications and cited references. *Journal of the Association for Information Science and Technology*, 66(11):2215–2222, 2015. ISSN 2330-1643. doi: 10.1002/asi.23329.
- [5] Derek J. De Solla Price. *Little Science, Big Science*. Columbia University Press, March 1963. ISBN 978-0-231-88575-1. doi: 10.7312/pric91844.
- [6] Derek J. de Solla Price. *Science since Babylon*. Yale University Press, New Haven, enl. ed edition, 1975. ISBN 978-0-300-01797-7.
- [7] Derek J. de Solla Price. Networks of Scientific Papers. *Science*, 149(3683):510–515, July 1965. doi: 10.1126/science.149.3683.510.
- [8] National Science and Technology Council. Materials Genome Initiative for Global Competitiveness. Government, White House Office of Science and Technology Policy, United States of America, June 2011.
- [9] Subcommittee on the Materials Genome Initiative Committee on Technology. Materials Genome Initiative Strategic Plan. Government, National Science and Technology Council, United States of America, November 2021.
- [10] Juan J. de Pablo, Nicholas E. Jackson, Michael A. Webb, Long-Qing Chen, Joel E. Moore, Dane Morgan, Ryan Jacobs, Tresa Pollock, Darrell G. Schlom, Eric S. Toberer, James Analytis, Ismaila Dabo, Dean M. DeLongchamp, Gregory A. Fiete, Gregory M. Grason, Geoffroy Hautier, Yifei Mo, Krishna Rajan, Evan J. Reed, Efrain Rodriguez, Vladan Stevanovic, Jin Suntivich, Katsuyo Thornton, and Ji-Cheng Zhao. New frontiers for the materials genome initiative. *npj Computational Materials*, 5(1):1–23, April 2019. ISSN 2057-3960. doi: 10.1038/s41524-019-0173-4.
- [11] Jun Zhou, Lei Shen, Miguel Dias Costa, Kristin A. Persson, Shyue Ping Ong, Patrick Huck, Yunhao Lu, Xiaoyang Ma, Yiming Chen, Hanmei Tang, and Yuan Ping Feng. 2DMatPedia, an open computational database of two-dimensional materials from top-down and bottom-up approaches. *Scientific Data*, 6(1):86, June 2019. ISSN 2052-4463. doi: 10.1038/s41597-019-0097-3.
- [12] Stefano Curtarolo, Wahyu Setyawan, Gus L. W. Hart, Michal Jahnatek, Roman V. Chepulskii, Richard H. Taylor, Shidong Wang, Junkai Xue, Kesong Yang, Ohad Levy, Michael J. Mehl, Harold T. Stokes, Denis O. Demchenko, and Dane Morgan. AFLOW: An automatic framework for high-throughput materials discovery. *Computational Materials Science*, 58:218–226, June 2012. ISSN 0927-0256. doi: 10.1016/j.commatsci.2012.02.005.
- [13] Stefano Curtarolo, Wahyu Setyawan, Shidong Wang, Junkai Xue, Kesong Yang, Richard H. Taylor, Lance J. Nelson, Gus L. W. Hart, Stefano Sanvito, Marco Buongiorno-Nardelli, Natalio Mingo, and Ohad Levy. AFLOWLIB.ORG: A distributed materials properties repository from high-throughput ab initio calculations. *Computational Materials Science*, 58:227–235, June 2012. ISSN 0927-0256. doi: 10.1016/j.commatsci.2012.02.002.
- [14] Morten Niklas Gjerding, Alireza Taghizadeh, Asbjørn Rasmussen, Sajid Ali, Fabian Bertoldo, Thorsten Deilmann, Nikolaj Rørbaek Knøsgaard, Mads Kruse, Ask Hjorth Larsen, Simone Manti, Thomas Garm Pedersen, Urko Petralanda, Thorbjørn Skovhus, Mark Kamper Svendsen, Jens Jørgen Mortensen, Thomas Olsen, and Kristian Sommer Thygesen. Recent progress of the

- computational 2D materials database (C2DB). *2D Materials*, 8(4):044002, July 2021. ISSN 2053-1583. doi: 10.1088/2053-1583/ac1059.
- [15] Sten Hastrup, Mikkel Strange, Mohnish Pandey, Thorsten Deilmann, Per S. Schmidt, Nicki F. Hinsche, Morten N. Gjerding, Daniele Torelli, Peter M. Larsen, Anders C. Riis-Jensen, Jakob Gath, Karsten W. Jacobsen, Jens Jørgen Mortensen, Thomas Olsen, and Kristian S. Thygesen. The Computational 2D Materials Database: High-throughput modeling and discovery of atomically thin crystals. *2D Materials*, 5(4):042002, September 2018. ISSN 2053-1583. doi: 10.1088/2053-1583/aacfc1.
- [16] David D. Landis, Jens S. Hummelshøj, Svetlozar Nestorov, Jeff Greeley, Marcin Dulak, Thomas Bligaard, Jens K. Nørskov, and Karsten W. Jacobsen. The Computational Materials Repository. *Computing in Science Engineering*, 14(6):51–57, November 2012. ISSN 1558-366X. doi: 10.1109/MCSE.2012.16.
- [17] Kamal Choudhary, Kevin F. Garrity, Andrew C. E. Reid, Brian DeCost, Adam J. Biacchi, Angela R. Hight Walker, Zachary Trautt, Jason Hattrick-Simpers, A. Gilad Kusne, Andrea Centrone, Albert Davydov, Jie Jiang, Ruth Pachter, Gowoon Cheon, Evan Reed, Ankit Agrawal, Xiaofeng Qian, Vinit Sharma, Houlong Zhuang, Sergei V. Kalinin, Bobby G. Sumpter, Ghanshyam Pili- nia, Pinar Acar, Subhasish Mandal, Kristjan Haule, David Vanderbilt, Karin Rabe, and Francesca Tavazza. The joint automated repository for various integrated simulations (JARVIS) for data- driven materials design. *npj Computational Materials*, 6(1):1–13, November 2020. ISSN 2057-3960. doi: 10.1038/s41524-020-00440-1.
- [18] Anubhav Jain, Shyue Ping Ong, Geoffroy Hautier, Wei Chen, William Davidson Richards, Stephen Dacek, Shreyas Cholia, Dan Gunter, David Skinner, Gerbrand Ceder, and Kristin A. Persson. Commentary: The Materials Project: A materials genome approach to accelerating materials innovation. *APL Materials*, 1(1):011002, July 2013. doi: 10.1063/1.4812323.
- [19] Claudia Draxl and Matthias Scheffler. The NOMAD laboratory: From data sharing to artificial intelligence. *Journal of Physics: Materials*, 2(3):036001, May 2019. ISSN 2515-7639. doi: 10.1088/2515-7639/ab13bb.
- [20] Scott Kirklin, James E. Saal, Bryce Meredig, Alex Thompson, Jeff W. Doak, Muratahan Aykol, Stephan Rühl, and Chris Wolverton. The Open Quantum Materials Database (OQMD): Assessing the accuracy of DFT formation energies. *npj Computational Materials*, 1(1):1–15, December 2015. ISSN 2057-3960. doi: 10.1038/npjcompumats.2015.10.
- [21] Tjeerd van der Ploeg, Peter C. Austin, and Ewout W. Steyerberg. Modern modelling techniques are data hungry: A simulation study for predicting dichotomous endpoints. *BMC Medical Research Methodology*, 14(1):137, December 2014. ISSN 1471-2288. doi: 10.1186/1471-2288-14-137.
- [22] Llew Mason, Jonathan Baxter, Peter Bartlett, and Marcus Frean. Boosting Algorithms as Gradient Descent. In *Advances in Neural Information Processing Systems*, volume 12. MIT Press, 2000.
- [23] Trevor Hastie, Robert Tibshirani, and J. H. Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer Series in Statistics. Springer, New York, 2nd ed edition, 2009. ISBN 978-0-387-84857-0.
- [24] Tianqi Chen and Carlos Guestrin. XGBoost: A Scalable Tree Boosting System. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 785–794, August 2016. doi: 10.1145/2939672.2939785.
- [25] Heng Liang, Kun Jiang, Tong-An Yan, and Guang-Hui Chen. XGBoost: An Optimal Machine Learning Model with Just Structural Features to Discover MOF Adsorbents of Xe/Kr. *ACS Omega*, 6(13):9066–9076, April 2021. doi: 10.1021/acsomega.1c00100.
- [26] Nadya Asanul Husna, Alhadi Bustamam, Arry Yanuar, Devvi Sarwinda, and Oky Hermansyah. The comparison of machine learning methods for prediction study of type 2 diabetes mellitus’s drug design. *AIP Conference Proceedings*, 2264(1):030010, 2020. doi: 10.1063/5.0024161.
- [27] Peter D. Ivatt and Mathew J. Evans. Improving the prediction of an atmospheric chemistry transport model using gradient-boosted regression trees. *Atmospheric Chemistry and Physics*, 20(13):8063–8082, July 2020. ISSN 1680-7316. doi: 10.5194/acp-20-8063-2020.

- [28] Daiguo Deng, Xiaowei Chen, Ruochi Zhang, Zengrong Lei, Xiaojian Wang, and Fengfeng Zhou. XGraphBoost: Extracting Graph Neural Network-Based Features for a Better Prediction of Molecular Properties. *Journal of Chemical Information and Modeling*, 61(6):2697–2705, June 2021. ISSN 1549-9596. doi: 10.1021/acs.jcim.0c01489.
- [29] Jörg Behler and Michele Parrinello. Generalized Neural-Network Representation of High-Dimensional Potential-Energy Surfaces. *Physical Review Letters*, 98(14):146401, April 2007. doi: 10.1103/PhysRevLett.98.146401.
- [30] 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):145301, April 2018. ISSN 0031-9007, 1079-7114. doi: 10.1103/PhysRevLett.120.145301.
- [31] Rhys E. A. Goodall and Alpha A. Lee. Predicting materials properties without crystal structure: Deep representation learning from stoichiometry. *Nature Communications*, 11(1):6280, December 2020. ISSN 2041-1723. doi: 10.1038/s41467-020-19964-7.
- [32] Jörg Behler. Four Generations of High-Dimensional Neural Network Potentials. *Chemical Reviews*, 121(16):10037–10072, August 2021. ISSN 0009-2665. doi: 10.1021/acs.chemrev.0c00868.
- [33] Kun Yao, John E. Herr, David W. Toth, Ryker Mckintyre, and John Parkhill. The TensorMol-0.1 model chemistry: A neural network augmented with long-range physics. *Chemical Science*, 9(8):2261–2269, February 2018. ISSN 2041-6539. doi: 10.1039/C7SC04934J.
- [34] Julia Westermayr, Michael Gastegger, and Philipp Marquetand. Combining SchNet and SHARC: The SchNarc machine learning approach for excited-state dynamics. *The Journal of Physical Chemistry Letters*, 11(10):3828–3834, May 2020. ISSN 1948-7185, 1948-7185. doi: 10.1021/acs.jpcclett.0c00527.
- [35] Kristof T. Schütt, Pieter-Jan Kindermans, Huziel E. Sauceda, Stefan Chmiela, Alexandre Tkatchenko, and Klaus-Robert Müller. SchNet: A continuous-filter convolutional neural network for modeling quantum interactions. *arXiv:1706.08566 [physics, stat]*, December 2017.
- [36] Alessandra Toniato, Philippe Schwaller, Antonio Cardinale, Joppe Geluykens, and Teodoro Laino. Unassisted noise reduction of chemical reaction datasets. *Nature Machine Intelligence*, 3(6):485–494, June 2021. ISSN 2522-5839. doi: 10.1038/s42256-021-00319-w.
- [37] Philippe Schwaller, Riccardo Petraglia, Valerio Zullo, Vishnu H. Nair, Rico Andreas Haeuselmann, Riccardo Pisoni, Costas Bekas, Anna Iuliano, and Teodoro Laino. Predicting retrosynthetic pathways using transformer-based models and a hyper-graph exploration strategy. *Chemical Science*, 11(12):3316–3325, March 2020. ISSN 2041-6539. doi: 10.1039/C9SC05704H.
- [38] Philippe Schwaller, Teodoro Laino, Théophile Gaudin, Peter Bolgar, Christopher A. Hunter, Costas Bekas, and Alpha A. Lee. Molecular Transformer: A Model for Uncertainty-Calibrated Chemical Reaction Prediction. *ACS Central Science*, 5(9):1572–1583, September 2019. ISSN 2374-7943. doi: 10.1021/acscentsci.9b00576.
- [39] Philippe Schwaller, Théophile Gaudin, Dávid Lányi, Costas Bekas, and Teodoro Laino. “Found in Translation”: Predicting outcomes of complex organic chemistry reactions using neural sequence-to-sequence models. *Chemical Science*, 9(28):6091–6098, July 2018. ISSN 2041-6539. doi: 10.1039/C8SC02339E.
- [40] Alain C. Vaucher, Philippe Schwaller, Joppe Geluykens, Vishnu H. Nair, Anna Iuliano, and Teodoro Laino. Inferring experimental procedures from text-based representations of chemical reactions. *Nature Communications*, 12(1):2573, May 2021. ISSN 2041-1723. doi: 10.1038/s41467-021-22951-1.
- [41] Jane Panteleev, Hua Gao, and Lei Jia. Recent applications of machine learning in medicinal chemistry. *Bioorganic & Medicinal Chemistry Letters*, 28(17):2807–2815, September 2018. ISSN 0960-894X. doi: 10.1016/j.bmcl.2018.06.046.
- [42] Pieter Gijsbers, Erin LeDell, Janek Thomas, Sébastien Poirier, Bernd Bischl, and Joaquin Vanschoren. An Open Source AutoML Benchmark. *arXiv:1907.00909 [cs, stat]*, July 2019.

- [43] Quanming Yao, Mengshuo Wang, Yuqiang Chen, Wenyuan Dai, Yu-Feng Li, Wei-Wei Tu, Qiang Yang, and Yang Yu. Taking Human out of Learning Applications: A Survey on Automated Machine Learning. *arXiv:1810.13306 [cs, stat]*, December 2019.
- [44] Xin He, Kaiyong Zhao, and Xiaowen Chu. AutoML: A survey of the state-of-the-art. *Knowledge-Based Systems*, 212:106622, January 2021. ISSN 0950-7051. doi: 10.1016/j.knosys.2020.106622.
- [45] Trang T Le, Weixuan Fu, and Jason H Moore. Scaling tree-based automated machine learning to biomedical big data with a feature set selector. *Bioinformatics*, 36(1):250–256, January 2020. ISSN 1367-4803. doi: 10.1093/bioinformatics/btz470.
- [46] Randal S. Olson, Ryan J. Urbanowicz, Peter C. Andrews, Nicole A. Lavender, La Creis Kidd, and Jason H. Moore. Automating Biomedical Data Science Through Tree-Based Pipeline Optimization. In Giovanni Squillero and Paolo Burelli, editors, *Applications of Evolutionary Computation*, Lecture Notes in Computer Science, pages 123–137, Cham, 2016. Springer International Publishing. ISBN 978-3-319-31204-0. doi: 10.1007/978-3-319-31204-0_9.
- [47] Randal S. Olson, Nathan Bartley, Ryan J. Urbanowicz, and Jason H. Moore. Evaluation of a Tree-based Pipeline Optimization Tool for Automating Data Science. In *Proceedings of the Genetic and Evolutionary Computation Conference 2016*, GECCO '16, pages 485–492, New York, NY, USA, July 2016. Association for Computing Machinery. ISBN 978-1-4503-4206-3. doi: 10.1145/2908812.2908918.
- [48] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [49] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. dAlché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- [50] Maryam Amir Haeri, Mohammad Mehdi Ebadzadeh, and Gianluigi Folino. Statistical genetic programming for symbolic regression. *Applied Soft Computing*, 60:447–469, November 2017. ISSN 1568-4946. doi: 10.1016/j.asoc.2017.06.050.
- [51] Kenneth E. Kinneer, William B. Langdon, Lee Spector, Peter J. Angeline, and Una-May O’Reilly. *Advances in Genetic Programming*. MIT Press, 1994. ISBN 978-0-262-19423-5.
- [52] Michael Schmidt and Hod Lipson. Distilling Free-Form Natural Laws from Experimental Data. *Science*, 324(5923):81–85, April 2009. doi: 10.1126/science.1165893.
- [53] David R. Stoutemyer. Can the Eureqa symbolic regression program, computer algebra and numerical analysis help each other? *arXiv:1203.1023 [cs]*, March 2012.
- [54] James Dean, Michael G. Taylor, and Giannis Mpourmpakis. Unfolding adsorption on metal nanoparticles: Connecting stability with catalysis. *Science Advances*, 5(9):eaax5101, 2019. doi: 10.1126/sciadv.aax5101.
- [55] Kaiyang Tan, Mudit Dixit, James Dean, and Giannis Mpourmpakis. Predicting Metal–Support Interactions in Oxide-Supported Single-Atom Catalysts. *Industrial & Engineering Chemistry Research*, 58(44):20236–20246, November 2019. ISSN 0888-5885. doi: 10.1021/acs.iecr.9b04068.
- [56] Runhai Ouyang, Stefano Curtarolo, Emre Ahmetcik, Matthias Scheffler, and Luca M. Ghiringhelli. SISSO: A compressed-sensing method for identifying the best low-dimensional descriptor in an immensity of offered candidates. *Physical Review Materials*, 2(8):083802, August 2018. ISSN 2475-9953. doi: 10.1103/PhysRevMaterials.2.083802.

- [57] Runhai Ouyang, Emre Ahmetcik, Christian Carbogno, Matthias Scheffler, and Luca M. Ghiringhelli. Simultaneous learning of several materials properties from incomplete databases with multi-task SISSO. *Journal of Physics: Materials*, 2(2):024002, March 2019. ISSN 2515-7639. doi: 10.1088/2515-7639/ab077b.
- [58] Achintha Ihalage and Yang Hao. Analogical discovery of disordered perovskite oxides by crystal structure information hidden in unsupervised material fingerprints. *npj Computational Materials*, 7(1):1–12, May 2021. ISSN 2057-3960. doi: 10.1038/s41524-021-00536-2.
- [59] Christopher J. Bartel, Christopher Sutton, Bryan R. Goldsmith, Runhai Ouyang, Charles B. Musgrave, Luca M. Ghiringhelli, and Matthias Scheffler. New tolerance factor to predict the stability of perovskite oxides and halides. *Science Advances*, 5(2):eaav0693, February 2019. doi: 10.1126/sciadv.aav0693.
- [60] Stephen R. Xie, Parker Kotlarz, Richard G. Hennig, and Juan C. Nino. Machine learning of octahedral tilting in oxide perovskites by symbolic classification with compressed sensing. *Computational Materials Science*, 180:109690, July 2020. ISSN 0927-0256. doi: 10.1016/j.commatsci.2020.109690.
- [61] Carlos Mera Acosta, Runhai Ouyang, Adalberto Fazio, Matthias Scheffler, Luca M. Ghiringhelli, and Christian Carbogno. Analysis of Topological Transitions in Two-dimensional Materials by Compressed Sensing. *arXiv:1805.10950 [cond-mat]*, May 2018.
- [62] Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5):206–215, May 2019. ISSN 2522-5839. doi: 10.1038/s42256-019-0048-x.
- [63] Sabine Körbel, Miguel A. L. Marques, and Silvana Botti. Stability and electronic properties of new inorganic perovskites from high-throughput ab initio calculations. *Journal of Materials Chemistry C*, 4(15):3157–3167, April 2016. ISSN 2050-7534. doi: 10.1039/C5TC04172D.
- [64] Claudia Draxl and Matthias Scheffler. NOMAD: The FAIR Concept for Big-Data-Driven Materials Science. *arXiv:1805.05039 [cond-mat, physics:physics]*, May 2018.
- [65] Hironao Yamada, Chang Liu, Stephen Wu, Yukinori Koyama, Shenghong Ju, Junichiro Shiomi, Junko Morikawa, and Ryo Yoshida. Predicting Materials Properties with Little Data Using Shotgun Transfer Learning. *ACS Central Science*, 5(10):1717–1730, October 2019. ISSN 2374-7943. doi: 10.1021/acscentsci.9b00804.
- [66] Lukasz Mentel. Mendeleev – A Python resource for properties of chemical elements, ions and isotopes, ver. 0.9.0, 2014.
- [67] Shyue Ping Ong, William Davidson Richards, Anubhav Jain, Geoffroy Hautier, Michael Kocher, Shreyas Cholia, Dan Gunter, Vincent L. Chevrier, Kristin A. Persson, and Gerbrand Ceder. Python Materials Genomics (pymatgen): A robust, open-source python library for materials analysis. *Computational Materials Science*, 68:314–319, February 2013. ISSN 09270256. doi: 10.1016/j.commatsci.2012.10.028.
- [68] John R Rumble, Thomas J Bruno, and Maria J Doa. *CRC Handbook of Chemistry and Physics: A Ready-Reference Book of Chemical and Physical Data*. CRC Press, Boca Raton, one hundred second edition, 2021. ISBN 978-0-367-71260-0.
- [69] 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, August 2016. ISSN 2057-3960. doi: 10.1038/npjcompumats.2016.28.
- [70] Logan Ward, Alexander Dunn, Alireza Faghaninia, Nils E. R. Zimmermann, Saurabh Bajaj, Qi Wang, Joseph Montoya, Jiming Chen, Kyle Bystrom, Maxwell Dylla, Kyle Chard, Mark Asta, Kristin A. Persson, G. Jeffrey Snyder, Ian Foster, and Anubhav Jain. Matminer: An open source toolkit for materials data mining. *Computational Materials Science*, 152:60–69, September 2018. ISSN 0927-0256. doi: 10.1016/j.commatsci.2018.05.018.

- [71] A. Salinas-Sanchez, J.L. Garcia-Muñoz, J. Rodriguez-Carvajal, R. Saez-Puche, and J.L. Martinez. Structural characterization of R₂BaCuO₅ (r = y, lu, yb, tm, er, ho, dy, gd, eu and sm) oxides by x-ray and neutron diffraction. *Journal of Solid State Chemistry*, 100(2):201–211, 1992. ISSN 0022-4596. doi: 10.1016/0022-4596(92)90094-C.
- [72] P. P. Ewald. Die Berechnung optischer und elektrostatischer Gitterpotentiale. *Annalen der Physik*, 369(3):253–287, 1921. ISSN 1521-3889. doi: 10.1002/andp.19213690304.
- [73] Jmol development team. Jmol, October 2016.
- [74] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A Next-generation Hyperparameter Optimization Framework. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD '19*, pages 2623–2631, New York, NY, USA, July 2019. Association for Computing Machinery. ISBN 978-1-4503-6201-6. doi: 10.1145/3292500.3330701.
- [75] James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for Hyper-Parameter Optimization. In *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011.
- [76] James Bergstra, Daniel Yamins, and David Cox. Making a Science of Model Search: Hyperparameter Optimization in Hundreds of Dimensions for Vision Architectures. In *Proceedings of the 30th International Conference on Machine Learning*, pages 115–123. PMLR, February 2013.
- [77] Lisha Li, Kevin Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. Hyperband: A novel bandit-based approach to hyperparameter optimization. *The Journal of Machine Learning Research*, 18(1):6765–6816, January 2017. ISSN 1532-4435.
- [78] P Linstrom, J. and W. G. Mallard, editors. *NIST Chemistry WebBook, NIST Standard Reference Database Number 69*. National Institute of Standards and Technology, Gaithersburg MD, 2021.
- [79] Martin Rahm, Roald Hoffmann, and N. W. Ashcroft. Atomic and Ionic Radii of Elements 1–96. *Chemistry – A European Journal*, 22(41):14625–14632, 2016. ISSN 1521-3765. doi: 10.1002/chem.201602949.
- [80] Martin Rahm, Roald Hoffmann, and N. W. Ashcroft. Corrigendum: Atomic and Ionic Radii of Elements 1–96. *Chemistry – A European Journal*, 23(16):4017–4017, 2017. ISSN 1521-3765. doi: 10.1002/chem.201700610.
- [81] A. K. Rappe, C. J. Casewit, K. S. Colwell, W. A. Goddard, and W. M. Skiff. UFF, a full periodic table force field for molecular mechanics and molecular dynamics simulations. *Journal of the American Chemical Society*, 114(25):10024–10035, December 1992. ISSN 0002-7863. doi: 10.1021/ja00051a040.
- [82] Michael W. Gaultois, Taylor D. Sparks, Christopher K. H. Borg, Ram Seshadri, William D. Bonificio, and David R. Clarke. Data-Driven Review of Thermoelectric Materials: Performance and Resource Considerations. *Chemistry of Materials*, 25(15):2911–2920, August 2013. ISSN 0897-4756. doi: 10.1021/cm400893e.
- [83] P. Villars, K. Cenzual, J. Daams, Y. Chen, and S. Iwata. Data-driven atomic environment prediction for binaries using the Mendeleev number: Part 1. Composition AB. *Journal of Alloys and Compounds*, 367(1):167–175, March 2004. ISSN 0925-8388. doi: 10.1016/j.jallcom.2003.08.060.
- [84] Pierre Geurts, Damien Ernst, and Louis Wehenkel. Extremely randomized trees. *Machine Learning*, 63(1):3–42, April 2006. ISSN 1573-0565. doi: 10.1007/s10994-006-6226-1.
- [85] Dipendra Jha, Logan Ward, Arindam Paul, Wei-keng Liao, Alok Choudhary, Chris Wolverton, and Ankit Agrawal. ElemNet: Deep Learning the Chemistry of Materials From Only Elemental Composition. *Scientific Reports*, 8(1):17593, December 2018. ISSN 2045-2322. doi: 10.1038/s41598-018-35934-y.
- [86] Dulal C. Ghosh. A new scale of electronegativity based on absolute radii of atoms. *Journal of Theoretical and Computational Chemistry*, 04(01):21–33, March 2005. ISSN 0219-6336. doi: 10.1142/S0219633605001556.

- [87] Pekka Pyykkö, Sebastian Riedel, and Michael Patzschke. Triple-Bond Covalent Radii. *Chemistry – A European Journal*, 11(12):3511–3520, 2005. ISSN 1521-3765. doi: 10.1002/chem.200401299.
- [88] Fabien Tran, Jan Doumont, Leila Kalantari, Peter Blaha, Tomáš Rauch, Pedro Borlido, Silvana Botti, Miguel A. L. Marques, Abhilash Patra, Subrata Jana, and Prasanjit Samal. Bandgap of two-dimensional materials: Thorough assessment of modern exchange–correlation functionals. *The Journal of Chemical Physics*, 155(10):104103, September 2021. ISSN 0021-9606. doi: 10.1063/5.0059036.
- [89] Arunkumar Chitteth Rajan, Avanish Mishra, Swanti Satsangi, Rishabh Vaish, Hiroshi Mizuseki, Kwang-Ryeol Lee, and Abhishek K. Singh. Machine-Learning-Assisted Accurate Band Gap Predictions of Functionalized MXene. *Chemistry of Materials*, 30(12):4031–4038, June 2018. ISSN 0897-4756. doi: 10.1021/acs.chemmater.8b00686.
- [90] Yu Zhang, Wenjing Xu, Guangjie Liu, Zhiyong Zhang, Jinlong Zhu, and Meng Li. Bandgap prediction of two-dimensional materials using machine learning. *PLOS ONE*, 16(8):e0255637, August 2021. ISSN 1932-6203. doi: 10.1371/journal.pone.0255637.
- [91] Kamal Choudhary and Brian DeCost. Atomistic Line Graph Neural Network for Improved Materials Property Predictions. *arXiv:2106.01829 [cond-mat]*, September 2021.
- [92] Chi Chen, Weike Ye, Yunxing Zuo, Chen Zheng, and Shyue Ping Ong. Graph Networks as a Universal Machine Learning Framework for Molecules and Crystals. *Chemistry of Materials*, 31(9):3564–3572, May 2019. ISSN 0897-4756. doi: 10.1021/acs.chemmater.9b01294.
- [93] Beatriz Cordero, Verónica Gómez, Ana E. Platero-Prats, Marc Revés, Jorge Echeverría, Eduard Cremades, Flavia Barragán, and Santiago Alvarez. Covalent radii revisited. *Dalton Transactions*, 1(21):2832–2838, May 2008. ISSN 1477-9234. doi: 10.1039/B801115J.
- [94] Artem R Oganov. *Modern Methods of Crystal Structure Prediction*. Wiley-VCH ; John Wiley [distributor, Weinheim, Germany; Chichester, 2011. ISBN 978-3-527-40939-6.
- [95] Dheeraj Kumar and Anil J. Elias. The Explosive Chemistry of Nitrogen. *Resonance*, 24(11):1253–1271, November 2019. ISSN 0973-712X. doi: 10.1007/s12045-019-0893-2.
- [96] Feng-Lei Fan, Jinjun Xiong, Mengzhou Li, and Ge Wang. On Interpretability of Artificial Neural Networks: A Survey. *IEEE Transactions on Radiation and Plasma Medical Sciences*, pages 1–1, 2021. ISSN 2469-7303. doi: 10.1109/TRPMS.2021.3066428.
- [97] Yu Zhang, Peter Tiño, Aleš Leonardis, and Ke Tang. A Survey on Neural Network Interpretability. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 5(5):726–742, October 2021. ISSN 2471-285X. doi: 10.1109/TETCI.2021.3100641.
- [98] A. Chaves, J. G. Azadani, Hussain Alsalman, D. R. da Costa, R. Frisenda, A. J. Chaves, Seung Hyun Song, Y. D. Kim, Daowei He, Jiadong Zhou, A. Castellanos-Gomez, F. M. Peeters, Zheng Liu, C. L. Hinkle, Sang-Hyun Oh, Peide D. Ye, Steven J. Koester, Young Hee Lee, Ph. Avouris, Xinran Wang, and Tony Low. Bandgap engineering of two-dimensional semiconductor materials. *npj 2D Materials and Applications*, 4(1):29, August 2020. ISSN 2397-7132. doi: 10.1038/s41699-020-00162-4.
- [99] Noel M. O’Boyle, Michael Banck, Craig A. James, Chris Morley, Tim Vandermeersch, and Geoffrey R. Hutchison. Open Babel: An open chemical toolbox. *Journal of Cheminformatics*, 3(1):33, October 2011. ISSN 1758-2946. doi: 10.1186/1758-2946-3-33.
- [100] Ask Hjorth Larsen, Jens Jørgen Mortensen, Jakob Blomqvist, Ivano E. Castelli, Rune Christensen, Marcin Du\textbackslashashlak, Jesper Friis, Michael N. Groves, Bjørk Hammer, Cory Hargus, Eric D. Hermes, Paul C. Jennings, Peter Bjerre Jensen, James Kermode, John R. Kitchin, Esben Leonhard Kolsbjerg, Joseph Kubal, Kristen Kaasbjerg, Steen Lysgaard, Jón Bergmann Maronsson, Tristan Maxson, Thomas Olsen, Lars Pastewka, Andrew Peterson, Carsten Rostgaard, Jakob Schiøtz, Ole Schütt, Mikkel Strange, Kristian S. Thygesen, Tejs Vegge, Lasse Vilhelmsen, Michael Walter, Zhenhua Zeng, and Karsten W. Jacobsen. The Atomic Simulation Environment— a Python Library for Working with Atoms. *Journal of Physics: Condensed Matter*, 29(27):273002, June 2017. ISSN 0953-8984. doi: 10.1088/1361-648X/aa680e.

- [101] Greg Landrum, Paolo Tosco, Brian Kelley, sriniker, gedeck, NadineSchneider, Riccardo Vianello, Ric, Andrew Dalke, Brian Cole, AlexanderSavelyev, Matt Swain, Samo Turk, Dan N, Alain Vaucher, Eisuke Kawashima, Maciej Wójcikowski, Daniel Probst, guillaume godin, David Cosgrove, Axel Pahl, JP, Francois Berenger, strets123, JLVarjo, Noel O'Boyle, Patrick Fuller, Jan Holst Jensen, Gianluca Sforza, and DoliathGavid. RDKit, 2021.
- [102] Bingbing Zhang, Xiaodong Zhang, Jin Yu, Ying Wang, Kui Wu, and Ming-Hsien Lee. First-Principles High-Throughput Screening Pipeline for Nonlinear Optical Materials: Application to Borates. *Chemistry of Materials*, 32(15):6772–6779, August 2020. ISSN 0897-4756. doi: 10.1021/acs.chemmater.0c02583.
- [103] Lorenz M Mayr and Dejan Bojanic. Novel trends in high-throughput screening. *Current Opinion in Pharmacology*, 9(5):580–588, October 2009. ISSN 1471-4892. doi: 10.1016/j.coph.2009.08.004.
- [104] James Dean, Michael J. Cowan, Jonathan Estes, Mahmoud Ramadan, and Giannis Mpourmpakis. Rapid Prediction of Bimetallic Mixing Behavior at the Nanoscale. *ACS Nano*, 14(7):8171–8180, July 2020. ISSN 1936-0851. doi: 10.1021/acsnano.0c01586.
- [105] Martin Uhrin, Sebastiaan P. Huber, Jusong Yu, Nicola Marzari, and Giovanni Pizzi. Workflows in AiiDA: Engineering a high-throughput, event-based engine for robust and modular computational workflows. *Computational Materials Science*, 187:110086, February 2021. ISSN 0927-0256. doi: 10.1016/j.commatsci.2020.110086.
- [106] Sebastiaan P. Huber, Spyros Zoupanos, Martin Uhrin, Leopold Talirz, Leonid Kahle, Rico Häuselmann, Dominik Gresch, Tiziano Müller, Aliaksandr V. Yakutovich, Casper W. Andersen, Francisco F. Ramirez, Carl S. Adorf, Fernando Gargiulo, Snehal Kumbhar, Elsa Passaro, Conrad Johnston, Andrius Merkys, Andrea Cepellotti, Nicolas Mounet, Nicola Marzari, Boris Kozinsky, and Giovanni Pizzi. AiiDA 1.0, a scalable computational infrastructure for automated reproducible workflows and data provenance. *Scientific Data*, 7(1):300, September 2020. ISSN 2052-4463. doi: 10.1038/s41597-020-00638-4.
- [107] Timur Bazhurov. Data-centric online ecosystem for digital materials science. *arXiv:1902.10838 [cond-mat, physics:physics]*, February 2019.
- [108] Timur Bazhurov and E. X. Abot. Fast and accessible first-principles calculations of vibrational properties of materials. *arXiv:1808.10011 [cond-mat, physics:physics]*, August 2018.
- [109] Protik Das, Mohammad Mohammadi, and Timur Bazhurov. Accessible computational materials design with high fidelity and high throughput. *arXiv:1807.05623 [cond-mat, physics:physics]*, July 2018.
- [110] Protik Das and Timur Bazhurov. Electronic properties of binary compounds with high fidelity and high throughput. *Journal of Physics: Conference Series*, 1290:012011, October 2019. ISSN 1742-6588, 1742-6596. doi: 10.1088/1742-6596/1290/1/012011.
- [111] Alexander Zech and Timur Bazhurov. CateCom: A practical data-centric approach to categorization of computational models. *arXiv:2109.13452 [cond-mat]*, September 2021.
- [112] 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. ISSN 1097-461X. doi: 10.1002/qua.24917.
- [113] Lauri Himanen, Marc O.J. Jäger, Eiaki V. Morooka, Filippo Federici Canova, Yashasvi S. Ranawat, David Z. Gao, Patrick Rinke, and Adam S. Foster. Dscribe: Library of descriptors for machine learning in materials science. *Computer Physics Communications*, 247:106949, February 2020. ISSN 00104655. doi: 10.1016/j.cpc.2019.106949.
- [114] Matthias Rupp, Alexandre Tkatchenko, Klaus-Robert Müller, and O. Anatole von Lilienfeld. Fast and Accurate Modeling of Molecular Atomization Energies with Machine Learning. *Physical Review Letters*, 108(5):058301, January 2012. doi: 10.1103/PhysRevLett.108.058301.
- [115] Guillaume Lemaître, Fernando Nogueira, and Christos K. Aridas. Imbalanced-learn: A python toolbox to tackle the curse of imbalanced datasets in machine learning. *Journal of Machine Learning Research*, 18(17):1–5, 2017.

S Supporting Information

S.1 XGBoost Feature Importance

The XGBoost approach allows us to assess the feature importance scores of the various model inputs. In this section, we report the importance scores for the 10 most-important features identified by each model. Feature importances are calculated using the default "Gain" metric available in XGBoost.

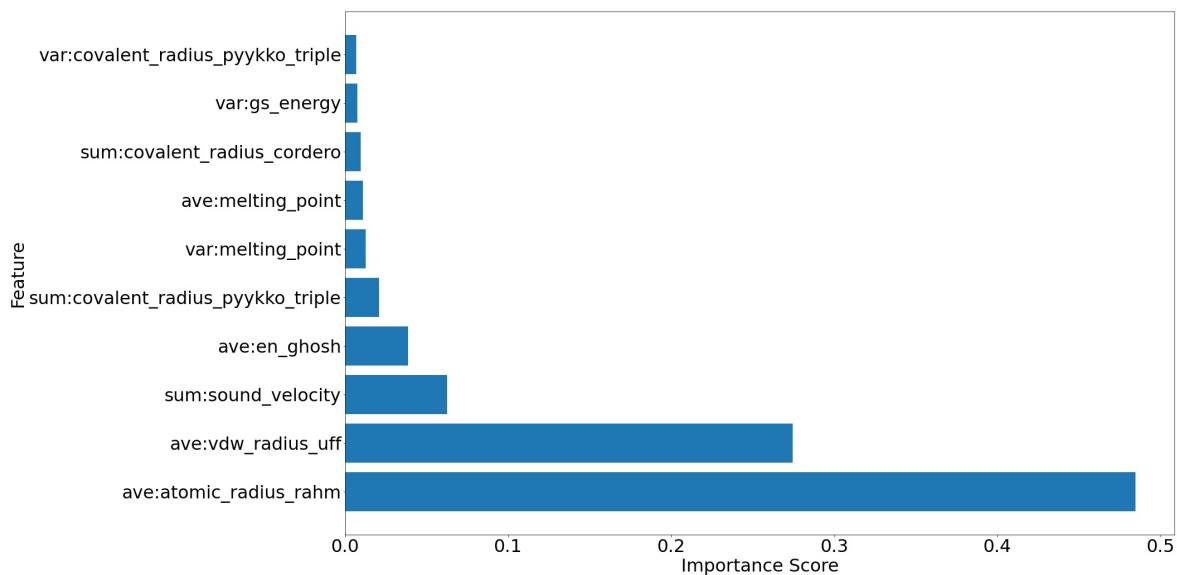


Figure S1: Importance metrics for the 10 most-important features identified by the XGBoost perovskite volume model using the "gain" feature importance metric.

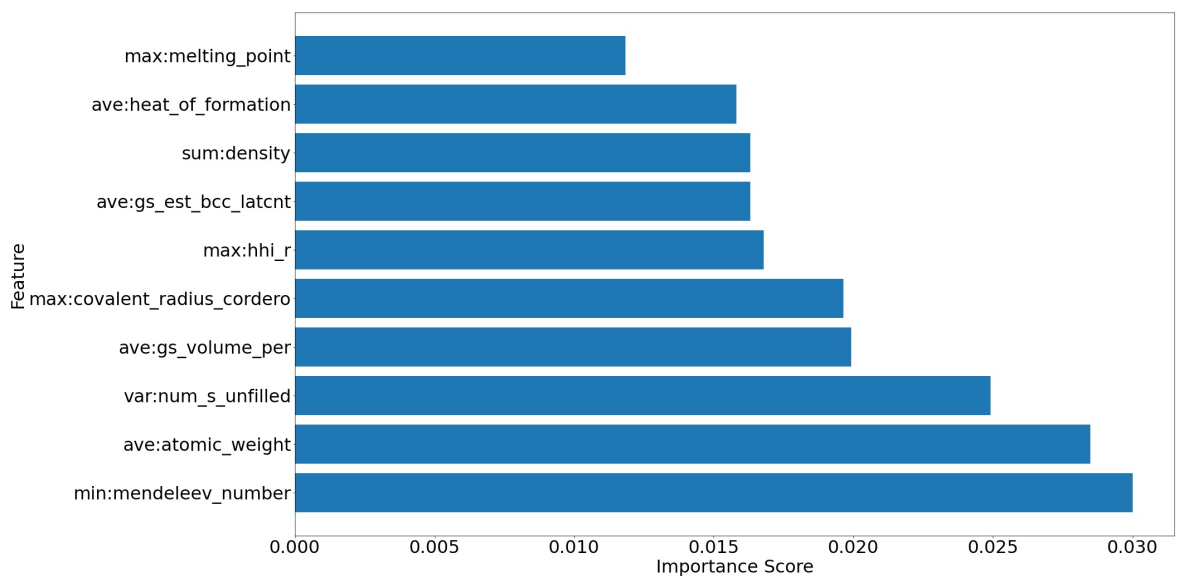


Figure S2: Importance metrics for the 10 most-important features identified by the XGBoost 2D material bandgap model using the "gain" feature importance metric.

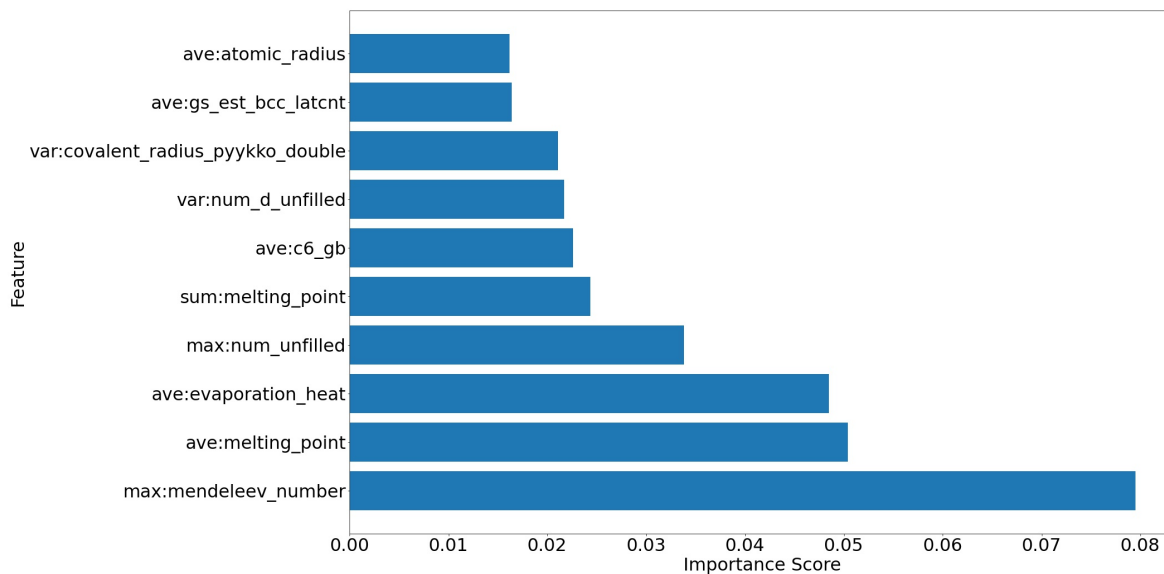


Figure S3: Importance metrics for the 10 most-important features identified by the XGBoost 2D material exfoliation model using the “gain” feature importance metric.

S.2 TPOT Model Components

To facilitate the reproducibility of our work, we have provided a list of the model components that TPOT could combine to generate models in S1. For more information on our specific methodology and other TPOT sections, see Section 2.4.2 in the manuscript.

Table S1: Search space of model components that are allowed by TPOT in its default regression configuration. Components are listed alongside the package they are sourced from, using the same class name defined in their respective package. Additionally, we list the role that the component.

Model Name	Source Package	Role
ElasticNetCV	Scikit-Learn	Estimator
ExtraTreesRegressor	Scikit-Learn	Estimator
GradientBoostingRegressor	Scikit-Learn	Estimator
AdaBoostRegressor	Scikit-Learn	Estimator
DecisionTreeRegressor	Scikit-Learn	Estimator
KNeighborsRegressor	Scikit-Learn	Estimator
LassoLarsCV	Scikit-Learn	Estimator
LinearSVR	Scikit-Learn	Estimator
XGBRegressor	XGBoost	Estimator
SGDRegressor	Scikit-Learn	Estimator
Binarizer	Scikit-Learn	Pre-Processing
FastICA	Scikit-Learn	Pre-Processing
FeatureAgglomeration	Scikit-Learn	Pre-Processing
MaxAbsScaler	Scikit-Learn	Pre-Processing
MinMaxScaler	Scikit-Learn	Pre-Processing
Normalizer	Scikit-Learn	Pre-Processing
Nystroem	Scikit-Learn	Pre-Processing
PCA	Scikit-Learn	Pre-Processing
PolynomialFeatures	Scikit-Learn	Pre-Processing
RBFSampler	Scikit-Learn	Pre-Processing
RobustScaler	Scikit-Learn	Pre-Processing
StandardScaler	Scikit-Learn	Pre-Processing
ZeroCount	TPOT	Pre-Processing
OneHotEncoder	TPOT	Pre-Processing
SelectFwe	Scikit-Learn	Feature Selection
SelectPercentile	Scikit-Learn	Feature Selection
VarianceThreshold	Scikit-Learn	Feature Selection
SelectFromModel	Scikit-Learn	Feature Selection
StackingEstimator	TPOT	Meta-Transformer

S.3 TPOT Feature Importance

Several of the models generated by TPOT can be used to assess how important a variable is. Although the linear support vector machine pipeline does not lend itself well to interpretation in the perovskite volume problem, the models created for the 2D material bandgap (Figure S4) and the 2D material exfoliation energy (Figure S5) lend themselves to some degree of interpretation.

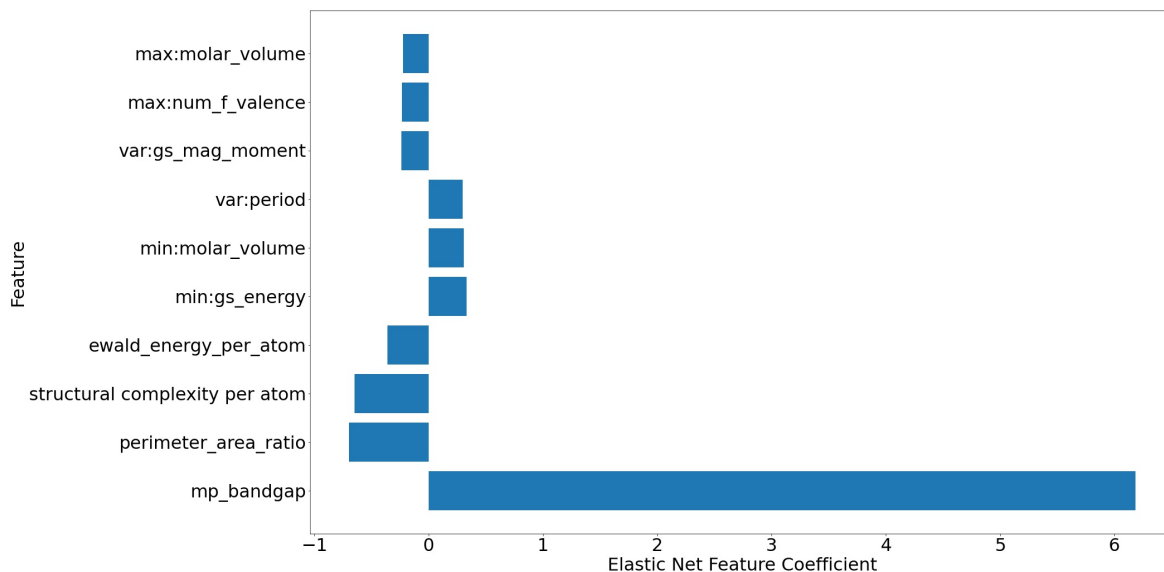


Figure S4: Importance metrics for the 10 most-important features for the 2D material bandgap problem, as identified by the Elastic Net model within the TPOT generated pipeline.

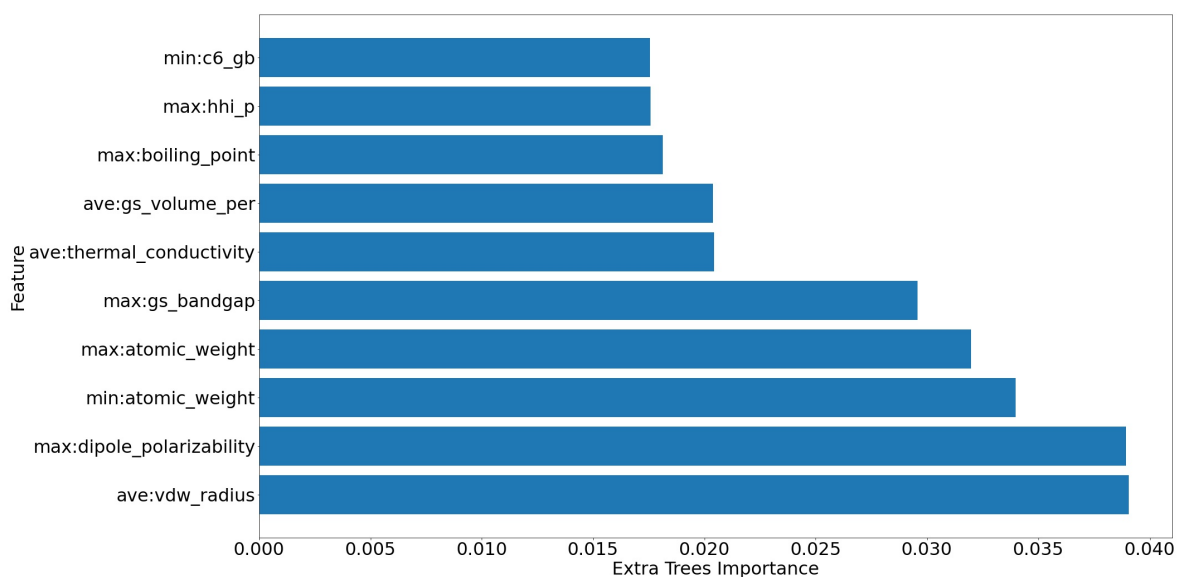


Figure S5: Importance metrics for the 10 most-important features on the 2D material exfoliation energy prediction problem, identified by the Extremely Randomized Trees model within the TPOT generated pipeline.

S.4 Prescreened features for SISSO

Here we present the set of thirty features used to find the SISSO models for each problem. Additionally, we describe the assumptions we made for the units of these descriptors.

Table S2: Results feature pre-screening for SISSO for the XenonPy compositional descriptors on the perovskite volume prediction problem. Also shown is the assumption we made for units of the XenonPy descriptors (see Section 2.4.4 for more details).

Feature	Units Assumption	Description
min:c6_gb	E_h/a_0^6	The min atomic C_6 dispersion coefficient
min:melting_point	K	The min elemental melting point
min:Polarizability	Å	The min elemental polarizability
min:atomic_number	—	The min atomic number
min:fusion_enthalpy	KJ / mol	The min fusion enthalpy
min:dipole_polarizability	Å ³	The min elemental dipole polarizability
min:boiling_point	K	The min elemental boiling point
min:gs_volume_per	cm ³ /mol	The min DFT volume per atom
min:covalent_radius_pyykko_triple	pm	The min single bond covalent radius
var:covalent_radius_pyykko_triple	pm	The var of the triple bond covalent radius
min:covalent_radius_slater	pm	The min covalent radius by Slater
min:covalent_radius_pyykko_double	pm	The min double bond covalent radius
min:period	—	The min period
min:covalent_radius_pyykko	pm	The min double bond covalent radius
min:covalent_radius_cordero	pm	The min covalent radius by Cordero et al
min:atomic_weight	Da	The min atomic weight
min:hhi_p	—	The min HHI production values
max:en_allen	eV	The max electronegativity by Allen et al
ave:covalent_radius_pyykko_triple	pm	The mean triple bond covalent radius
min:atomic_radius_rahm	pm	The min atomic radius by Rahm et al
ave:atomic_number	—	The mean atomic number
ave:gs_volume_per	cm ³ /mol	The mean DFT volume per atom
ave:heat_capacity_mass	J / g / K	The mean mass specific heat capacity
ave:covalent_radius_slater	pm	The mean covalent radius by Slater
ave:covalent_radius_pyykko	pm	The mean single bond covalent radius
max:mendeleev_number	—	The max Mendeleev number
min:vdw_radius_alvarez	pm	The min vdw radius by Alvarez
min:vdw_radius_mm3	pm	The min vdw radius according to the MM3 FF
max:en_ghosh	EN	The max electronegativity by Ghosh et al.
min:vdw_radius	pm	The min vdw radius

Table S3: Results feature pre-screening for SISSO for the compositional and structural descriptors defined in Section 2.2, for the 2D material bandgap prediction problem. Also shown is the assumption we made for units of the XenonPy descriptors (see Section 2.4.4 for more details).

Feature	Units Assumption	Description
mp:bandgap	eV	
min:c6_gb	E_h/a_0^6	The min atomic C_6 dispersion coefficient
min:covalent_radius_slater	pm	The min covalent radius by Slater
min:boiling_point	K	The min boiling point
min:covalent_radius_cordero	pm	The min covalent radius by Cordero et al
min:melting_point	K	The min melting point
min:dipole_polarizability	\AA^3	The min elemental dipole polarizability
min:covalent_radius_pyykko	pm	The min single bond covalent radius
min:Polarizability	\AA^3	The min elemental polarizability
ave:boiling_point	K	The mean boiling point
min:period	—	The min period
ave:density	g/cm^3	The mean density
ave:covalent_radius_cordero	pm	The mean covalent radius by Cordero et al
ave:covalent_radius_slater	pm	The mean covalent radius by Slater
max:en_allen	eV	The max electronegativity according to Allen
sum:first_ion_en	eV	The sum of the first ionization energies
min:vdw_radius_alvarez	pm	The min vdw radius by Alvarez
ave:covalent_radius_pyykko	pm	The mean single bond covalent radius
ave:period	—	The mean period
ave:melting_point	K	The mean melting point
min:vdw_radius	pm	the min vdw radius
max:mendeleviev_number	—	The max Mendeleev number
min:covalent_radius_pyykko_triple	pm	The min triple bond covalent radius
min:covalent_radius_pyykko_double	pm	The min double bond covalent radius
max:first_ion_en	eV	The max first ionization energy
min:atomic_radius_rahm	pm	The min atomic radius by Rahm et al
sum:specific_heat	J / g / K	The sum of the specific heats at 293 K
min:vdw_radius_mm3	pm	The min vdw radius according to the MM3 FF
sum:en_allen	eV	The sum of the Allen electronegativities
max:en_pauling	EN	The max of the Pauling electronegativities

Table S4: Results feature pre-screening for SISSO for the compositional and structural descriptors defined in Section 2.2, for the 2D material exfoliation energy prediction problem. Also shown is the assumption we made for units of the XenonPy descriptors (see Section 2.4.4 for more details).

Feature	Units Assumption	Description
min:bulk_modulus	GPa	The min bulk modulus
min:thermal_conductivity	W / m / K	The min thermal conductivity
decomposition_energy	eV/atom	The energy of decomposition of the material
min:evaporation_heat	KJ / mol	The min evaporation heat
min:boiling_point	K	The min boiling point
ave:electron_affinity	eV	The mean electron affinity
ave:thermal_conductivity	W / m / K	The mean thermal conductivity
max:electron_affinity	eV	The max electron affinity
var:heat_capacity_molar	J / mol / K	The var in the heat capacity
min:melting_point	K	The min melting point
min:heat_of_formation	KJ / mol	The minimum heat of Formation
max:heat_capacity_molar	J / mol / K	The max heat molar capacity
ave:icsd_volume	cm ³ / mol	The mean ICSD volume
sum:thermal_conductivity	W / m / K	The sum of all of the thermal conductivities
bandgap	eV	The band gap
min:vdw_radius_uff	pm	The min vdw radius from the UFF
max:thermal_conductivity	W / m / K	The max thermal conductivity
ave:heat_capacity_molar	J / mol / K	The mean molar heat capacity
var:electron_affinity	eV	The varaince in electron affinity
ave:vdw_radius_uff	pm	The mean vdw radius from the UFF
perimeter_area_ratio	Å ⁻¹	The ratio between the perimeter and area of the surface
ave:bulk_modulus	GPa	The mean bulk modulus
ave:atomic_volume	cm ³ /mol	The mean atomic volume
max:icsd_volume	cm ³ /mol	The max ICSD atomic volume
ave:evaporation_heat	KJ / mol	The avearge evaporation heat
ave:mendeleev_number	—	The mean Mendeleev number
ave:hhi_r	—	The mean HHI reserve values
ave:num_p_valence	—	The mean number of valance p electrons
max:num_p_valence	—	The max number of valance p electrons
var:num_p_valence	—	The var of number of valance p electrons

S.5 Exfoliation Energy Comparison

To facilitate comparison at more-relevant exfoliation energies, we zoomed in the plot featured in Figure 3 to range up to 2 eV. We present the full parity plot in this section, in Figure S6.

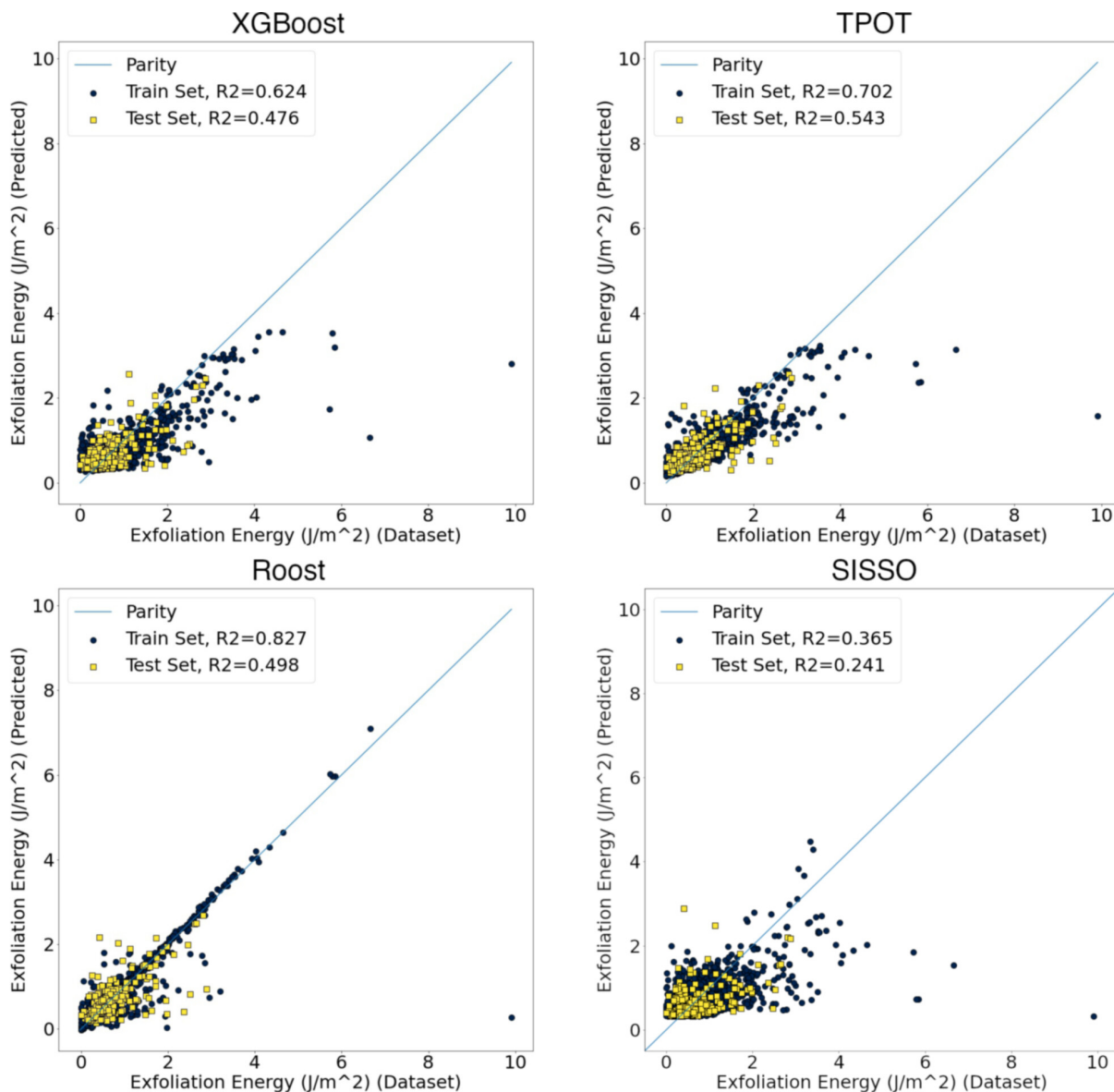


Figure S6: Parity plots for the XGBoost, TPOT, Roost, and SISSO models on the 2D material exfoliation energy problem. Included are the training and testing sets. A diagonal line indicating parity is drawn as a guide to the eye. The SISSO model we report is the rung 2, 2-term model. Regression statistics for the models shown on this plot can be found in Table 4.

S.6 Metal Classification with XGBoost

S.6.1 Structural Fingerprints

For the purposes of 2D material metallicity classification, we leveraged the the Sine Matrix^[112] fingerprint as implemented by the DDescribe package^[113]. The intuition of this fingerprint is that, although it has little to no physical interpretation, it is periodic unlike the Coulomb Matrix^[114], and captures the infinite energy resulting from two atoms overlapping^[113]. The descriptor creates an $N \times N$ matrix, where N is the number of atoms in the cell. We then take the eigenspectrum (i.e. the eigenvalues sorted from largest to smallest) of this matrix to generate a feature vector of length N . Because the 2D structures considered for the metal classifier have between 1 and 40 atoms within the unit cell, we 0-pad the eigenspectrum such that it is always a vector of length 40.

Like the other models in this study, 10% of the data was held out as a testing set. While performing the train/test split, we stratify the data to ensure the same proportion of metal / nonmetal data is in the training and testing set. Finally, in order to we leveraged K-Means Synthetic Minority Oversampling Technique (SMOTE) as implemented by Imbalanced-Learn^[115], with K set to 4 neighbors.

S.6.2 Classification Methodology

The Sine Matrix Eigenspectrum described in Section S.6.1 was used as the input feature to predict whether a material was a metal or nonmetal.

Early stopping was more-aggressively applied, such that if 5 rounds passed without improvement in the Receiver Operating Characteristic (ROC) Area Under the Curve (AUC) of the validation set, XGBoost halted further training. Multiple evaluation metrics were considered in the training process. XGBoost was set to internally optimize the logistic regression score of the model, and Optuna was set to optimize the F1 score of the model, as implemented in SciKit-Learn^[48].

S.6.3 Regression Methodology

An XGBoost model was trained using the metal/nonmetal predictions of the Metallicity Classifier described in section S.6.2. The same training and testing set used for the metal classifier was applied here as well, with all systems predicted by the classifier to be metals removed. Features used included the XenonPy descriptors described in section 2.2.1 and the structural descriptors from section 2.2.2. Because the data being fed into this model was roughly Poisson-distributed, we applied Poisson-based error metrics. XGBoost was set to optimize the squared error, and Optuna was set to minimize the mean Poisson deviance. The TPE sampler with Hyperband pruning was again applied here. For XGBoost predictions on the perovskite volume, we used only the compositional features from XenonPy.

Confusion matrices to summarize the model’s performance on the training and test sets can be found in Table S5 and Table S6 respectively.

Table S5: Confusion matrix for the XGBoost Metal Classification training set

		Predicted	
		Nonmetal	Metal
Actual	Nonmetal	2550	333
	Metal	353	2479

Table S6: Confusion matrix for the XGBoost Metal Classification test set

		Predicted	
		Nonmetal	Metal
Actual	Nonmetal	248	73
	Metal	92	223

Metrics summarizing model performance can be found in Table S7.

Table S7: Model performance metrics for the XGBoost Metal Classifier

Metric	Training Set	Test Set
True Positive Rate (TPR)	0.875	0.708
False Positive Rate (FPR)	0.116	0.227
Accuracy	0.880	0.741
F1 Score	0.878	0.730
ROC AUC	0.947	0.821

The ROC curves of the classifier can be found in Figure S7 and Figure S8.

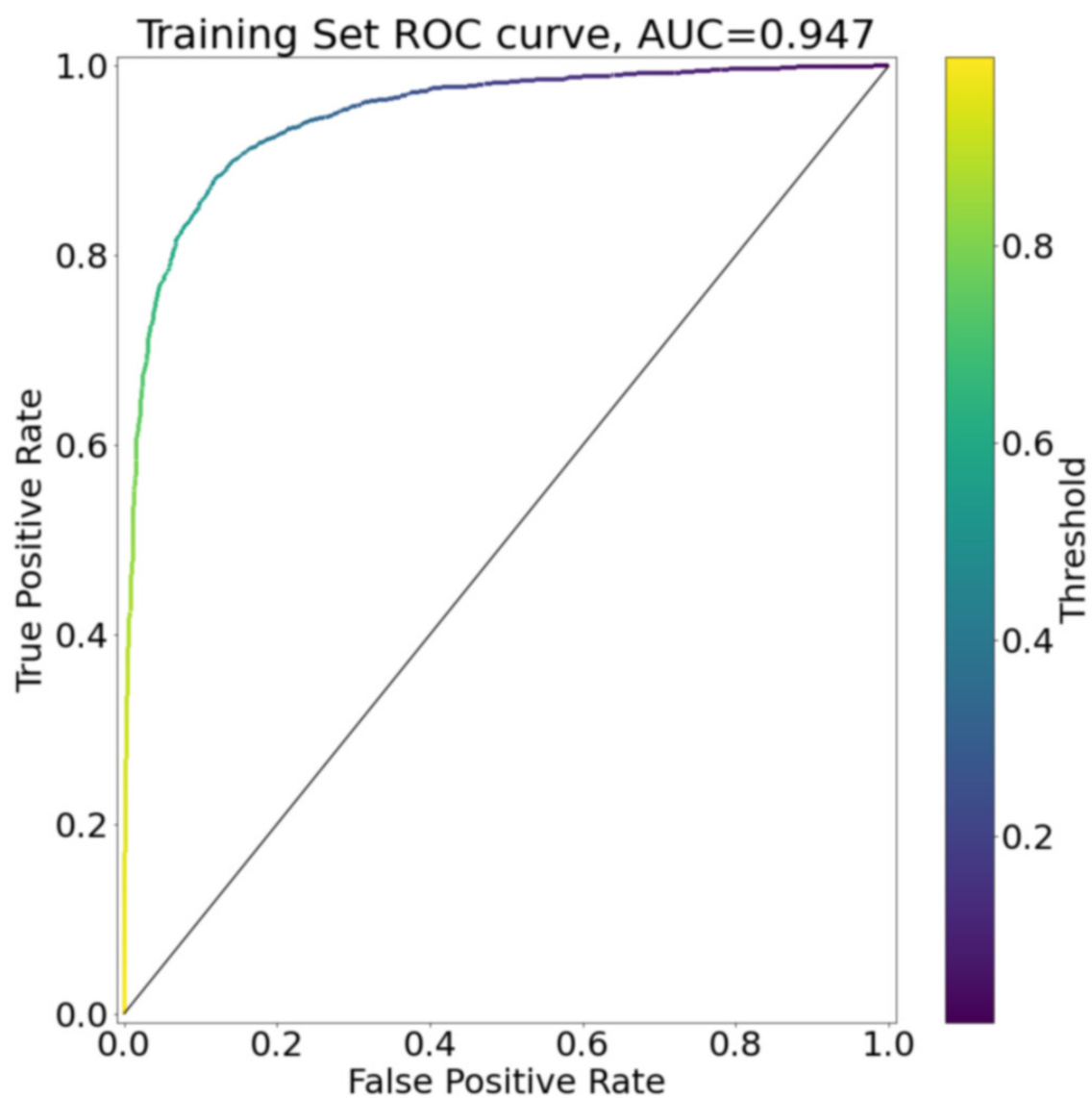


Figure S7: Training-Set ROC curve for XGBoost model classifying whether a material is metallic or not.

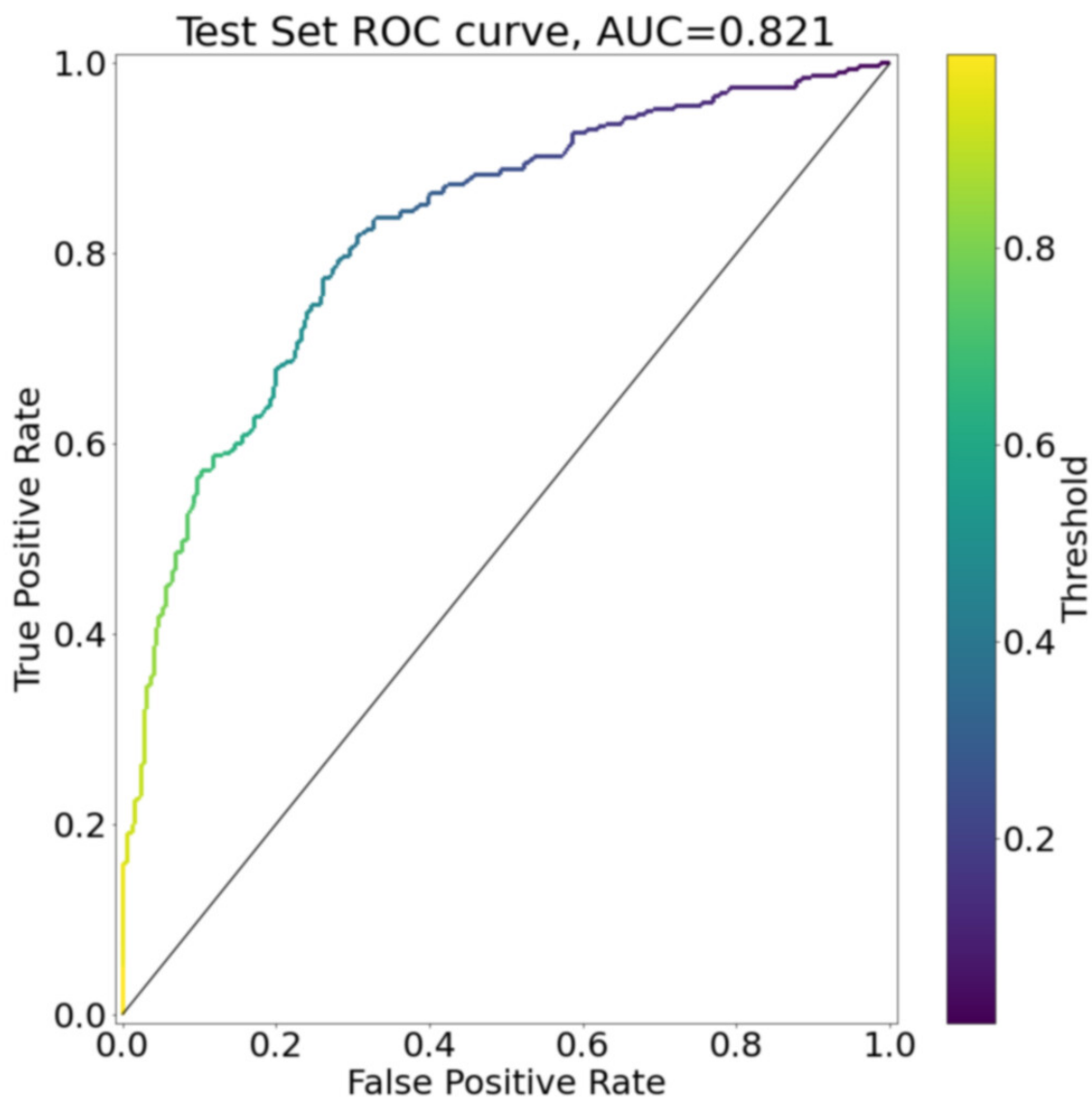


Figure S8: Test-Set ROC curve for XGBoost model classifying whether a material is metallic or not.

Finally, once we had developed an automated way to separate the dataset into metals and nonmetals, we trained an XGBoost regression model on the systems predicted to be nonmetals by the classifier. Statistics for model performance can be found in Table S8, and a parity plot for model performance can be found in Figure S9.

Table S8: Error metrics for the XGBoost regression model trained on systems predicted to be nonmetals by the XGBoost classifier.

Error Metric	Training Set	Test-Set
MAE	0.376	0.645
RMSE	0.530	0.840
Max Error	5.893	3.259
R ²	0.878	0.698

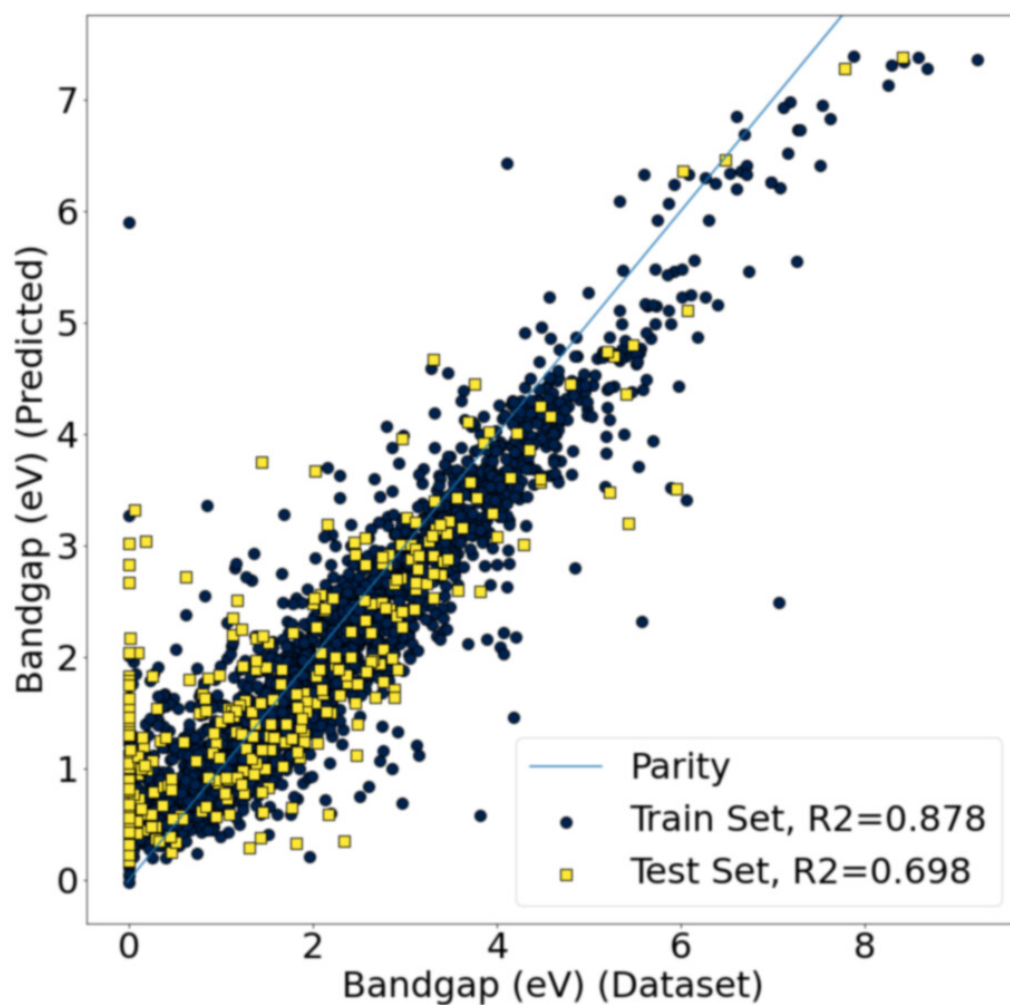


Figure S9: Parity plot showing XGBoost regression model performance on the prediction of 2D material bandgaps.

We can also leverage the regression model to determine which features are the most important to predicting bandgap (Figure S10).

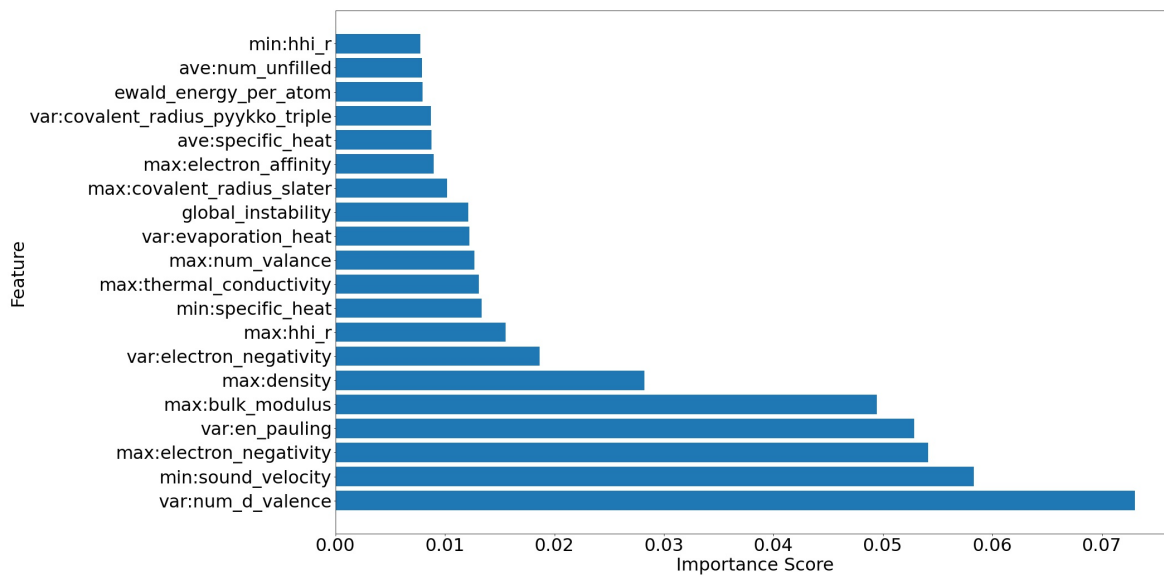


Figure S10: Importance metrics for the XGBoost bandgap regression model on the nonmetal systems identified by the XGBoost Metal Classifier.