Machine Learning and Energy Minimization Approaches for Crystal Structure Predictions: A Review and New Horizons

Jake Graser[†], Steven K. Kauwe[†], and Taylor D. Sparks^{†*}

[†]Department of Material Science and Engineering, University of Utah, Salt Lake City, Utah 84112, United States KEYWORDS Review, Machine Learning, Crystal Structure, Density Functional Theory, Genetic Algorithm, Simulated Annealing, Random Forest, Support Vector Machine, Artificial Neural Network, Confusion Matrix

ABSTRACT: Predicting crystal structure has always been a challenging problem for physical sciences. Recently, computational methods have been built to predict crystal structure with success but have been limited in scope and computational time. In this paper, we review computational methods such as density functional theory and machine learning methods used to predict crystal structure. We also explored the breadth versus accuracy of building a model to predict across any crystal structure using machine learning. We extracted 24,913 unique chemical formulae existing between 290 K and 310 K from the Pearson Crystal Database. Of these 24,913 formulae, there exists 10,711 unique crystal structures referred to as entry prototypes. Common entries might have hundreds of chemical compositions while the vast majority of entry prototypes are represented by fewer than ten unique compositions. To include all data in our predictions, entry prototypes that lacked a minimum number of representatives were relabeled as 'Other'. By selecting the minimum numbers to be 150, 100, 70, 40, 20, and 10, we explored how limiting class sizes affected performance. Using each minimum number to reorganize the data, we looked at the classification performance metrics: accuracy, precision, and recall. Accuracy ranged from 97±2% to 85±2%, average precision ranged from 86±2% to 79±2%, while average recall ranged from 73±2% to 54±2% for minimum-class representatives from 150 to 10, respectively.

I INTRODUCTION

2 Scientific exploration into chemical whitespace has always 3 been a challenging process due to the high risk, high reward 4 nature of research into untested territory. Materials discovery 5 and characterization is a very time intensive process. Synthesis 6 of untested materials requires a large amount of trial and error 7 to determine optimum synthesis conditions with some chemical

- reactions taking days to weeks to perform. Many of these untested materials use exotic elements or compounds which can
- 10 be expensive. In addition to the cost of reagents, samples must
- 11 then be characterized for crystal structure and microstructure.
- 12 Techniques such as diffraction, spectroscopy, and electron mi-
- 13 croscopy can be very time intensive.

28

29

30

Once a material is finally synthesized and characterized, its 14 properties can be evaluated in the engineering design process. 15 However, most applications require an optimization of multiple 16 17 properties which may be interrelated. If we look at the field of thermoelectrics, for example, materials are compared to one an-18 other using a figure of merit, $zT = \sigma S^2 \kappa^{-1}T$, where S is the 19 Seebeck coefficient, σ is the electrical conductivity, κ is the 21 thermal conductivity, and T is temperature. The material prop-22 erties σ , κ , and S are all interrelated. For example, electrical 23 conductivity requires high carrier concentration whereas See-24 beck coefficient requires low carrier concentration to increase 25 zT. In addition, thermal conductivity also increases with carrier concentration which in turn decreases zT. Therefore, optimization of thermoelectric materials requires a compromise between 27

31 The need to discover new materials is not unique to the field of

a better intrinsic balance in these properties.

these properties. Some of the most significant advances in this

field have come from identifying new compounds which exhibit

- 32 thermoelectrics. Similar challenges are seen across many mate-
- rial science fields such as superconductivity^{1, 2}, lithium ion batteries^{3, 4}, solid oxide fuel cells ^{5, 6}, catalysts⁷, high strength ma-
- steries, solid oxide fuel cells s, s, catalysts', high strength materials, and others. In these fields, a relatively small number
- 36 of materials are being actively investigated compared to the tens
- 37 of thousands of known potential compounds in databases such

as the Inorganic Crystal Structure Database (ICSD). In many instances, some of the most exciting and promising new mate-

40 rials have been discovered via fortuity and luck. Critical engi-

neering materials such as vulcanized rubber¹⁰, Teflon¹¹, and synthetic plastics¹² to everyday luxuries such artificial dyes¹³,

43 super glue¹⁴, and synthetic sweeteners¹⁵ were all discovered

44 though chance.

This current approach to materials discovery and deployment is far too slow and expensive to meet the demands that we face in the 21st century. Instead, we need a rational and structured method to explore chemical whitespace. This new method not only needs to be economical but quick, precise, and accurate as well.

Consider The National Academy of Engineering's Grand Challenges. These include such challenges as making solar energy economical, providing access to clean water, or developing carbon capture and sequestration methods, among others. 16 Solutions to these challenges will undoubtedly require radically improved materials to be developed as quickly as possible. To this end, the President of the United States implemented The Materials Genome Initiative (MGI) in 2011 in order to deploy new materials "twice as fast at a fraction of the cost." The MGI proposes to achieve this goal by enhancing collaboration between experimental and computational materials scientists. Computational resources can screen and reduce the total number of experiments necessary rather than experimentally testing every composition. Although the MGI has only been in existence for a short time, we are already seeing key successes from techniques rooted in MGI principles. For example, the Ford Motor Company has employed an MGI-based approach known as Integrated Computational Materials Engineering (ICME) to reduce the time for deploying engine aluminum casting, saving them a hundred million USD18. Other examples include QuesTek Innovations using ICME to develop new aviation components such as high strength steel for landing struts or helicopter rotors¹⁹; GE Aviation developing gas turbine compo-

nents without rhenium to reduce cost²⁰; Ford and General Mo-133 tors researching materials to improve powertrain castings²¹; and 134 75 investigation into distortions caused by welding in ship build-76 77 Critical to most MGI-based techniques is knowing the crystal 137 78 structure of a candidate material *a priori* and then using this 139 structure to calculate performance for a given property. Indeed, 140 79 80 understanding the specific relationships between crystal struc-81 ture, processing, and materials properties such as electrical, op-142 tical, mechanical performance is at the heart of the materials 143 science discipline. However, predicting crystal structure itself 144 for any given composition has been a surprisingly vexing chal-145 lenge for materials scientists, chemists, and physicists for over 146 86 87 a century. 147 While some general rules have been identified which offer in-88 while some general rules have been identified which offer insight, such as Pauling's five rules for crystal structures²³, there
150
are numerous exceptions to these rules and predicting the struc150
ture of some simple and most complex compounds still chal151
lenges scientists today²⁴. Pauling's rules are as follows: 1) The
152
radii ratio and radius sum rule for polyhedra formation, 2) The
154 89 90 91 92 93

electrostatic valence principle for electroneutrality related to the 94 coordination number of the cation, 3) The stability of the crystal 155 95 related to polyhedra sharing of corners, edges, and faces, 4) The lask of sharing of polyhedra when multiple actions with lense 157 96 lack of sharing of polyhedra when multiple cations with large 15% valence and small coordination number are present in the crys-97 98 tal, and 5) The multiplicity of constituents in the crystal will be 159 99 small. Later, in the early 1980's, Pierre Villars built multiple 160 100 three-dimensional stability diagrams by the determination of 161 three specific atomic properties to help separate binary and ter-162 101 102 nary alloys. By using the difference of Zunger's pseudopoten-103 tial radii sums, a difference in Martynov-Bastsanov electroneg-104 ativity, and the sum of the valence electrons, Villars was able to 165 105 build a predictive model to predict thousands of binary and ter166
nary compounds²⁵⁻²⁸. Modern day prediction techniques now
167
rely heavily on computational materials science and take many
168 106 107 108

In this paper, we have done a literature review of multiple algo- $\frac{171}{172}$ 111 rithms with an overview of the basics of each algorithm, a brief. 112 focus on the history, key breakthroughs, and modern examples 174 of crystal structure prediction. We will then discuss new ap-113

109

110

sity functional theory29,30

forms such as simulated annealing, genetic algorithms, and den-

proaches for crystal structure predictions based on machine, 77

learning. This paper will give an overview of the promise and 117 challenges in using machine learning to predict structures. 178

FOR^{179} 1.1 ENERGY BASED **ALGORITHMS** 118 PREDICTING CRYSTAL STRUCTURE 119

Density functional theory (DFT), simulated annealing, and ge-181 120 netic algorithms all require a crystal structure suggestion or a182 121 randomly generated atomic configuration as a starting point to 183 122 begin calculations from first principles²⁹. The algorithms then 184 123 124 search for the lowest energy structure using energetic potentials 185 unique to each algorithm²⁹. The lowest energy states are as-186 126 sumed to be the ground energy states and thus a compound's 187 most likely thermodynamically stable crystal structure. It's im-188 127 portant to note that due to these algorithms focus on energy min-189 128 129 imization, only ground state structures can be calculated. These 190 algorithms can't be used to determine metastable states or struc-191 tures that require external temperature or pressure to remain sta-192 132 ble.

1.1.1 Density Functional Theory

Density functional theory is the most well-known predictive algorithm currently used by material scientists and researchers. Density functional calculations investigate the electronic structure of many-body systems at the ground state. Rather than simulating the interaction between every subatomic particle it uses approximations. These approximations are nucleonic potentials for atoms and use an electron density rather than calculating each individual electron interaction. DFT achieves relatively high accuracy though the quantum mechanical modeling of the spatially dependent electron density in a system. Modern DFT is based on non-interacting electrons moving in a system-wide electronic potential. The potential is constructed using the structure and elemental composition of the system with their interelectronic interactions. The potential is evaluated to determine the energy cost for each state, or configuration, of the system. The correct ground state electron density is determined when the energy of the system has been minimized. DFT requires knowing a candidate crystal structure before making a calculation. Therefore, researchers will create a list of possible structures and then use DFT quantum calculations to determine each structure's energy at zero kelvin31. The lowest energy is then determined to be the most stable, and thus most likely structure³¹. Therefore, this is a zero-kelvin approximation and does not work for high temperatures which can include room temperature³². Researchers have merged different techniques, such as molecular dynamics³³, to overcome this issue.

The pioneer of DFT was Douglas Rayner Hartree when he created a self-consistent field for electrons to solve for the wave function using the field around the nucleus³⁰. This was expanded by two of his graduate students, Fock and Slater, in 1930 by replacing the equation with a determinant³⁰. This became the Hartree-Fock equation and Slater determinants. Earlier, in 1927, Thomas and Fermi developed a model to calculate atomic properties³⁰. This used a local density (LD) approximation for kinetic energy³⁰. This set the foundation for solving the wave equation for atoms using density functionals. In 1964, Hohenberg and Kohn introduced a variational principle for energy, which showed a relationship between the ground state density of electrons and the wave function30. This was a huge accomplishment and the start of modern DFT. Formalism and refinement of the densities helped refine the algorithm. Yet, acceptance of this algorithm didn't occur until around 1990 due to wariness within the field of chemistry³⁰. The creation of simple to use software packages greatly improved DFT acceptance and use³⁰. Between 1990 and 2015 there has been over 160,000 publications in the field of chemistry alone³⁰.

Yet, DFT is not without its shortcomings. Computational costs remain a large issue for each DFT calculation while numerous tests must be run to determine the proper energy functional as different approximations will affect the outcome³⁰. DFT also struggles with highly correlated electron systems, large scale systems, modeling weak or Van Der Waals forces, time-dependent dynamics and properties that are not observable at the ground state such as excited states or room temperature bandgap energy³². Critically, DFT also cannot calculate disordered structures necessary for many unique material properties (partial cation occupancy, oxygen vacancies, etc). Recently, better programmed algorithms and approximations coupled with the introduction of machine learning has helped offset the computational cost and time requirement limitations^{32, 34}.

Density functional theory has had many successes in material253 194 property predictions ranging from batteries^{35, 36}, capacitors³⁷,254 195

thermoelectrics³⁸⁻⁴⁰, superconductors^{41, 42}, photovoltaics^{37, 43},255 196

197 and catalysts⁴⁴. DFT calculations continue to improve accuracy256

and upwards of 15,000 density functional theory papers are257 198

published every year⁴⁵. 199

1.1.2 FORCE FIELD MODELS

200

Another method to calculate the energy of atomic configura-261 201 tions are empirical force field models. These simulate intera-262 202 tomic interactions in unique ways depending on the model^{46, 47}, 263 203

These models are used when the system grows in complexity²⁶⁴ 204 where ab initio calculations become too computationally de-265 205

manding 48. These models are considered critical for nanoparti-266 206 207

cle structure predictions to reduce computational time due to the 267

large size of each predicted system46, 47. 208

However, the accuracy of the calculation is heavily dependent 269 209 on which model is used as well as the accuracy of the model's²⁷⁰ 210

approximations⁴⁸. Swamy et al.⁴⁸ used two separate force field²⁷¹ 211

models to predict all known polymorphs of TiO₂ with varying²⁷² 212 success depending on the specific polymorph. It was shown that 273 213

one model could predict low pressure polymorphs accurately²⁷⁴ 214

but struggled with high pressure polymorphs while the other 275

model could predict high pressure polymorphs accurately but²⁷⁶ 216 277 had difficulty with low pressure polymorphs.

OPTIMIZATION²⁷⁹ **ENERGY GLOBAL** 218 1.1.3

278

219 **ALGORITHMS**

Simulated annealing is a computational approach that is based₂₈₂ 220 on the process of physical annealing wherein a material is 283 221 heated to modify its crystal or microstructure. Modifying the 284 222 crystal and microstructure requires atomic mobility by over-285 223 coming energy barriers. This is made possible by increasing the 286 224 temperature. These atoms then settle into their lowest energy_{2,87} 225 state in the crystal at the given temperature. This allows the at-288 226 227 oms to move and adjust the crystal structure to reach the mini-289 mization of the Gibbs function assuming constant temperature 290 228 and volume. In simulated annealing, the same concept is ap-291 229 plied. We allow these atoms to settle into their lowest energy₂₉₂ 230 state as their simulated energy, commonly labeled temperature, 293 231 is slowly decreased in the model. It is important to note that this 294 232 temperature is an arbitrary energy unit and not related to actual 295 233 temperatures⁴⁹. Minimum energy is found by randomizing the₂₉₆ 234 motion of simulated atoms using either Monte Carlo statistics, 297 235 236 molecular dynamics, or other techniques²⁹. The initial tempera-298 237 ture is selected so that the kinetic energy of each atom is high 200

so global minima can be found²⁹. Simulated annealing has had success in predicting inorganic structures and can make predic-301 240

enough to allow the system to overcome local energy barriers 300

tions of partially disordered materials- something DFT cannot₃₀₂ 241 242

First introduced by Kirkpatrick in 1983, a relation of physical 305 243

annealing and statistical mechanic relations lead him to intro-306 244 duce the algorithm⁵¹. By using the algorithm he showed pre-307 245

dictions for the classic traveling salesman problem to the phys-308 246

ical determination of wiring and cooling within a computer 51.309 247

The algorithm was quickly adopted due to its robust nature and 310 248 ease of use^{29, 52}. The algorithm was expanded upon with new₃₁₁ 249

250 techniques to increase the accuracy of the algorithm such as au-

tomated assembly of secondary building units and a hybrid312 251

252 method using Monte Carlo basin hopping²⁹.

However, simulated annealing is computationally intensive, which led to many researchers developing workarounds such as parallelization processing techniques⁵². The most concerning problems with simulated annealing are: slow energy landscape exploration, the inability to focus on specific problems of interest due to the randomization of atomic placement, and computational limitations^{29, 53}. In theory, simulated annealing can explore the entire energy landscape and find the energy minimum regardless of starting position but the time and computational resources required for this complete exploration can be excessive⁵⁴. The algorithm must start in a specific point and it could miss low energy regions as the algorithm reduces temperature. To offset this limitation, many simulated annealing runs are done sequentially and different starting spots are selected to better explore the energy landscape²⁹.

Genetic algorithms are another energy-based technique used by researchers to predict crystal structures. Created by John Holland in 1975, genetic algorithms are a subset of evolutionary algorithms⁵⁹. These algorithms are based on the concept of evolution where the strong can procreate offspring while the weaker will not. An initial population of structures is generated with constrained but randomized atomic placement, or prearranged atomic configurations⁵⁵. The algorithm then selects parents from the population to exchange crystallographic information. This is done using a random selection code, such as tournaments or a roulette wheel, where the percentage chance of selection is dependent on how well it meets certain criteria specified by the algorithm, referred to as a fitness function²⁹. The parent structures share information in either a random pattern or by mixingtogether, which results in an offspring structure⁵⁵. This process repeats until the desired number of offspring is achieved. Mutations can be introduced to add diversity to the population by changing random properties of the crystal structures. Mutations and offspring then have their individual energies minimized and are added to a new dataset called a generation. This new dataset also includes the samples of the dataset that are lower in energy than the new offspring and mutants. The old generation is replaced and the process is continued until the energy converges to some final criteria or by reducing the training set at each generation until only a single crystallographic structure remains^{29, 55}. The practice of applying genetic algorithms to material science has been growing rapidly in the last few years⁵⁵⁻⁵⁸. After their initial introduction, multiple variations of genetic algorithms have been introduced such as adaptive genetic algorithm58. Specific examples for crystal structure predictions exists such as global space-group optimization known as GSGO57, and the genetic algorithm for structure and phase prediction also known as GASP⁵⁶.

A popular technique for the creation of offspring structures is the 'cut and splice' method⁶¹. This method creates a new generation by splitting a chosen structure at an arbitrary plane. Members of the generation can be sliced and combined to form a new randomized structure, or offspring. The cut and splice method can also be used to generate mutations by rotating a section of the given structure at an arbitrary angle. These new structures have their energy minimized and are added to the new generation⁶¹. Repeating this process allows us to optimize the crystal structure with each new generation eventually reaching a minimum.

Like all optimization techniques, the selection of algorithm parameters will affect performance. For example, convergence to a global minimum can be discovered if the initial population size is large enough, the creation of the offspring population is 316 set at an appropriate rate to explore space but not to over satu-377 rate a local minimum, and a mutation rate is significant enough 378 317

to shift the algorithm out of local minimums. However, this 379 318

leads to large computational times. Thus, a promising region380 319

320 should be determined with a small population size to reduce381

computational time, if selected incorrectly, this can lead to find-382

321

322 ing only a local minimum⁶⁰.

323

There are many simulation software packages available for 385

these algorithms mentioned above. The most common software 386 324 325 package used for DFT in crystal structure predictions is the Vi-387

enna Ab Initio Simulation Package, or VASP for short. Histor-388 326

ically, it has had success in crystal structure predictions^{31, 34, 56,}389 327

58, 62, 63. As mentioned above, DFT needs a starting structure to 390 328 calculate energy. To generate this starting structure for DFT391 329

there are multiple methods and algorithms available. USPEX392 330

(Universal Structure Predictor: Evolutionary crystallography)393 331 has been used to determine high pressure phases of CaCO₃⁶⁴,394 332

333 while CALYPSO (Crystal Structure Analysis by Particle395

Swarm Optimization) is another software package that recently 396 334

335 has been used to determine structures of noble gases at high397

336 pressure⁶³, and finally the AIRSS (Ab Initio Random Structure

337 Searching) method which has been used for determining the 398

338 crystal structure for lithium based crystal structures^{65, 66}.

1.2 MACHINE LEARNING ALGORITHMS FOR401 339 PREDICTING CRYSTAL STRUCTURE 340

400

All the previous techniques involve calculating and comparing $\frac{403}{100}$ 341

energies for crystal structures to make predictions about which

crystal structures are most thermodynamically stable. This sta-343

344 bility is with respect to certain given conditions, defined by ei-345 ther the algorithm or user, such as constant entropy or constant

volume. A fundamentally different approach exists which re-346

lies on machine learning. Machine learning is a heavily data-347

348 centric technique where large amounts of data are collected and 349 analyzed. A predictive model is then trained on this large col-

350 lection of data. This model can then be fed inputs similar to the

351 collected data to predict probabilistic results. These inputs can

352 be categorical as well as numerical. Thus, all supervised learn-353

ing algorithms can be characterized as classification or regres-354 sion problems. This ability to segregate crystal structure data of

355 all types is key for allowing us to predict crystal structure.

Companies such as Amazon and Netflix already collect an enor-356

357 mous amount of data related to consumer interest, browsing

habits, viewing history etc. and are already incorporating ma-358 359 chine learning into their websites as highly efficient ways to

360 recommend products or entertainment options to consumers.

361 These algorithms do not need to know exactly why the con-

sumer is interested, it only needs to predict the probabilistic 362

363 likelihood of a consumer being interested. The mechanistic details of the relationships, which are essential to a technique such 406 364

as DFT, are not even necessarily known in machine learning₄₀₇ 365

algorithms. Instead, only the probabilistic determination of a₄₀₈ 366

given outcome from input data matters. Yet, these algorithms₄₀₉ 367

should be viewed as tools to assist rather then to replace₄₁₀ 368 experimental and computational materials scientists.411

369 Ultimately, the algorithm will only make predictions and some 412 370

of these could be correspond to completely unstable or even₄₁₃ 371

physically impossible compounds. Therefore, chemical₄₁₄ intuition will still need to be utilized to determine what is415

374 valuable and what to ignore.

Machine learning algorithms rely on building predictive models₄₁₇

from empirical data or calculated data^{31, 67}. The data used for 376

supervised machine learning is organized in tables referred to as training datasets. Each row describes a single entity or observation, and each column represents a commonly shared feature or attribute. We label these columns as features. These commonly shared features include a column which contains the key property you are attempting to predict, such as crystal structure. The column features can be numerical or categorical. A sample of features used to predict crystal structure could include compositional thresholds, bond character, or average number of valence electrons among many others. For data to be useful in machine learning, each row needs a value for the features that will be included in the model. When features are missing too many values, they are generally discarded, although there are methods to estimate the missing values. For example, imputation is the procedure of replacing empty values in a data set⁶⁸. Imputation is typically handled by filling empty cells with the mean for continuous numerical data, the mode for categorical number data, or the most common string for written categorical data. There are also more sophisticated procedures which involve building nested predictive models to fill in the missing attributes⁶⁸.

Like the energy-based algorithms mentioned above, multiple machine learning algorithms exist such as random forest algorithms, support vector machines, and artificial neural networks. They all share the ability to use a collection of data to build a predictive algorithm. Each of these algorithms build prediction models in different ways and the data requirements, such as size and formatting, differ per algorithm.

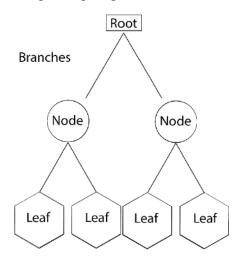


Figure 1: Graphical description of a decision tree with the root, node, and leaf sections labeled. The data is passed from one section to another along branches.

A popular machine learning technique used for predictive model building is the random forest algorithm⁶⁹. The random forest algorithm utilizes many independent decision trees trained from collected data. Training is the process of using the input data to create a criteria-based prediction model that has predictive power. A decision tree is trained by using a subset of features from the data. The training process compares feature values for all inputs and attempts to segregate the input data. The features separate the input data at different feature values creating successively more homogeneous groups⁶⁹.

As discussed above, there are two types of commonly used decision trees: Classification, and Regression⁷⁰. Classification de418 cision trees create predictions that attempt to classify categori-481 419 cal data, an example being crystal structure such as fluorite, spi-482 nel, etc. Regression decision trees predict continuous numerical 420 outputs, such as thermal conductivity. In the random forest al-421 422 gorithm, all the trees in the "forest" have different structure because they sample different data and random features⁷⁰. The 423 trees are composed of unique nodes and branches. The nodes 424 425 are a way to represent splitting points in the data. The initial node is referred to the root of the tree. Feature values from the data are used to separate an input data group. The groups that 427 428 result from separation are called branches. Each subsequent node receives an input group from the branch above it. That 429 separation is output to nodes below it until all the groups are 430 431 homogeneous. These final nodes are referred to as leaves. The tree structures that are built to separate the experimental data 432 433 can then be used as a model for separating future data⁶⁹, an ex-434 ample is shown in Figure 1.

Random forest has already been used as a high-throughput ma-435 terial screening process for thermal conductivity or energeti-436 cally favorable compositions^{71, 72}. For example, Oliynyk et al. 73 437 built a model that predicted whether 21 ternary compositions 438 were either full-Heusler, inverse Heusler, or not Heusler. Of 439 these 21, 19 were confirmed though experimental results. The 440 challenge with differentiating these classes is that they all look $_{483}$ 441 nearly indistinguishable via powder diffraction and single crys-484 442 tals are difficult to grow and therefore rarely used to character-485 443 ize structure. Even with these difficulties, the algorithm had an₄₈₆ 444 accuracy of 94%. Balachandran et al. also used decision trees₄₈₇ 445 as well as support vector machines, which we will discuss be-488 446 low, to predict wide band gap binary structures as well as tran-447 sition metal intermetallic compounds. These algorithms had ac-489 448 curacies ranging from 86.7% to 96.7% for the decision trees and 490 449 86.7% to 93.3% for the support vector machines 74.

A support vector machine (SVM) is another machine learning₄₉₃ 451 algorithm based on the segregation of data. SVMs can segregate494 452 regression or classification data like the random forest model₄₉₅ 453 discussed above. SVMs accomplish this task by plotting the 496 454 data into an n-dimensional space. The algorithm then attempts₄₉₇ 455 to create a hyperplane to segregate the data. This hyperplane is 498 456 457 determined by maximizing the vector normal to the hyperplane, usually labeled W, and the closest data point to create the larg-499 458 459 est gap possible⁷⁵. By graphing the data with respect to different₅₀₀ 460 physical properties, the algorithm can compare the hyperplane501 461 separation with respect to physical properties. The hyperplane 502 462 with the largest split is the defining feature relative to the phys-503 463 ical properties and thus the most important feature to segregate 504 464 the data. To help with this process, an error function can be in-505 troduced to allow the algorithm to ignore a few data points to 506 465 plot the hyperplane. The distance of this separation can be used 466 as a method to optimize the algorithm 76. A graphical example is 507 467 shown in Figure 2. The strength of SVM algorithms is the abil-508 468 469 ity to optimize itself by adjusting its kernel function and thus 509 470 adjusting the dimensionality of the problem to help with segre-510 gation^{77,78}. The kernel function is a symmetric and continuous511 471 472 function were, if the restrictions of Mercer's theorem is met,512 473 can define the dot product in a specific space⁷⁹. This allows the 513 474 program to increase dimensionality without calculating the dot514 475 product continuously, allowing the algorithm to expand its di-515 476 mensionality without impacting computational time. Yet, this 516 477 leads to the major disadvantage of SVMs. The choice of the 517 478 kernel function is a critical and challenging process and unlike 518 479 certain other machine learning techniques, SVMs can still be

computationally demanding depending on the dimensionality of

480

the problem with respect to other machine learning algorithms⁷⁵.

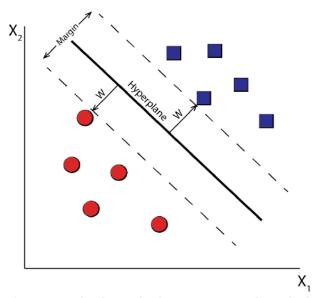


Figure 2: Graphical example of support vector machine. The data is being separated by the hyperplane that maximizes the vector W.

SVMs have been used in the prediction of protein sequences⁸⁰, residue-position importance⁸¹, domain boundaries in protein structures⁸², microstructure imaging⁸³, and atmospheric corrosion of materials⁸⁴. SVMs have also been used in crystal predictions of binary and intermetallic compounds with success⁷⁴.

With regards to predicting crystal structure specifically, Oliynyk *et al.* has built a support vector machine to predict the crystal structure of binary compounds as well as ternary compounds. The binary algorithm achieved an accuracy of 93.2% with a training set of 706 compounds. The authors provided further experimental validation by synthesizing one compound, rhodium cadmium, which was predicted to have the cesium chloride structure, and confirming via X-ray diffraction that the predicted structure was correct⁷⁷. The ternary algorithm achieved an accuracy of 96.9% with 1556 unique compounds⁸⁵.

Artificial neural networks (ANNs) are another machine learning method used in material informatics. ANNs are capable of modeling essentially any complex relationship given enough data⁸⁶. ANN's tend to perform well for large amounts of data, experiencing performance saturation later than other machine learning models. They are particularly capable of dealing with data that must do with spatial and temporal relations, such as images and text processing.

Artificial neural networks are based on a collection of connected units called neurons. Neural networks rely on layering of neurons to allow for processing of complex patterns. Most ANN models consist of an input layer, hidden neuron layers, and an output layer. ANN models work by processing input values though massive connected networks called hidden layers. Each connection in the network is called a synapse, and it transmits information to the neurons downstream. Information is passed starting from the input layer until the output neurons are reached. The neurons derive value from the synapse connections while also converting the data into non-linear space if needed. ANN models are trained by adjusting the weights of

each synapse until the output of the network is close to the out-554 put of the training data⁸⁷. A graphical example is shown in Fig-555 ure 3

520

521

522

523

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

541

542

543

544

545

546

547

548

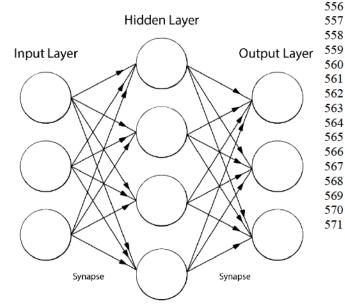


Figure 3: Graphical definitions of Artificial Neural Networks with the input layer, the hidden layer, and the output layer. Each layer connects to each other layer though multiple path called synapse.

As of the writing of this work, neural networks are not used in crystal structure determination. There have been a few preliminary classifications related to structure such as protein secondary structure investigations during folding^{88, 89}. When it comes to inorganic crystal structures, neural networks mainly focus on data interpretation and categorizing. Recently, Timoshenko *et al.* built an artificial neural network to decipher metallic nano-572 particle structures from experimental data⁹⁰. Due to the large 573 requirement of information required for artificial neural net-574 works, they created a dataset of simulations that were verified 575 against experimental data. The accuracy of artificial neural net-576 works allowed them to predict an average coordination number 577 up to the fourth coordination shell, and thus the size and shape 578 of the nanoparticle.

Regardless of the type of machine learning algorithm utilized, $^{580}_{581}$ success is measured by the ability to forecast and predict accu- $^{582}_{582}$ rately. There are many different and unique ways to test the ac- $^{583}_{583}$ curacy of a machine learning algorithm. One of the most com- $^{584}_{585}$ mon and simple methods is the k-fold cross-validation. The idea $^{585}_{585}$ behind k-fold validation is the separation of the training data $^{586}_{585}$ into k equal datasets. The algorithm is trained on k-1 datasets $^{587}_{587}$ and is used to predict the dataset that wasn't used in the training. $^{588}_{587}$ real values are compared to the predicted values. For $^{589}_{589}$ real valued property, root mean squared error is the commonly reported error metric.

Root Mean Squared Error =
$$\sqrt{\frac{\sum_{i=1}^{n} (\hat{y}_i - y_i)^2}{n}}$$
 591
592
593

549 Where \hat{y}_i is the observed value, y_i is the predicted value, and n595550 is the total amount of predictions performed by the algorithm.596 551 A quick way to examine your algorithm is a by comparing the597 552 results to a random guessing algorithm. If the probability of the598 553 algorithm is the same or worse than a random selection of each

class, then rebuilding the algorithms, focusing on features or shifting to a different type of algorithm, is required.

When model predictions are categorical as opposed to real valued, a more useful accuracy evaluation tool is the confusion matrix. The confusion matrix compares the known values from the training set and the predicted values from the algorithm. From this you can compare how often the algorithm classifies the data correctly to determine its accuracy, in-class precision, recall, and false positives/negatives. The overall model accuracy is the ratio of the total number of correct predictions versus the total predictions. In class precision is the accuracy of a specific prediction in the model. In class recall is the number of times a class was predicted correct over the total number of instances of that class. A false positive or negative is when the algorithm is wrong in its prediction. We can define these equations of accuracy, in class recall, and in class precision as seen below. For clarity we will define True Positive as TP, True Negative as TN, False Positive as FP, and False Negative as FN.

$$In Class Precision = \frac{TP}{TP + FP}$$

$$In Class Recall = \frac{TP}{TP + FN}$$

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

To help illustrate the idea of a confusion matrix, let us consider an algorithm that predicts if the crystal structure will be cesium chloride type structure (CsCl) or another category we shall label 'Other' (See Figure 4). For false positives/negatives, we will define that the CsCl type structure outcome is positive, and the 'Other' outcome is negative. Out of 200 compounds, let us have 55 CsCl type structures compounds and 145 'Other' compounds in our training set. Let us assume the algorithm predicts 60 CsCl type structure compounds and 140 'Other' compounds. The accuracy of this algorithm would be the sum of the number of times it guessed correctly over the total training set. That would be 50 correct CsCl type structure compounds plus 135 correct 'Other' compounds divided by the total training set of 200 to give us 92.5% accuracy. The in-class precision for CsCl type structure would be 50 divided by 60 or 83.33%. In-class recall for CsCl type structure would be 50 divided by 55 or 90.91%. CsCl type structure would have a false positive of 10 while 'Other' would have a false negative of 5.

1.3 SYNERGY VS COMPETITION IN ENERGY-BASED VS MACHINE LEARNING APPROACHES

Researchers have started using machine learning techniques to explore chemical whitespace focused on crystal structure with success^{71, 73, 91, 92}. The databases of information online, such as the International Crystal Structure Database or the Pearson's Crystal Database, give large amounts of physical parameters that can be used to build training datasets. These can then be used to build a prediction algorithm. This is very attractive to

researchers for high throughput material exploration. Yet, prob-627 lems still exist within physical sciences with machine learning 628 algorithms because large and diverse training sets are required 629 as well as knowledge of coding and algorithm deployment. The building of training sets requires a large amount of time and ef-631 fort while energy-based algorithms still struggle with calcula-632 tion time and cost. It's not surprising then that most researchers do a combination of each technique to offset the weakness of 634 each type of algorithm. Using machine learning, researchers can 635 search chemical whitespace quickly and single out interesting or promising materials. These are then fed into energy-based functionals or calculations to create a more refined prediction 638 of the material properties and characterisitics^{31, 34, 93-95}. Others₆₃₉ use these energy-based algorithms such as DFT to generate da-640 tasets for unknown or ill-defined chemical compounds to train 641 their machine learning dataset upon 96,97. 642

6		^
v	1	J

616

617

618

619

620

621

622

623

624

625

626

ć

599

600

601

602

603

604

605

607

608

609

610

611

612

613

614

a)		Predicted CsCl	Predicted 'Other'	
	Actual CsCl	<u>True Positive</u> 50	<u>False Negative</u> 5	CsCI Total 55
	Actual 'Other'	<u>False Positive</u> 10	<u>True Negative</u> 135	'Other' Total 145
		Total CsCl Prediction 60	Total 'Other' Prediction 140	

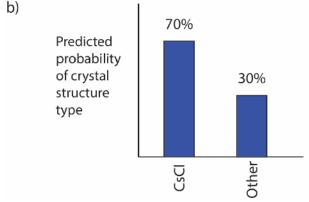


Figure 4: a) Example of a confusion matrix with CsCl defined as positive outcome and 'Other' as a negative outcome. b) example of predicted probability of a specific data point from the algorithm. Due to CsCl type structure having the larger percentage, the algorithm would categorize this data point as CsCl type structure.

As described earlier with Heusler and basic binary structure pre-651 dictions, machine learning has been used for a very select few 652 specific crystal structure predictions. However, a general, uni-653 versal structure type prediction algorithm has never been de-654 ployed using machine learning. Therefore, in this paper we ex-655 tend previous efforts to determine the extent to which machine 656 learning could predict any crystal structure type. We accom-657 plish this by training off all crystal structure data available in 658 Pearson's Crystal Database to predict the structure for any com-659 position.

2. METHODS

643 644

645

648

662

663

In this work, we use a machine learning algorithm from the open source program H2O FLOW. A database of 24,215 unique formulae, and associated entry prototype (EP), Pearson symbol, space group number, phase prototype, etc. was assembled from the original 24,913 entries obtained from the Pearson Crystal Database. This was the result of removing formula with exotic elements such as Polonium, Astatine, Protactinium, etc. These exotic elements lacked sufficient elemental properties for our machine learning method. These were organized into identification numbers in order of decreasing size. A graphical representation of the specific entry to the number of representatives per entry is shown in Figure 5. This was then screened for materials near room temperature (290-310K) with duplicates removed. The chemical formula for each entry was then separated into composition-weighted elemental and atomic properties to allow the model to explore any chemical composition as all features were elementally based. These formulae-based features were then uploaded into H2O FLOW and were then used as a training set in the random forest machine learning model. This model was selected due to its ease of use and scalability for the size of the data. Error metrics were calculated in accordance to the k-fold validation methods discussed above.

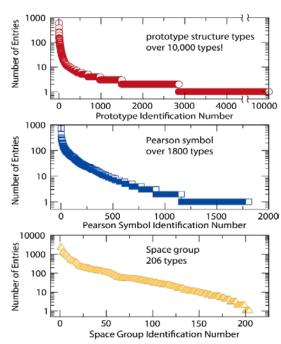


Figure 5: Histogram of entry size versus entry prototype, Pearson's symbol, and space group. The size of each category drops quickly with the majority of each category having only a few entries.

To make a given structure prediction there must be multiple example compositions or instances having that structure type to correlate composition to a structure type. In our dataset, there were 10,711 unique entry prototypes and 97.5% of the entry prototypes had fewer than 10 instances. However, a mere 2.5% of the entry prototypes, those most common structures such as perovskite or spinel, encompassed 28.5% of the data. This led us to question at what point, in terms of number of prototype entries, can we build accurate models. Moreover, since machine learning model accuracy generally increases with number of instances per class type, up to a point, we can study the expected tradeoff between model breadth and model accuracy. Specifically, to handle the uneven distribution of entry prototypes, a

minimum number of instances was set at an arbitrary cutoff. This cutoff was then varied for different models. Entry prototypes with fewer than the required number of instances were categorized into a single class named 'Other'. When the minimum cutoff value was varied between 150 and 10 the 'Other' class encompassed between 92.5% and 64.1% of the input data, respectively. The database was prepared multiple times with separate cutoff values with minimum-bin size of: 150, 100, 70, 40, 20, and 10.

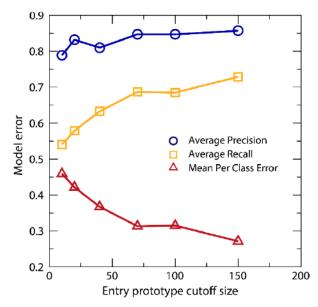


Figure 6: Model error with respect to cutoff size. Each point is a specific cutoff with guidelines inserted between points. As cutoff size increases, the model's overall accuracy increases as expected. Error bars (2%) are smaller than marker size.

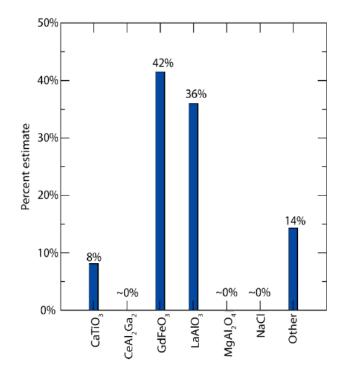


Figure 7: Graphical representation of the algorithms probabilities for entry prototype. With GdFeO₃ having the largest probability, it will be selected as the algorithms entry prototype prediction.

Although 97.5% the entry prototypes exist below the cutoff limit of 10, we still find the classification of 'Other' to be useful information. With a cutoff limit of 10 entry prototypes, a prediction of 'Other' leads to a rare crystal structure with less than 10 known similar compounds listed within the Pearson's Crystal Database. If a researcher is looking for a very rare crystal structure for a specific property, that crystal structure most likely will exist within the 'Other' category. Moreover, the model is able to make accurate predictions with moderate recall for the most common crystal structure types.

The prepared datasets were all analyzed with a distributed random forest algorithm. All the algorithms had a limitation of 150 trees with a maximum depth of 40. Each prepared dataset resulted in a unique model. Error metrics were calculated using 5-fold cross-validation in accordance to the k-fold validation methods discussed above. Predictions were plotted as a confusion matrix.

For error analysis, each cutoff model was built with five different random seeds. The error metric we compared is the mean per class error. Mean class error is defined as one minus recall as seen in Figure 6. The difference between the largest and smallest metric is calculated to determine error range.

3. RESULTS AND DISCUSSION

Before describing model accuracy, we first remark on the model speed. As described in the introduction, one of the key advantages of machine learning is the speed of prediction. In this model, we trained our algorithm on 25,000 different entries with 90 columns of metadata each. Therefore, our overall dataset exceeded 2.2 million data values. Nevertheless, training the model on 25,000 entry prototypes only took up to 2 hours depending on the minimum bin size. Once the model was trained, 15,000 entry prototypes were predicted in under 10 seconds on a laptop

(Intel-I7, 2.6GHz processor 16GB RAM, 64-bit Windows771 708 709 10). These composition-based predictions included the as-772 signed entry prototype (structure with highest probability) and 773 710 711 a breakdown of the probabilities for each other entry proto-774 712 type. The most probable class is selected as the final prediction.775 713 For example, the model output for Cr_{0.12}Ru_{0.88}SrO₃, which has 776 the entry prototype GdFeO₃ type structure, would have a distri-777 714 bution across all possible entry prototypes. GdFeO3 type struc-715 ture has the highest probability so it would be selected as the 778 717 entry prototype. The graphical representation of this data can be 779 seen in Figure 7. For compounds with multiple crystal struc-780 718 tures possible, the model predicts the most common structure at 781 719 720 room temperature due to the model and training set having been 782 721 built at room temperature. 722

To determine the error range of the model, each cutoff model785 was built with a different random seed five separate times. The786 range was determined by the difference in the largest and smallest error. This difference ranged from 0.5% for accuracy, 2% for recall, and 1.8% for precision. To be conservative with the error range, the largest error was adopted for all our percentage errors discussed below.

723

724

725

726

727

728 729

730 As expected, as the minimum cutoff size was reduced for each 731 entry prototype, from 150 to 10, fewer data were available and 732 the mean in class error increased slowly. This mean class error 733 ranged from 27±2% to 46±2% for a minimum-bin size ranging from 150 to 10, respectively. In comparison, random guessing 734 735 mean-class error ranged from 99.8% to 83.3%. If the algorithm 736 only selected the 'Other' category, its mean class error ranged 737 99.9% to 86.7%. We can see that regardless of cutoff, all 738 showed an overall error far lower than random guessing or only selecting 'Other'. Although mean in-class accuracy describes 739 740 the overall performance of the model, the reliability of a prediction is better understood by evaluating the accuracy for predict-742 ing individual entry prototypes. For the 150-cutoff section, the largest class error of the entry prototypes was CaTiO3 type 743 structure at 52% while others are much more accurate such as 744 745 CeAl2Ga2 type structure and MgAl2O4 type structure with classification errors of only 14% and 19%, respectively. To clarify, 746 747 when we discuss classification errors, we are describing the percent of the time the algorithm categorized a prediction in the 749 wrong category. For example, if the algorithm predicts a 750 CaTiO₃ type structure prediction as a GdFeO₃ type structure that 751 would be a classification error. An alternative metric used is 752 precision, or one minus the classification error. Overall the in-753 class error is surprisingly low even when we only include as few 754 as 10 entry prototypes in training data with classes such as 755 BaNiO₃ type structure with only six entries has 100% precision. 756 However, some specific classes with only one or two entry pro-757 totypes predictions have zero precision. In other words, we see evidence that when the model thinks a composition belongs to 758 759 a given class it will predict it with very good precision but in a 760 significant number of cases where only one or two data points 761 exist it will just call it 'Other'.

762 Some entry prototype predictions are more consistently correct than others. When these high-accuracy prototypes are predicted, we can have high confidence that the prediction is correct. If we consider the entry prototype cutoff of 10 we can see 765 766 examples of these high-accuracy entry prototypes including CeAl₂Ga₂ type structure, MgCuAl₂ type structure, and CeNiSi₂ 767 768 type structure which all perform at a precision above 90% 769 which is 20% above the average model precision. Similarly, we 770 can express doubt for predictions involving entry prototypes

that are frequently predicted incorrect in the model. Low-accuracy entry prototypes are rarely valuable as predictions, some examples include TiNiSi type structure, Th₃P₄ type structure, and BiF₃ type structure which scored 20% below the model's average precision. Some classes with very few entries have precisions of 20% or lower. Further confirming the benefit of larger amounts of representatives in data sets.

Although the average precision was stable, the average recall dropped off steadily with smaller cutoff sizes. The average recall ranged from 73±2% to 54±2% showing that as the number of classes increased, the algorithms ability to classify the known training data diminished. This is due to certain classes having only one or two entries after the removal of exotic elements. These are usually categorized as 'Other' again showing the necessity of classes with many entries. An outline of the errors is shown in Figure 6.

	Ca(Ca _{cs} Nd _{cs}) ₂ NbO _c	Ca ₂ Nb ₂ O ₇	CaTi0 ₃	CeAl ₂ Ga ₂	ij	CuZrSiAs	FeAs	GdFeO ₃	K¸NiF₄	LaAIO₃	MgAl₂O₄	MgCu ₂	NaCl	NaFeO ₂	Tinisi	Other	
	Ü																Recall
Ca(Ca _{0.5} Nd _{0.5}) ₂ NbO ₆	94	0	0	0	0	0	0	0	0	0	0	0	0	0	0	43	0.686
Ca ₂ Nb ₂ O ₇	0	95	0	0	0	0	0	0	0	0	0	0	0	0	0	50	0.655
CaTiO ₃	0	0	133	0	0	0	0	12	0	5	0	0	0	0	1	105	0.522
CeAl ₂ Ga ₂	0	0	0	161	0	0	0	0	0	0	0	0	0	0	0	28	0.847
Cu	0	0	0	0	56	0	0	0	0	0	0	0	0	0	0	70	0.444
CuZrSiAs	0	0	0	0	0	93	0	0	0	0	0	0	0	0	0	21	0.816
FeAs	0	0	0	0	0	0	88	0	0	0	0	0	0	0	0	15	0.854
GdFeO ₃	0	0	9	0	0	0	0	454	0	19	0	0	1	0	0	120	0.753
K ₂ NiF ₄	0	0	0	0	0	0	0	3	81	2	0	0	0	0	0	56	0.570
LaAlO ₃	0	0	2	0	0	0	0	33	1	92	0	0	0	0	0	27	0.594
MgAl ₂ O ₄	0	0	0	0	0	0	0	0	0	0	315	0	0	0	0	69	0.820
MgCu₂	0	0	0	0	0	0	0	0	0	0	0	58	0	0	0	53	0.523
NaCl	0	0	0	0	0	0	0	1	0	0	1	0	140	1	0	81	0.625
NaFeO ₂	0	0	0	0	0	0	0	0	0	0	0	0	0	105	0	34	0.755
TiNiSi	0	0	0	1	0	0	3	0	0	0	0	0	0	0	45	65	0.395
Other	0	5	29	5	37	2	18	59	21	16	31	14	21	12	8	20984	0.986
Total	105	100	173	167	93	95	109	562	103	134	103	72	162	118	54	21821	0.950
Precision	0.895	0.950	0.769	0.964	0.602	0.979	0.807	0.808	0.786	0.687	0.908	0.806	0.864	0.890	0.833	0.962	

Figure 8: Confusion matrix of algorithm with a cutoff size of 100. A perfect confusion matrix would have all non-diagonal sections zero. Precision and recall have been rounded to three decimal places.

A confusion matrix was generated for each model. The confu-824 sion matrix for the algorithm with a cutoff of 100 is shown in 825 Figure 8. The error matrix for each class was trained with entry826 prototype imbalances in mind. This was done by normalizing827 the predicted value by the total number of predictions in the 828 class. In this paper, we will focus on the confusion matrix with 829 a bin size of 100 due to the large confusion matrix generated for 830 smaller bin sizes. The algorithm with a cutoff of 100 showed a831 mean precision of 85±2% with a mean recall of 68±2%. In other832 words, the average ability for the algorithm to correctly predict833 a certain structure was 85±2% while the average ability for the 834 algorithm to predict a true positive rate was 68±2%. To clarify835 further the idea behind recall and precision, let us look at836 CaTiO₃ type structure in the entry prototype dataset with a cut-837 off of 150. Out of 162 guesses, the algorithm classified CaTiO₃838 type structures correctly 123 times. This gives the precision of 123/162 or 76%. The recall is the amount of times actual CaTiO₃ type structures was classified correctly. For example, out of 255 known CaTiO3 entry prototypes only 123 were correctly classified giving a recall of 123/255 or 48%.

787

788

789 790

791

792

793

794

795

797

799

800

801 802

803

804 805

806

807 To further understand the misclassification issues, we have 808 compared CaTiO3, LaAlO3, and GdFeO3 which are shown in Figure 9. All lattices show remarkable structural similarities. 809 810 While they are all variations of the standard cubic structure of 811 perovskites, CaTiO3 and GdFeO3 are distorted orthorhombic structures while LaAlO3 is a trigonal structure. The essential 812 structures are similar in terms of polyhedra, bond distances, and 813 814 polyhedral connectivity but vary in terms of polyhedral tilting 815 or rotation. CaTiO3 and GdFeO3 experience this tilting of the 816 octahedra due to calcium and gadolinium being too small to form the cubic structure. GdFeO3 also shows more distortion 817 along the c axis for gadolinium then CaTiO3 does with calcium. 818 LaAlO₃ deviates from the ideal cubic structure by experiencing 820 a rotation of the octahedra due to the length of the aluminum oxygen bond98. 821

822 Future work would be to take individual structures that are quite 823 similar and group them by "family" which would increase the size of representative entries per family while reducing 'Other' category percentage. This would help spread the data across multiple classes and create a more balanced training set. We believe this would help increase recall in the algorithm. Other possible future work would be to extend this approach to identify what synthetic routes would result in different structure types. Finding specific information on synthesis methods can be a difficult proposal due to the vast amount of research papers using different methods to achieve the same goal. However, Kim *et al.* recently published a machine learning paper to collect and organize this data as an interesting approach to overcoming this challenge⁹⁹. However, their work only covered select oxide materials and required 640,000 papers to build a learning model. Some materials of interest may not have sufficiently established literature publications.

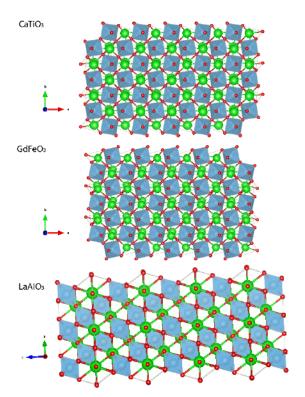


Figure 9: Comparison of misclassification of GdFeO3 with CaTiO3 and LaAlO3. Calcium, gadolinium, and lanthanum are represented in green. Titanium, iron, and aluminum are represented as blue. Oxygen is represented as red.

4. CONCLUSION

840

841

842

843

849

853

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869 870

871

The ability to predict crystal structures remains a challenging problem. The capability to engineer specific materials with cer-904 tain properties will require the ability to predict crystal struc-905 844 tures. Currently, characterization techniques such as diffraction 906 845 or spectroscopy are the standard for assessing a compound's 907 846 crystal structure, but these require pre-made physical sample to 908 847 848 measure. First principle calculations to predict crystal structure, on the other hand, could be used to screen materials prior to 909 synthesis. These approaches have grown in recent years but are 910 850 hindered by long computational times, limited scope, and cost.911 851 Machine learning offers a fundamentally new approach that can 852 operate in concert with the experimental and first principle approaches mentioned earlier by rapidly offering probabilistic 913 854 predictions of crystal structure rather than calculations.

Previous publications have introduced the possibility of ma-916 chine learning-based crystal structure predictions but have been very limited in scope. For example, previous publications dealt⁹¹⁷ only with a range from 3 to 208 specific crystal structures.918 These were limited to binary structures, ternary structures, org19 Heusler/inverse Heusler compounds^{73, 74, 77, 85}. Moreover, previ-920 ous work-built machine learning models which only incorporated training and validation sets limited from 55 to 1948 entries 921 73, 74, 77, 85. In contrast, consider that large inorganic material da-922 tabases such as PCD or ICSD feature around a quarter of a mil-923 lion entries dispersed over more than 10,000 unique crystal₉₂₄ structures at room temperature alone. Therefore, in this contri-bution we have built machine learning models that not only extend far beyond previous work but also begin to address what 926 are the limitations and trade-offs in predicting any arbitrary927 crystal structure. To do so, we have incorporated 24,215 of the928 24,913 structure entries we obtained from the Pearson Crystal

Database. The 24,215 entries were the result of simplifying feature development aspects of the machine learning process and included over 10,000 unique entry prototypes. With these models, we explored the implication of massively imbalanced entry prototype distributions and quantify the model performance associated with compromising model breadth for accuracy. The most notable trade-off is recall, which dropped quickly with a range from 73±2% to 54±2% for minimum-class sizes ranging from 150 to 10, respectively. These values drastically outperform simple metrics such as random guessing, which has a mean-class error ranged from 83.3% to 99.8% and fixing the prediction to the dominant class 'Other', which results in a mean-class error from 86.7% to 99.9%. Reducing the scope of the model had little effect on average precision or accuracy, which was consistent across all the algorithms with a range of $86\pm2\%$ to $79\pm2\%$ and from $97\pm2\%$ to $85\pm2\%$, respectively. Although the model struggled to exhaustively capture all members of a crystal structure, particularly with decreasing class size, the consistently high prediction accuracy is notable.

Successful performance in predicting crystal structure validates this machine learning approach for the exploration of chemical whitespace. We have created a tool that rapidly and efficiently predicts one of the critical factors for physical phenomenon in a material. The output of our machine learning based model is useful to influence or validate other crystal structure approaches. We see particular value when used synergistically with other machine learning algorithms based around physical property prediction.

4. AUTHOR INFORMATION

873

874

875 876

877

878

879

884

885 886

887

888

889

890

892

893

894

895

896

897

898

899

900

901

902

sparks@eng.utah.edu, corresponding author

ACKNOWLEDGEMENTS

The authors gratefully acknowledge support from the NSF CAREER Award DMR 1651668. Special thanks to Anton O. Oliynyk for his assistance in feature creation and counselling as well as Bryce Meredig and Christopher H. Borg from Citrine Informatics for helpful discussions and insight.

- Gurevich, A., Challenges and opportunities for applications of unconventional superconductors. Annu. Rev. Condens. Matter Phys. 2014, 5 (1), 35-56.
- Foltyn, S.; Civale, L.; MacManus-Driscoll, J.; Jia, Q.; Maiorov, B.; Wang, H.; Maley, M., Materials science challenges for high-temperature superconducting wire. Nature materials 2007, 6 (9), 631-642.
- Tarascon, J.-M.; Armand, M., Issues and challenges facing rechargeable lithium batteries. Nature 2001, 414 (6861), 359-367.
- Tarascon, J.-M., Key challenges in future Libattery research. Philosophical Transactions of the Royal Society of London A: Mathematical, Physical and Engineering Sciences 2010, 368 (1923), 3227-3241.
- Mori, D.; Hirose, K., Recent challenges of hydrogen storage technologies for fuel cell vehicles. International journal of hydrogen energy 2009, 34 (10), 4569-4574.

- Shao, Y.; Yin, G.; Wang, Z.; Gao, Y., Proton 929 6.
- 930 exchange membrane fuel cell from low temperature 982
- 931 to high temperature: material challenges. *Journal of* 983
- Power Sources 2007, 167 (2), 235-242. 932
- 933 Kärkäs, M. D.; Åkermark, B., Water oxidation 985
- using earth-abundant transition metal catalysts: 934
- opportunities and challenges. Dalton Transactions 935
- **2016**, *45* (37), 14421-14461. 936
- Kaner, R. B.; Gilman, J. J.; Tolbert, S. H., 937
- Designing superhard materials. Science 2005, 308 938
- (5726), 1268-1269. 939
- 940 9. Zhao, Z.; Xu, B.; Tian, Y., Recent advances in
- superhard materials. Annual Review of Materials 941
- Research 2016, 46, 383-406. 942
- 10. Babcock, E., WHAT IS VULCANIZATION? 943
- Industrial & Engineering Chemistry 1939, 31 (10), 944
- 945 1196-1199.
- 11. Garrett, A. B., The flash of genius. Van 946
- 947 Nostrand: 1963.
- 948 12. Baekeland, L., The Bakelizer: National
- 949 Museum of American History, Smithsonian
- Institution: a National Historic Chemical Landmark, 1002 950
- November 9, 1993. American Chemical Society: 19931003 951
- 952 13. Ball, P., Chemistry: Perkin, the mauve maker.1004
- Nature 2006, 440 (7083), 429-429. 953
- 14. Champagne, C., Serendipity, Super Glue and 1006 954
- Surgery: Cyanoacrylates as Hemostatic Aids in the 955
- Vietnam War. 2009. 956
- 957 15. The Inventor of Saccharine. Munn & Co:
- Scientific American, 1886; Vol. LV, p 36. 958
- 959 16. Shapiro, M.; Reed, C.; Korsmeyer, R.,
- Removal of Tritium from the Molten Salt Breeder 960
- Reactor Fuel. ORNL-MIT-117 1970. 961
- Patel, P., Materials Genome Initiative and 962 17.
- energy. MRS Bulletin 2011, 36 (12), 964-966. 963
- 18. Seshadri, R.; Sparks, T. D., Perspective: 964
- Interactive material property databases through 965
- aggregation of literature data. APL Materials 2016, 41018 966
- 967 (5), 053206. 1019
- 19. Sangid, M. D.; Matlik, J. F., ANALYSIS A better 1020 968
- 969 way to engineer aerospace components. AEROSPACE1021
- 970 AMERICA 2016, 54 (3), 40-43.
- Fink, P. J.; Miller, J. L.; Konitzer, D. G., 971
- Rhenium reduction—alloy design using an 972
- economically strategic element. JOM 2010, 62 (1), 973
- 974 55-57.
- 975 21. Joost, W. J., Reducing vehicle weight and
- improving US energy efficiency using integrated 976
- computational materials engineering. Jom 2012, 64 1029 977
- (9), 1032-1038. 978
- 22. Rieger, T.; Gazdag, S.; Prahl, U.; Mokrov, O.; 1031 979
- 980

- distortion in ship building. Advanced Engineering Materials **2010**, *12* (3), 153-157.
- 23. Pauling, L., The principles determining the structure of complex ionic crystals. Journal of the american chemical society **1929**, 51 (4), 1010-1026.
- Schön, J. C.; Jansen, M., First Step Towards 986 24.
- 987 Planning of Syntheses in Solid-State Chemistry:
- Determination of Promising Structure Candidates by 988
- Global Optimization. Angewandte Chemie 989
- International Edition 1996, 35 (12), 1286-1304. 990
- 991 25. Villars, P., A three-dimensional structural
- 992 stability diagram for 998 binary AB intermetallic
- compounds. Journal of the Less Common Metals 993
- **1983,** *92* (2), 215-238. 994

981

984

1005

- Villars, P., A three-dimensional structural 995 26.
- stability diagram for 1011 binary AB2 intermetallic 996
- compounds: II. Journal of the Less Common Metals 997
- **1984,** *99* (1), 33-43. 998
- 999 27. Villars, P., A semiempirical approach to the
- prediction of compound formation for 3486 binary 1000
- 1001 alloy systems. Journal of the Less Common Metals
 - **1985,** *109* (1), 93-115.
 - Villars, P., A semiempirical approach to the
 - prediction of compound formation for 96446 ternary alloy systems: II. Journal of the Less Common Metals
 - **1986**, *119* (1), 175-188.
- 29. Woodley, S. M.; Catlow, R., Crystal structure 1007
 - prediction from first principles. Nature materials
- 1009 **2008,** 7 (12), 937-946.
- Jones, R. O., Density functional theory: Its 1010 30.
- 1011 origins, rise to prominence, and future. Reviews of
- modern physics **2015**, 87 (3), 897. 1012
- 31. Fischer, C. C.; Tibbetts, K. J.; Morgan, D.; 1013
- Ceder, G., Predicting crystal structure by merging 1014
- data mining with quantum mechanics. Nature 1015
- 1016 materials 2006, 5 (8), 641-646.
- 1017 32. Jain, A.; Shin, Y.; Persson, K. A.,
 - Computational predictions of energy materials using density functional theory. Nature Reviews Materials
 - **2016,** *1,* 15004.
 - Ong, S. P.; Andreussi, O.; Wu, Y.; Marzari, N.;
- 1022 Ceder, G., Electrochemical windows of room-
- temperature ionic liquids from molecular dynamics 1023
- and density functional theory calculations. Chemistry 1024
- of Materials 2011, 23 (11), 2979-2986. 1025
- 1026 Hautier, G.; Fischer, C. C.; Jain, A.; Mueller,
- T.; Ceder, G., Finding nature's missing ternary oxide 1027
- 1028 compounds using machine learning and density
 - functional theory. Chemistry of Materials 2010, 22
- 1030 (12), 3762-3767.
 - 35. Kang, K.; Meng, Y. S.; Bréger, J.; Grey, C. P.;
- Rossiter, E.; Reisgen, U., Simulation of welding and 1032 Ceder, G., Electrodes with high power and high

```
capacity for rechargeable lithium batteries. Science 1086
                                                                          Pittaway, F.; Paz-Borbón, L. O.; Johnston, R.
1033
                                                                  46.
      2006, 311 (5763), 977-980.
                                                                  L.; Arslan, H.; Ferrando, R.; Mottet, C.; Barcaro, G.;
1034
                                                            1087
1035
      36.
              Anasori, B.; Xie, Y.; Beidaghi, M.; Lu, J.;
                                                            1088
                                                                  Fortunelli, A., Theoretical studies of palladium-gold
      Hosler, B. C.; Hultman, L.; Kent, P. R.; Gogotsi, Y.;
                                                                  nanoclusters: Pd- Au clusters with up to 50 atoms.
                                                            1089
1036
      Barsoum, M. W., Two-dimensional, ordered, double 1090
                                                                  The Journal of Physical Chemistry C 2009, 113 (21),
      transition metals carbides (MXenes). Acs Nano 2015,1091
                                                                  9141-9152.
1038
      9 (10), 9507-9516.
1039
                                                            1092
                                                                  47.
                                                                          Shao, G.-F.; Tu, N.-N.; Liu, T.-D.; Xu, L.-Y.;
1040
      37.
              Sharma, V.; Wang, C.; Lorenzini, R. G.; Ma, R.1093
                                                                  Wen, Y.-H., Structural studies of Au-Pd bimetallic
      Zhu, Q.; Sinkovits, D. W.; Pilania, G.; Oganov, A. R.;
                                                                  nanoparticles by a genetic algorithm method. Physica
1041
                                                           1094
      Kumar, S.; Sotzing, G. A., Rational design of all
                                                                  E: Low-dimensional Systems and Nanostructures
1042
                                                            1095
      organic polymer dielectrics. Nature communications 1096
                                                                  2015, 70, 11-20.
1043
1044
      2014, 5, 4845.
                                                            1097
                                                                  48.
                                                                          Swamy, V.; Gale, J. D.; Dubrovinsky, L.,
              Yan, J.; Gorai, P.; Ortiz, B.; Miller, S.; Barnett, 1098
                                                                  Atomistic simulation of the crystal structures and
1045
      38.
      S. A.; Mason, T.; Stevanović, V.; Toberer, E. S.,
                                                                  bulk moduli of TiO2 polymorphs. Journal of Physics
1046
                                                            1099
      Material descriptors for predicting thermoelectric
                                                                  and Chemistry of Solids 2001, 62 (5), 887-895.
                                                            1100
1047
      performance. Energy & Environmental Science 2015, 1101
                                                                          Harris, K. J.; Foster, J. M.; Tessaro, M. Z.;
1048
      8(3), 983-994.
                                                                  Jiang, M.; Yang, X.; Wu, Y.; Protas, B.; Goward, G. R.,
1049
                                                            1102
      39.
1050
              Zhu, H.; Hautier, G.; Aydemir, U.; Gibbs, Z.
                                                                  Structure Solution of Metal-Oxide Li Battery
                                                           1103
1051
      M.; Li, G.; Bajaj, S.; Pöhls, J.-H.; Broberg, D.; Chen, W.1104
                                                                  Cathodes from Simulated Annealing and Lithium
                                                                  NMR Spectroscopy. Chemistry of Materials 2017.
1052
     Jain, A., Computational and experimental
                                                            1105
                                                                          Naserifar, S.; Zybin, S.; Ye, C.-C.; Goddard III,
1053
      investigation of TmAgTe 2 and XYZ 2 compounds, a 1106
      new group of thermoelectric materials identified by 1107
                                                                  W. A., Prediction of structures and properties of 2, 4,
1054
      first-principles high-throughput screening. Journal of 1108
                                                                  6-triamino-1, 3, 5-triazine-1, 3, 5-trioxide (MTO) and
1055
1056
      Materials Chemistry C 2015, 3 (40), 10554-10565.
                                                            1109
                                                                  2, 4, 6-trinitro-1, 3, 5-triazine-1, 3, 5-trioxide
      40.
              Madsen, G. K., Automated search for new
                                                           1110
                                                                  (MTO3N) green energetic materials from DFT and
1057
      thermoelectric materials: the case of LiZnSb. Journal 1111
                                                                  ReaxFF molecular modeling. Journal of Materials
1058
      of the American Chemical Society 2006, 128 (37),
                                                                  Chemistry A 2016, 4 (4), 1264-1276.
1059
                                                            1112
1060
      12140-12146.
                                                            1113
                                                                  51.
                                                                          Kirkpatrick, S.; Gelatt, C. D.; Vecchi, M. P.,
                                                                  Optimization by simulated annealing. science 1983,
1061
              Kolmogorov, A.; Shah, S.; Margine, E.; Bialon,1114
      A.; Hammerschmidt, T.; Drautz, R., New
                                                                  220 (4598), 671-680.
1062
                                                            1115
      superconducting and semiconducting Fe-B
                                                            1116
                                                                  52.
                                                                          Azencott, R., Simulated annealing:
1063
      compounds predicted with an ab initio evolutionary 1117
                                                                  parallelization techniques. Wiley-Interscience: 1992;
1064
      search. Physical review letters 2010, 105 (21),
                                                                  Vol. 27.
1065
                                                            1118
      217003.
                                                                          Ingber, L., Simulated annealing: Practice
1066
                                                            1119
                                                                  53.
      42.
              Li, Y.; Hao, J.; Liu, H.; Li, Y.; Ma, Y., The
                                                                  versus theory. Mathematical and computer
1067
                                                            1120
1068
      metallization and superconductivity of dense
                                                            1121
                                                                  modelling 1993, 18 (11), 29-57.
      hydrogen sulfide. The Journal of chemical physics
                                                            1122
                                                                  54.
                                                                          Bertsimas, D.; Tsitsiklis, J., Simulated
1069
      2014, 140 (17), 174712.
                                                            1123
                                                                  annealing. Statistical science 1993, 8 (1), 10-15.
1070
              Yan, F.; Zhang, X.; Yonggang, G. Y.; Yu, L.;
                                                                  55.
                                                                          Lloyd, L. D.; Johnston, R. L.; Salhi, S.,
1071
      43.
                                                            1124
      Nagaraja, A.; Mason, T. O.; Zunger, A., Design and
                                                                  Strategies for increasing the efficiency of a genetic
1072
                                                            1125
1073
      discovery of a novel half-Heusler transparent hole
                                                            1126
                                                                  algorithm for the structural optimization of nanoalloy
                                                                  clusters. Journal of computational chemistry 2005, 26
      conductor made of all-metallic heavy elements.
                                                            1127
1074
1075
      Nature communications 2015, 6.
                                                            1128
                                                                  (10), 1069-1078.
              Greeley, J.; Jaramillo, T. F.; Bonde, J.;
                                                                  56.
                                                                          Singh, A. K.; Revard, B. C.; Ramanathan, R.;
1076
                                                            1129
      Chorkendorff, I.; Nørskov, J. K., Computational high- 1130
1077
                                                                  Ashton, M.; Tavazza, F.; Hennig, R. G., Genetic
1078
      throughput screening of electrocatalytic materials for 131
                                                                  algorithm prediction of two-dimensional group-IV
      hydrogen evolution. Nature materials 2006, 5 (11), 1132
                                                                  dioxides for dielectrics. Physical Review B 2017, 95
1079
1080
      909-913.
                                                            1133
                                                                  (15), 155426.
      45.
              Lejaeghere, K.; Bihlmayer, G.; Björkman, T.; 1134
                                                                  57.
                                                                          Trimarchi, G.; Freeman, A. J.; Zunger, A.,
1081
      Blaha, P.; Blügel, S.; Blum, V.; Caliste, D.; Castelli, I. E.1135
                                                                  Predicting stable stoichiometries of compounds via
1082
      Clark, S. J.; Dal Corso, A., Reproducibility in density 1136
                                                                  evolutionary global space-group optimization.
1083
      functional theory calculations of solids. Science 2016;137 Physical Review B 2009, 80 (9), 092101.
1084
```

351 (6280), aad3000.

```
58.
              Wu, S.; Ji, M.; Wang, C.-Z.; Nguyen, M. C.;
                                                           1190
                                                                 70.
                                                                         A., R. N.; G., W. D., Phonon Transport in
1138
      Zhao, X.; Umemoto, K.; Wentzcovitch, R.; Ho, K.-M., 1191
                                                                 Asymmetric Sawtooth Nanowires. ASME-JSME
1139
                                                                 Thermal Engineering Joint Conference, March 13-17,
1140
      An adaptive genetic algorithm for crystal structure
                                                           1192
      prediction. Journal of Physics: Condensed Matter
                                                                 2010.
                                                           1193
1141
1142
      2013, 26 (3), 035402.
                                                           1194
                                                                 71.
                                                                         Meredig, B.; Agrawal, A.; Kirklin, S.; Saal, J. E.;
      59.
              Kumar, M.; Husian, M.; Upreti, N.; Gupta, D.,1195
                                                                 Doak, J.; Thompson, A.; Zhang, K.; Choudhary, A.;
1143
      Genetic algorithm: Review and application.
1144
                                                           1196
                                                                 Wolverton, C., Combinatorial screening for new
      International Journal of Information Technology and 1197
                                                                 materials in unconstrained composition space with
1145
      Knowledge Management 2010, 2 (2), 451-454.
                                                                 machine learning. Physical Review B 2014, 89 (9),
1146
                                                           1198
              Pham, D.; Karaboga, D., Intelligent
                                                                 094104.
1147
      60.
                                                           1199
                                                                 72.
      optimisation techniques: genetic algorithms, tabu
                                                           1200
                                                                         Carrete, J.; Li, W.; Mingo, N.; Wang, S.;
1148
1149
      search, simulated annealing and neural networks.
                                                           1201
                                                                 Curtarolo, S., Finding unprecedentedly low-thermal-
      Springer Science & Business Media: 2012.
                                                                 conductivity half-Heusler semiconductors via high-
                                                           1202
1150
      61.
              Froltsov, V. A.; Reuter, K., Robustness of 'cut1203
                                                                 throughput materials modeling. Physical Review X
1151
      and splice'genetic algorithms in the structural
                                                           1204
                                                                 2014, 4 (1), 011019.
1152
      optimization of atomic clusters. Chemical Physics
                                                           1205
                                                                 73.
                                                                         Oliynyk, A. O.; Antono, E.; Sparks, T. D.;
1153
      Letters 2009, 473 (4-6), 363-366.
                                                                 Ghadbeigi, L.; Gaultois, M. W.; Meredig, B.; Mar, A.,
1154
                                                           1206
              Brgoch, J.; Hermus, M., Pressure-Stabilized
                                                           1207
                                                                 High-throughput machine-learning-driven synthesis
1155
1156
      Ir3—in a Superconducting Potassium Iridide. The
                                                           1208
                                                                 of full-Heusler compounds. Chemistry of Materials
1157
      Journal of Physical Chemistry C 2016, 120 (36),
                                                                 2016, 28 (20), 7324-7331.
                                                           1209
                                                                         Balachandran, P. V.; Theiler, J.; Rondinelli, J.
1158
      20033-20039.
                                                           1210
                                                                 74.
      63.
              Miao, M.-s.; Wang, X.-l.; Brgoch, J.; Spera, F.;1211
                                                                 M.; Lookman, T., Materials prediction via
1159
                                                                 classification learning. Scientific reports 2015, 5.
1160
      Jackson, M. G.; Kresse, G.; Lin, H.-q., Anionic
                                                           1212
1161
      chemistry of noble gases: formation of Mg-NG (NG=1213
                                                                         Clarke, B.; Fokoue, E.; Zhang, H. H., Principles
     Xe, Kr, Ar) compounds under pressure. Journal of the 1214
                                                                 and theory for data mining and machine learning.
1162
      American Chemical Society 2015, 137 (44), 14122-
                                                           1215
                                                                 Springer Science & Business Media: 2009.
1163
                                                                         Boser, B. E.; Guyon, I. M.; Vapnik, V. N. In A
      14128.
                                                           1216
                                                                 76.
1164
1165
      64.
              Oganov, A. R.; Glass, C. W.; Ono, S., High-
                                                           1217
                                                                 training algorithm for optimal margin classifiers,
1166
      pressure phases of CaCO3: crystal structure
                                                           1218
                                                                 Proceedings of the fifth annual workshop on
      prediction and experiment. Earth and Planetary
                                                                 Computational learning theory, ACM: 1992; pp 144-
1167
                                                           1219
1168
      Science Letters 2006, 241 (1-2), 95-103.
                                                           1220
                                                                 152.
                                                                 77.
1169
      65.
              Morris, A. J.; Grey, C.; Pickard, C. J.,
                                                           1221
                                                                         Oliynyk, A. O.; Adutwum, L. A.; Harynuk, J. J.;
      Thermodynamically stable lithium silicides and
                                                                 Mar, A., Classifying Crystal Structures of Binary
1170
                                                           1222
      germanides from density functional theory
                                                                 Compounds AB through Cluster Resolution Feature
1171
                                                           1223
      calculations. Physical Review B 2014, 90 (5), 054111. 1224
                                                                 Selection and Support Vector Machine Analysis.
1172
1173
      66.
              See, K. A.; Leskes, M.; Griffin, J. M.; Britto, S.;1225
                                                                 Chemistry of Materials 2016, 28 (18), 6672-6681.
                                                                         Chen, N., Support vector machine in
      Matthews, P. D.; Emly, A.; Van der Ven, A.; Wright, D1226
                                                                 78.
1174
      S.; Morris, A. J.; Grey, C. P., Ab Initio Structure Search 227
                                                                 chemistry. World Scientific: 2004.
1175
      and in Situ 7Li NMR Studies of Discharge Products in 1228
                                                                 79.
                                                                         Theodoridis, S.; Koutroumbas, K., Pattern
1176
      the Li-S Battery System. Journal of the American
                                                                 recognition. Academic press London: 1999.
1177
                                                           1229
1178
      Chemical Society 2014, 136 (46), 16368-16377.
                                                           1230
                                                                         Wang, L.; Brown, S. J. In Prediction of RNA-
                                                                 binding residues in protein sequences using support
1179
      67.
              Curtarolo, S.; Morgan, D.; Persson, K.;
                                                           1231
1180
      Rodgers, J.; Ceder, G., Predicting crystal structures
                                                           1232
                                                                 vector machines, Engineering in Medicine and
      with data mining of quantum calculations. Physical 1233
                                                                 Biology Society, 2006. EMBS'06. 28th Annual
1181
      review letters 2003, 91 (13), 135503.
                                                                 International Conference of the IEEE, IEEE: 2006; pp
1182
                                                           1234
1183
              Batista, G. E.; Monard, M. C., An analysis of 1235
                                                                 5830-5833.
      four missing data treatment methods for supervised 1236
                                                                 81.
                                                                         Janda, J.-O.; Busch, M.; Kück, F.; Porfenenko,
1184
1185
      learning. Applied artificial intelligence 2003, 17 (5-6),1237
                                                                 M.; Merkl, R., CLIPS-1D: analysis of multiple sequence
      519-533.
                                                                 alignments to deduce for residue-positions a role in
1186
                                                           1238
                                                                 catalysis, ligand-binding, or protein structure. BMC
1187
      69.
              Hastie, T.; Tibshirani, R.; Friedman, J., The
                                                           1239
      Elements of Statistical Learning. Second Edition ed.; 1240
                                                                 bioinformatics 2012, 13 (1), 55.
1188
      Springer New York: New York, NY, 2009.
                                                                         Redfern, O. C.; Harrison, A.; Dallman, T.;
1189
                                                           1241
                                                                 Pearl, F. M.; Orengo, C. A., CATHEDRAL: a fast and
```

```
effective algorithm to predict folds and domain
                                                                 Nanoparticles. The Journal of Physical Chemistry
1243
                                                           1282
      boundaries from multidomain protein structures.
                                                           1283
                                                                 Letters 2017.
1244
1245
      PLoS computational biology 2007, 3 (11), e232.
                                                           1284
                                                                 91.
                                                                          Raccuglia, P.; Elbert, K. C.; Adler, P. D.; Falk,
      83.
              Sundararaghavan, V.; Zabaras, N.,
                                                                 C.; Wenny, M. B.; Mollo, A.; Zeller, M.; Friedler, S. A.;
1246
                                                           1285
1247
      Classification and reconstruction of three-
                                                           1286
                                                                 Schrier, J.; Norquist, A. J., Machine-learning-assisted
      dimensional microstructures using support vector
                                                                 materials discovery using failed experiments. Nature
1248
                                                           1287
      machines. Computational Materials Science 2005, 321288
1249
                                                                 2016, 533 (7601), 73-76.
1250
      (2), 223-239.
                                                           1289
                                                                 92.
                                                                          Oliynyk, A. O.; Sparks, T. D.; Gaultois, M. W.;
      84.
              Fang, S.; Wang, M.; Qi, W.; Zheng, F., Hybrid 1290
                                                                 Ghadbeigi, L.; Mar, A., Gd12Co5. 3Bi and Gd12Co5Bi,
1251
      genetic algorithms and support vector regression in 1291
                                                                 Crystalline Doppelgänger with Low Thermal
1252
      forecasting atmospheric corrosion of metallic
                                                           1292
                                                                 Conductivities. Inorganic chemistry 2016, 55 (13),
1253
1254
      materials. Computational Materials Science 2008, 441293
                                                                 6625-6633.
      (2), 647-655.
                                                                 93.
                                                                          Seko, A.; Maekawa, T.; Tsuda, K.; Tanaka, I.,
1255
                                                           1294
      85.
              Oliynyk, A. O.; Adutwum, L. A.; Rudyk, B. W.;1295
                                                                 Machine learning with systematic density-functional
1256
      Pisavadia, H.; Lotfi, S.; Hlukhyy, V.; Harynuk, J. J.;
                                                                 theory calculations: Application to melting
                                                           1296
1257
      Mar, A.; Brgoch, J., Disentangling Structural
                                                           1297
                                                                 temperatures of single-and binary-component solids.
1258
      Confusion through Machine Learning: Structure
                                                                 Physical Review B 2014, 89 (5), 054303.
1259
                                                           1298
      Prediction and Polymorphism of Equiatomic Ternary 1299
                                                                          Schmidt, J.; Shi, J.; Borlido, P.; Chen, L.; Botti,
1260
1261
      Phases ABC. Journal of the American Chemical
                                                           1300
                                                                 S.; Marques, M. A., Predicting the thermodynamic
                                                                 stability of solids combining density functional theory
1262
      Society 2017.
                                                           1301
1263
      86.
              Le, T.; Epa, V. C.; Burden, F. R.; Winkler, D. A.1302
                                                                 and machine learning. Chemistry of Materials 2017.
      Quantitative structure-property relationship
                                                           1303
                                                                 95.
                                                                          Lee, J.; Seko, A.; Shitara, K.; Nakayama, K.;
1264
      modeling of diverse materials properties. Chemical 1304
                                                                 Tanaka, I., Prediction model of band gap for inorganic
1265
1266
      reviews 2012, 112 (5), 2889-2919.
                                                           1305
                                                                 compounds by combination of density functional
      87.
              Basheer, I.; Hajmeer, M., Artificial neural
                                                                 theory calculations and machine learning techniques.
1267
                                                           1306
      networks: fundamentals, computing, design, and
                                                           1307
                                                                 Physical Review B 2016, 93 (11), 115104.
1268
      application. Journal of microbiological methods 20001308
                                                                          Mannodi-Kanakkithodi, A.; Pilania, G.; Huan,
1269
1270
      43 (1), 3-31.
                                                           1309
                                                                 T. D.; Lookman, T.; Ramprasad, R., Machine learning
      88.
1271
              Patel, M. S.; Mazumdar, H. S., Knowledge
                                                           1310
                                                                 strategy for accelerated design of polymer
      base and neural network approach for protein
                                                                 dielectrics. Scientific reports 2016, 6, 20952.
1272
                                                           1311
      secondary structure prediction. Journal of theoretical 312
                                                                 97.
                                                                          Pilania, G.; Wang, C.; Jiang, X.; Rajasekaran,
1273
1274
      biology 2014, 361, 182-189.
                                                           1313
                                                                 S.; Ramprasad, R., Accelerating materials property
              Holley, L. H.; Karplus, M., Protein secondary 1314
                                                                 predictions using machine learning. Scientific reports
1275
      structure prediction with a neural network.
                                                                 2013, 3.
1276
                                                           1315
      Proceedings of the National Academy of Sciences
                                                           1316
                                                                 98.
                                                                          Tilley, R. J., Perovskites: structure-property
1277
1278
      1989, 86 (1), 152-156.
                                                           1317
                                                                 relationships. John Wiley & Sons: 2016.
              Timoshenko, J.; Lu, D.; Lin, Y.; Frenkel, A. I., 1318
                                                                          Kim, E.; Huang, K.; Tomala, A.; Matthews, S.;
      90.
                                                                 99.
1279
      Supervised Machine Learning-Based Determination 1319
                                                                 Strubell, E.; Saunders, A.; McCallum, A.; Olivetti, E.,
1280
```

1320

1321

Machine-learned and codified synthesis parameters

of oxide materials. Scientific Data 2017, 4.

of Three-Dimensional Structure of Metallic