

# Discovering Chemically Novel, High-Temperature Superconductors

Colton C. Seegmiller\*, Sterling G. Baird†, Hasan M. Sayeed‡, and Taylor D. Sparks§¶

\*Department of Engineering, Utah Valley University

Orem, UT 84058, USA

Email: 10803013@uvu.edu

†Department of Materials Science and Engineering, University of Utah

Salt Lake City, UT 84108, USA

Email: sterling.baird@utah.edu

‡Department of Materials Science and Engineering, University of Utah

Salt Lake City, UT 84108, USA

Email: hasan.sayeed@utah.edu

§Department of Materials Science and Engineering, University of Utah

Salt Lake City, UT 84108, USA

Email: sparks@eng.utah.edu

¶Chemistry Department, University of Liverpool

Liverpool L7 3NY, United Kingdom

**Abstract**—One of the biggest unsolved problems in condensed matter physics is what mechanism causes high-temperature superconductivity and if there is a material that can exhibit superconductivity at both room temperature and atmospheric pressure. Among the many important properties of a superconductor, the critical temperature ( $T_c$ ) or transition temperature is the point at which a material transitions into a superconductive state. In this implementation, machine learning is used to predict the critical temperatures of chemically unique compounds in an attempt to identify new chemically novel, high-temperature superconductors. The training data set (SuperCon) consists of known superconductors and their critical temperatures, and the testing data set (NOMAD) consists of around 700,000 novel chemical formulae. The chemical formulae in these data sets are first passed through a collection of rapid screening tools, **SMACT**, to check for chemical validity. Next, the **DiSCoVeR** algorithm is used to train on the SuperCon data to form a model, and then screens through batches of the formulae in the NOMAD data set. Having a combination of a chemical distance metric, density-aware dimensionality reduction, clustering, and a regression model, the **DiSCoVeR** algorithm serves as a tool to identify and assess these superconducting compositions [1]. This research and implementation resulted in the screening of chemically novel compositions exhibiting critical temperatures upwards of 150 K, which correlates to superconductors in the cuprate class. This implementation demonstrates a process of performing machine learning-assisted superconductor screening (while exploring chemically distinct spaces) which can be utilized in the materials discovery process.

**Index Terms**—machine learning, materials informatics, high-temperature superconductor, critical temperature, transition temperature, cuprates

## I. INTRODUCTION

Superconductivity has been a major focus in research since its discovery in 1911 [2]. The discovery of a material that exhibits superconductivity at operating temperatures above 273 K and at atmospheric pressure (101 kPa) would have

an enormous technological impact. It would absolutely revolutionize the fields of digital electronics and the electric power industry. For many years, all known superconductors were thought to exist within the bounds of Bardeen-Cooper-Schrieffer (BCS) theory, which stated that the superconductivity of materials could not exist above temperatures of 30 K [3]. It wasn't until 1986 when Johannes G. Bednorz and Karl A. Müller discovered a new class of superconductor in the cuprate family that exceeded this BCS theory threshold. [4]. As explained in [5], “the superconducting cuprates are very different from conventional superconductors, in the fact that they are not traditional metals, but instead doped oxides that behave like bad metals. Often, the pairing for superconduction does not happen with electrons, but instead with the doped holes – which act as quasiparticles that pair up and behave like the Cooper pairs, but with opposite charge. It is still not fully known what drives the pairing mechanism to get superconductivity in these materials.” Materials with these properties are deemed in the category of a type-II superconductor. Other types of superconductors have since been discovered beyond the cuprate family such as heavy-fermion-based, buckminsterfullerene-based, carbon-allotrope, iron-pnictogen-based, nickel-based, and strontium-ruthenate superconductors among others.

It was also a cuprate that was discovered with a critical temperature above the boiling point of liquid nitrogen (77 K). This led to the realization that applications of superconductivity were looking more realistic and feasible in the near future [6]. Superconductors with a critical temperature above the boiling point of liquid nitrogen are called high-temperature superconductors. It is important to note that all high-temperature superconductors are type-II superconductors. To date, cuprate superconductors hold the record for the

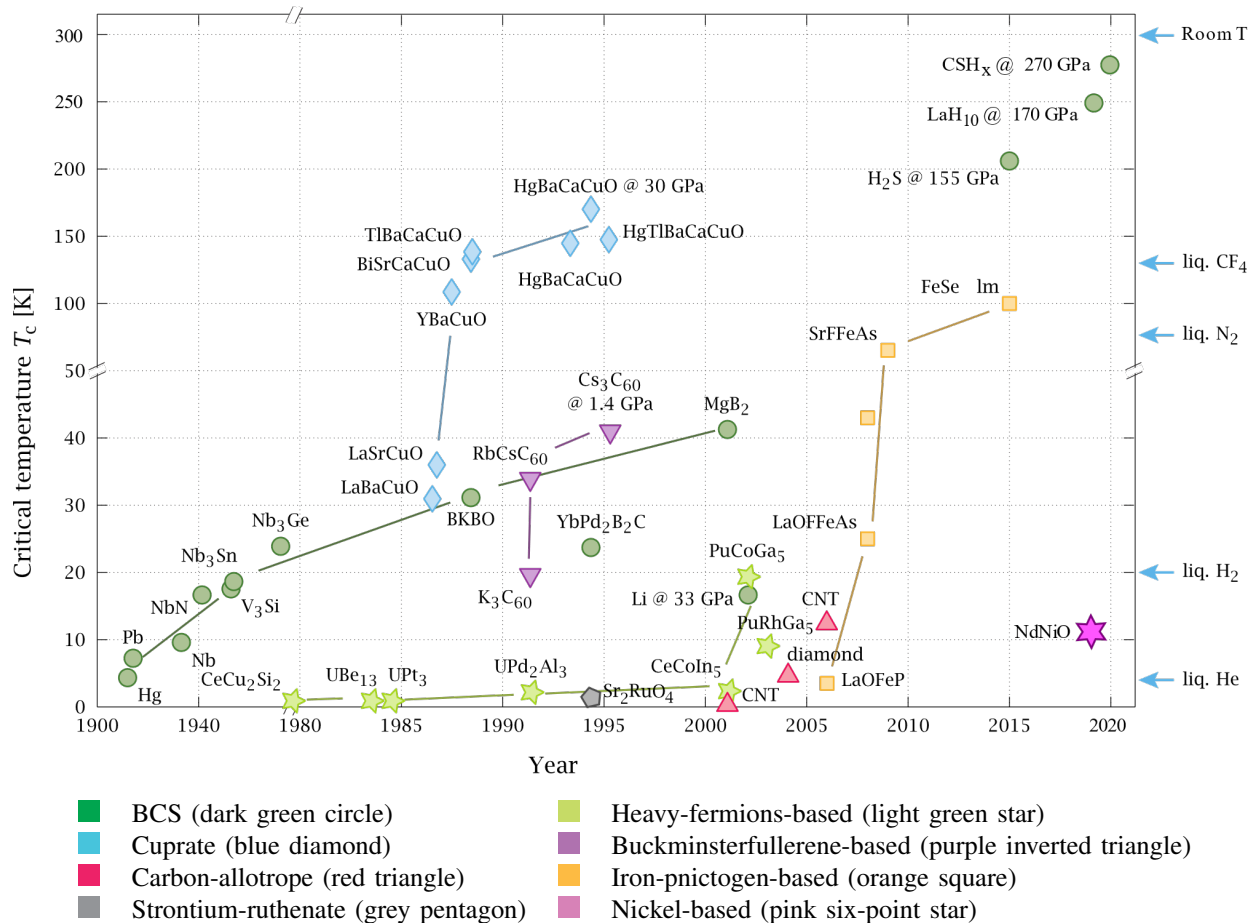


Fig. 1: Timeline of superconductors as adapted from [5]. Colors represent different classes of materials. Note the change in axes around 1980 and 50 K

highest critical temperature at atmospheric pressure. In the past few years, other materials have demonstrated high critical temperatures, but only at extremely high pressures [7]. Fig. 1 shows a timeline of the discovery of superconductors and their critical temperatures.

A critical technological need will be to bridge the increasingly high-temperature performance with ambient pressure [8]. A useful tool to accomplish this goal of accelerated superconductive materials discovery is through machine learning, where there has been various implementations. For example, superconducting phase diagrams were predicted using text mining [9], superconducting hydrogen compounds were found using a genetic algorithm and genetic programming [10], critical temperature and pressure were predicted for hydrides [11], critical temperatures of doped Fe-based superconductors were predicted based on structural and topological parameters [12], and critical temperature was predicted on a structure based model using a structural descriptor [13], and superconductor materials and properties have been automatically extracted from literature [14]. An ML-guided discovery will hopefully replace the “serendipitous discovery paradigm” that has existed in this last century of superconductor research [15].

In this work, we use the SuperCon data set for training, similar to what has been done in other implementations [13], [16]–[20]. Unique, reduced chemical formulae are curated [21] from the NOMAD data set [22] and used for testing. The chemical formulae in the training and testing data are first screened through SMACT [23] for validity and then trained and predicted using the DiSCoVeR algorithm [1]. This results in the screening of novel, chemically valid formulae with predicted critical temperatures.

## II. METHODS

### A. Data

Materials informatics has shown that the cuprate class of superconductors contains a highly unexplored materials space that has yet to be explored [19]. This is why formulae from the SuperCon database are used for the training of our model. Of these formulae, “roughly 5,700 compounds are cuprates and 1,500 are iron-based (about 35 and 9 percent, respectively), reflecting the significant research efforts invested in these two families. The remaining set of about 8,000 is a mix of various materials, including conventional phonon-driven superconductors (e.g., elemental superconduc-

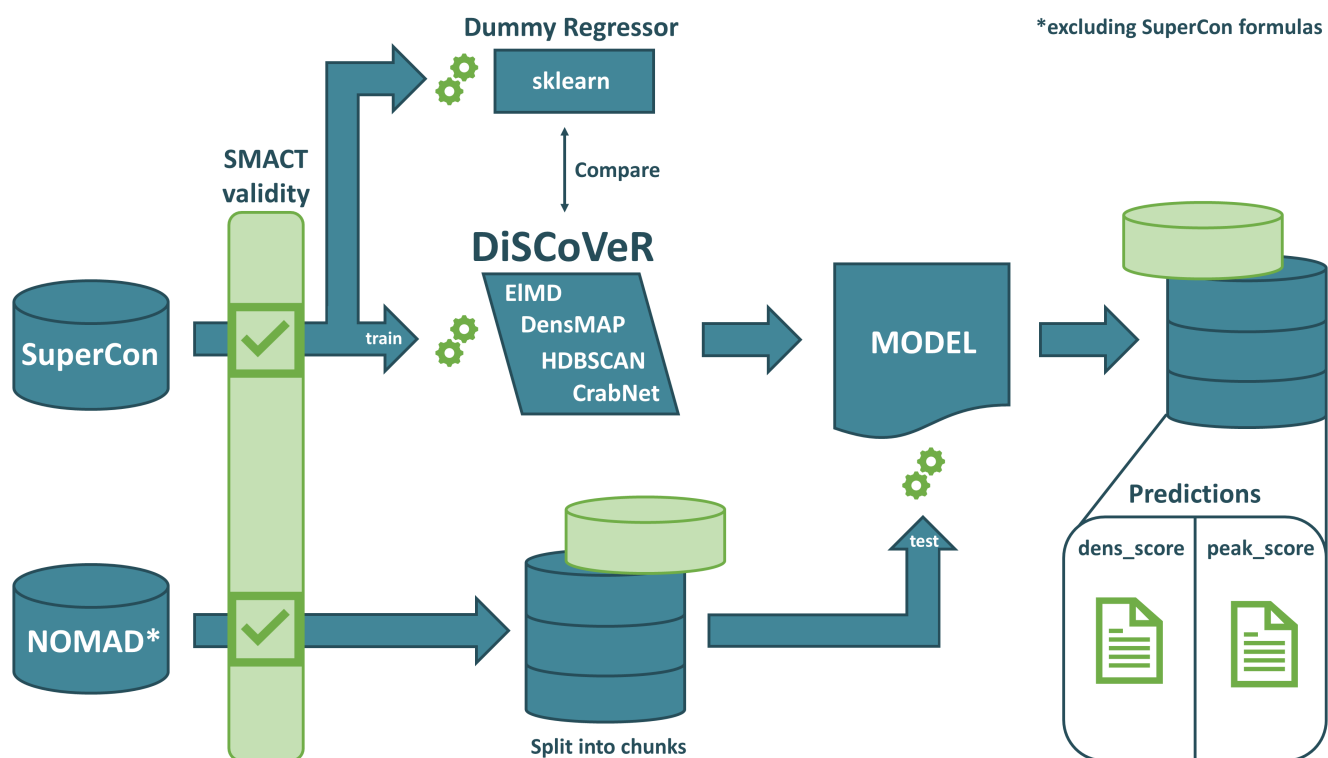


Fig. 2: Workflow of implementation. 12,415 formulae from SuperCon and 694,398 formulae from NOMAD are first verified through SMACT. After featurization and training of the SuperCon data through DiSCoVeR, chunks of the NOMAD data are screened through the model

tors, A15 compounds), known unconventional superconductors like the layered nitrides and heavy fermions, and many materials for which the mechanism of superconductivity is still under debate (such as bismuthates and borocarbides)” [16]. The compositions in SuperCon were reduced using the `get_reduced_composition_and_factor()` method from the `pymatgen.core.Composition` class. After some data cleaning, the SuperCon training data was reduced to 12,415 formulae of superconductors and their critical temperatures. Fig. 3a shows the distribution of the SuperCon data set after cleaning.

Compositions from the Novel Materials Discovery (NOMAD) data set are used in the prediction of our model. This repository contains computational materials science data that is allowed to be curated [22]. For this implementation, we used a specific curated data set of unique reduced chemical formulae [21]. This data was restricted to density functional theory (DFT) calculations and does not include noble gases or radioactive elements. It is also directly usable with the `pymatgen.core.Composition` class [24], which is what this implementation exploits.

Additional curating was done on the NOMAD data set. Next, the formulae in NOMAD that overlapped with formulae in the SuperCon training data were removed for better accuracy while predicting. The NOMAD data was reduced to 694,398 compositions. The training data and test data are then

screened through SMACT for validity (see Fig. 2).

### B. SMACT

SMACT is a composition-based screening tool [23]. It generates a search space, or a set of element combinations, that is screened using chemical filters. Oxidation states, charge neutrality, and electronegativity can be considered to screen for candidates that make “chemical sense.” If the overall charge of a composition is neutral, then SMACT will consider it valid. The original checker however does not consider the countless combinations of oxidation states for metal alloys. To account for this, materials composed of all metal elements are assumed valid in the checker. To perform this, we implement a function called `smact_validity()`. The Boolean values for each of the predicted compositions are under the `is_valid` column in Table II

### C. DiSCoVeR

DiSCoVeR stands for Descending from Stochastic Clustering Variance Regression. This algorithm is a conglomerate of multiple tools (as shown in Fig. 2) that are ultimately used for the screening and assessment of the superconductive compositions. “DiSCoVeR screens candidates that have a high probability of success while enforcing – through the use of novel loss functions – that the candidates exist beyond typical materials landscapes *and* have high performance. In other

TABLE I: Methods used in the DiSCoVeR algorithm. Reproduced from [1] with permission from the Royal Society of Chemistry.

Method	What is it?	What is its role in DiSCoVeR?
CrabNet [25]	Composition-based property regression	Predict performance for proxy scores
EIMD [26]	Composition-based distance metric	Supply distance matrix to DensMAP
DensMAP [27]	Density-aware dimensionality reduction	Obtain densities for density proxy
HDBSCAN [28]	Density-aware clustering	Create chemically homogenous clusters
Peak proxy	High performance relative to nearby compounds	Proxy for surprising high performance
Density proxy	Sparsity relative to nearby compounds	Proxy for chemical novelty
Peak proxy score	Weighted sum of performance and peak proxy	Used to rank compounds
Density proxy score	Weighted sum of performance and density proxy	Used to rank compounds
Pareto front	Optimal performance/uniqueness trade-offs	Visually screen compounds (no weights)

words, DiSCoVeR acts as a multi-objective screening where the promise of a compound depends on both having desirable target properties and existing in sparsely populated regions of the cluster to which it's assigned. This approach then favors discovery of novel, high-performing chemical families as long as embedded points which are close together or far apart exhibit chemical similarity or chemical distinctiveness, respectively" [1]. Table I describes the methods used in DiSCoVeR and explains each of their roles.

The training data for the CrabNet model consists only of the compositions and the measured property values (in this case, critical transition temperature). `peak_proxy` and `dens_proxy` come into play when assessing the trade-offs between performance and novelty. This can be visualized via a Pareto front (performance vs. novelty), which exists only as a visualization tool in this work, and is a complement to the scaled and weighted sums of performance novelty.

The training data is also trained using sklearn's `DummyRegressor` and the mean average error (MAE) is compared alongside that of DiSCoVeR's to serve as a metric. During testing, the NOMAD data set is partitioned into chunks to help with computation. The predictions for high-performing compositions are appended and organized after being screened through the trained model.

### III. RESULTS AND DISCUSSION

There are many essential properties to consider in the search for a novel superconductor such as pressure information, the material's critical magnetic field, its critical current density, phase diagram information, and additional structural data. When considering the entire materials discovery process, synthesizing and screening candidate materials for superconductivity is the final objective. Critical temperature is the most reasonable superconductor property to predict since pressure, critical magnetic field, and critical current density are more difficult, intensive, and expensive to measure and less evident in current data. In regard to extrapolation performance for superconductor discovery, Meredig et al. states that "novel materials discovery would be enabled by running a model against a large database of candidate compounds and simply ranking them by predicted  $T_c$ " [29], which is what is done

in this implementation. For this specific implementation, a composition-based approach is used to test the limits of this algorithm by predicting a single property: a material's critical temperature.

#### A. *peak\_score* and *dens\_score*

The DiSCoVeR algorithm has two expected outputs, `peak_score` and `dens_score`, which are two different ways of evaluating the joint performance and novelty of a compound. `peak_score` uses peak proxy as the underlying novelty proxy which favors compounds that have "surprising" high performance. In other words, these compounds stand out as a high-performance "peak" in the embedding space relative to surrounding compounds that are low performance. On the other hand, `dens_score` uses density proxy which favors compounds that are far from the training data points (i.e. low density in the embedding space), irrespective of the target values. The overall score is determined by considering both the performance (predicted property) and a novelty proxy (either `peak_score` or `dens_score` in this work); however, these scores can have different units and different distributions that make it difficult to compare the values fairly. For example, peak proxy has the same units as the predicted property (degrees Kelvin in this case) but is the difference between the predicted property and the average of the neighboring compounds, On the other hand, density proxy is a measure of density in an embedding space (i.e. non-physical units) and has no direct comparison with performance. Since the absolute values can vary by orders of magnitude, it's necessary to scale the values to a relative range prior to aggregating them (e.g. via a weighted sum). However, the choice of scaling is arbitrary. For example, one could use sklearn's `MinMaxScaler` or `RobustScaler`, or any number of common or custom scaling functions. In order to reduce the effect of outliers, `RobustScaler` was used.

After scaling, the default weighting is a 50/50 between the performance score (i.e. the predicted property) and the novelty proxy; however, this is also a parameter that can be adjusted by the user. In this implementation, `peak_score` and `dens_score` are both considered and set at the default weighting.

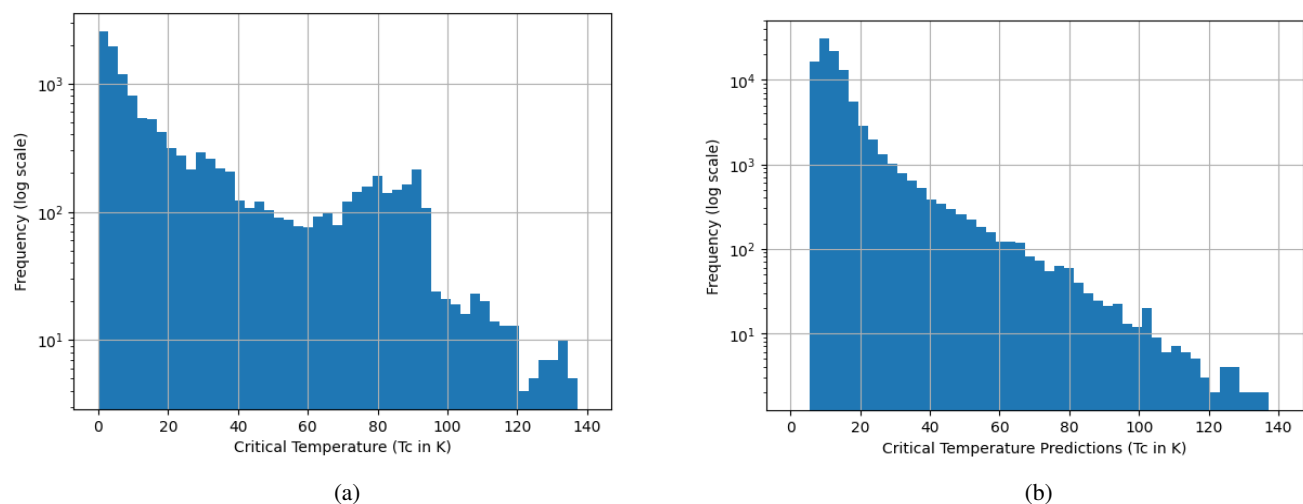


Fig. 3: (a) Distribution of SuperCon data set after cleaning, critical temperature ( $T_c$  in K) on x-axis and frequency in log scale on y-axis (b) Distribution of NOMAD data set after predictions, critical temperature predictions ( $T_c$  in K) on x-axis and frequency in log scale on y-axis

### B. $T_c$ prediction

Table II shows the 20 highest screened compositions after being sorted by  $T_c$  (shown in the prediction column). The columns for both scores aren't sorted since they are both evenly weighted, and both high-performing and chemically novel compositions are desired.

The distribution for the top 100,000 of these screened compositions is shown in Fig. 3b.

### C. Synthesizability prediction

Since the test data was obtained from the NOMAD database, a repository of computationally-generated materials, we aimed to evaluate the synthesizability of the superconductors with predicted critical temperatures. To assess their stability, we queried the materials from the Materials Project and obtained their energy above hull values. The lower the energy above hull, the more stable the compound is considered to be. These values are under the `predicted_e_above_hull` column in Table II. Additionally, we obtained the `is_theoretical` property, which indicates if a material has been reported in the International Crystal Structure Database (ICSD) (i.e. if it has been synthesized previously).

We trained a CrabNet [25] model using data from the Materials Project, optimizing its hyperparameters with the Adaptive Experimentation (Ax) Platform (<https://ax.dev>). After training, this model was used to predict the energy above hull and `is_theoretical` property for the superconductors in question. For the `is_theoretical` property, a value close to 1 indicates that similar compounds have not been synthesized previously, and their synthesis would represent a new exploration in the chemical space. A value close to 0, on the other hand, suggests that similar compounds have been synthesized before, and their synthesis would be considered exploitation. These values are under the `is_theoretical` column in Table II.

### D. Conditional thresholds

Considering each of the properties represented as columns in Table II, conditional thresholds were determined to identify the best formulae. This condition was determined by evaluating the uncertainties of the properties. The condition is met if `is_valid == TRUE & predicted_e_above_hull <= 0.1 & is_theoretical >= 0.95`. Compositions that meet these conditions are the ones shown in Table II.

## IV. CONCLUSION

12,415 known superconductors in the SuperCon database were first validated through SMACT, and then trained on the DISCOVeR algorithm. 694,398 curated, chemically-novel formulae were taken from the NOMAD repository, also validated through SMACT, and then screened through the trained model in chunks. Critical temperatures for each of the formulae in this NOMAD data set were predicted. A weighted uniqueness/performance ranking for each of the compositions was obtained. These sorted compositions also include a Boolean value to whether or not they are valid according to SMACT.

After screening these compositions, additional post-processing work was done to predict energy above hull and stability, which are useful metrics for synthesis. Finally, the compositions were filtered through a condition to get the best representation of formulae that are ready for synthesis.

This implementation reveals a process of performing ML-assisted superconductor screening using an algorithm that uniquely accounts for chemical similarity, and identifies and evaluates new high-performing, chemically distinct compositions. These predicted compositions are openly available in the hopes of being used in the materials discovery process. Since these validation formulae are sorted, they can now undergo additional post-processing and characterization.

TABLE II: Top 20 Screened Compositions

formula	prediction ( $T_c$ in K)	dens_score	peak_score	is_valid	predicted_e_above_hull (eV/atom)	is_theoretical
CaSbPb4	129.98	12.087	24.588	TRUE	0.02017	0.99966
Ba4Ca4Cu6Hg2O17	129.46	75.558	20.483	TRUE	0.04358	1.00037
Ba6Ca6Cu9Hg3O25	128.02	90.421	19.858	TRUE	0.03249	0.98616
Ba2CaTi	126.58	2.4024	24.832	TRUE	-0.00004	0.99990
Ba2Ca3TiCu4O11	125.17	51.912	20.975	TRUE	0.02111	0.99982
BaY7	119.31	2.7561	22.166	TRUE	0.09270	0.99994
Na(Cu3O4)2	112.20	11.140	20.507	TRUE	0.05103	0.99985
TiCuO2	110.69	11.381	20.107	TRUE	0.00912	1.00002
CrHO2	103.07	9.6154	19.126	TRUE	0.00159	1.00069
Ca3Ti2O6	102.97	40.688	17.146	TRUE	0.00299	0.99564
CrAuO2	101.84	18.792	18.565	TRUE	0.00012	1.00023
AlTiO2	100.97	10.007	19.484	TRUE	0.04094	1.00061
Ba3Sr(Cu2O5)2	100.83	102.87	13.320	TRUE	0.06490	0.99991
CuHgO2	100.30	19.942	18.027	TRUE	0.04921	1.00021
CdAgO2	99.445	20.981	17.237	TRUE	0.02421	1.00022
Ca6Al7O16F	98.903	11.387	18.215	TRUE	0.00013	1.00011
Ca3LaMn4O12	98.880	12.568	17.770	TRUE	0.06862	1.00001
TlAgO2	96.406	13.486	17.558	TRUE	0.00192	1.00004
Ca10Ti8NbAl(SiO5)10	95.957	10.442	17.674	TRUE	0.00445	0.99938
Ca33In4P28Pb3O112	95.874	37.240	15.794	TRUE	0.00152	1.00006

## ACKNOWLEDGMENT

This work was supported by the National Science Foundation (NSF). C.C.S. acknowledges support from NSF Grant No. DMR-1950589 and Utah Valley University's Undergraduate Research Scholarly and Creative Activities (URSCA) program. H.M.S. and T.D.S. acknowledge support from NSF Grant No. 1936383. S.G.B. and T.D.S. acknowledge support from NSF Grant No. 1651668. We acknowledge **Erick Lawrence** for the idea of compositional stability prediction.

## CREDIT STATEMENT

**Colton C. Seegmiller**: Conceptualization, Methodology, Software, Validation, Formal analysis, Investigation, Writing - Original Draft, Writing - Review & Editing. **Sterling Baird**: Supervision, Project administration, Conceptualization, Methodology, Software, Writing - Original Draft, Writing - Review & Editing. **Hasan M Sayeed**: Methodology, Software, Formal analysis, Investigation, Writing - Original Draft, Writing - Review & Editing. **Taylor D. Sparks**: Supervision, Project administration, Funding acquisition, Conceptualization, Formal analysis, Resources, Writing - Review & Editing.

## DATA AVAILABILITY

The raw data required to reproduce these findings are available to download from <https://github.com/vstanev1/Supercon>, [https://figshare.com/articles/dataset/NOMAD\\_Chemical\\_Formulas\\_and\\_Calculation\\_IDs/19319783](https://figshare.com/articles/dataset/NOMAD_Chemical_Formulas_and_Calculation_IDs/19319783). The processed data required to reproduce these findings are available to download from <https://github.com/cseeg/DiSCoVeR-SuperCon-NOMAD-SMACT>, this link also includes the full .csv files of screened compositions (with and without the conditional threshold applied).

## REFERENCES

- [1] Sterling G Baird, Tran Q Diep, and Taylor D Sparks. Discover: a materials discovery screening tool for high performance, unique chemical compositions. *Digital Discovery*, 1(3):226–240, 2022.
- [2] H Kamerlingh Onnes. The resistance of pure mercury at helium temperatures. further experiments with liquid helium. iv. In *Proceedings Koninklijke Akademie van Wetenschappen te Amsterdam*, volume 13, pages 1274–1276, 1911.
- [3] Mark Buchanan. Mind the pseudogap. *Nature*, 409(6816):8–12, 2001.
- [4] J. G. Bednorz and K. A. Muller. Possible high  $T_c$  superconductivity in the Ba-La-Cu-O system. *Z. Phys. B*, 64:189–193, 1986.
- [5] PJ Ray. Master's thesis: Structural investigation of  $La_{2-x}Sr_xCuO_{4+y}$  following staging as a function of temperature. *University of Copenhagen*, 2016.
- [6] Maw-Kuen Wu, Jo R Ashburn, Clj Torng, Pei-Herng Hor, Ri L Meng, Lo Gao, Z Jo Huang, YQ Wang, and aCW Chu. Superconductivity at 93 k in a new mixed-phase y-ba-cu-o compound system at ambient pressure. *Physical review letters*, 58(9):908, 1987.
- [7] Elliot Snider, Nathan Dasenbrock-Gammon, Raymond McBride, Xiaoyu Wang, Noah Meyers, Keith V Lawler, Eva Zurek, Ashkan Salamat, and Ranga P Dias. Synthesis of yttrium superhydride superconductor with a transition temperature up to 262 k by catalytic hydrogenation at high pressures. *Physical Review Letters*, 126(11):117003, 2021.
- [8] Simone Di Cataldo, Christoph Heil, Wolfgang von der Linden, and Lilia Boeri.  $La_{1-x}Ba_xHf_2$ : Towards high- $T_c$  low-pressure superconductivity in ternary superhydrides. *Physical Review B*, 104(2):L020511, 2021.
- [9] Callum J Court and Jacqueline M Cole. Magnetic and superconducting phase diagrams and transition temperatures predicted using text mining and machine learning. *npj Computational Materials*, 6(1):18, 2020.
- [10] Takahiro Ishikawa, Takashi Miyake, and Katsuya Shimizu. Materials informatics based on evolutionary algorithms: Application to search for superconducting hydrogen compounds. *Physical Review B*, 100(17):174506, 2019.
- [11] Michael J Hutcheon, Alice M Shipley, and Richard J Needs. Predicting novel superconducting hydrides using machine learning approaches. *Physical Review B*, 101(14):144505, 2020.
- [12] Yun Zhang and Xiaojie Xu. Predicting doped fe-based superconductor critical temperature from structural and topological parameters using machine learning. *International Journal of Materials Research*, 112(1):2–9, 2021.
- [13] Jingzi Zhang, Zhuoxuan Zhu, X-D Xiang, Ke Zhang, Shangchao Huang, Chengquan Zhong, Hua-Jun Qiu, Kailong Hu, and Xi Lin. Machine learning prediction of superconducting critical temperature through the structural descriptor. *The Journal of Physical Chemistry C*, 126(20):8922–8927, 2022.

- [14] Luca Foppiano, Pedro Baptista Castro, Pedro Ortiz Suarez, Kensei Terashima, Yoshihiko Takano, and Masashi Ishii. Automatic extraction of materials and properties from superconductors scientific literature. *Science and Technology of Advanced Materials: Methods*, 3(1):2153633, 2023.
- [15] Stephan R Xie, Gregory R Stewart, James J Hamlin, Peter J Hirschfeld, and Richard G Hennig. Functional form of the superconducting critical temperature from machine learning. *Physical Review B*, 100(17):174513, 2019.
- [16] Valentin Stanev, Corey Oses, A Gilad Kusne, Efrain Rodriguez, John-pierre Paglione, Stefano Curtarolo, and Ichiro Takeuchi. Machine learning modeling of superconducting critical temperature. *npj Computational Materials*, 4(1):29, 2018.
- [17] Zhong-Li Liu, Peng Kang, Yu Zhu, Lei Liu, and Hong Guo. Material informatics for layered high-*t<sub>c</sub>* superconductors. *APL Materials*, 8(6):061104, 2020.
- [18] Tomohiko Konno, Hodaka Kurokawa, Fuyuki Nabeshima, Yuki Sakishita, Ryo Ogawa, Iwao Hosako, and Atsutaka Maeda. Deep learning model for finding new superconductors. *Physical Review B*, 103(1):014509, 2021.
- [19] Rhys EA Goodall, Bonan Zhu, Judith L MacManus-Driscoll, and Alpha A Lee. Materials informatics reveals unexplored structure space in cuprate superconductors. *Advanced Functional Materials*, 31(52):2104696, 2021.
- [20] Elizabeth A Pogue, Alexander New, Kyle McElroy, Nam Q Le, Michael J Pekala, Ian McCue, Eddie Gienger, Janna Domenico, Elizabeth Hedrick, Tyrel M McQueen, et al. Closed-loop machine learning for discovery of novel superconductors. *arXiv preprint arXiv:2212.11855*, 2022.
- [21] Sterling G. Baird. NOMAD Chemical Formulas and Calculation IDs. 3 2022.
- [22] Claudia Draxl and Matthias Scheffler. The nomad laboratory: from data sharing to artificial intelligence. *Journal of Physics: Materials*, 2(3):036001, 2019.
- [23] Daniel W Davies, Keith T Butler, Adam J Jackson, Jonathan M Skelton, Kazuki Morita, and Aron Walsh. Smact: Semiconducting materials by analogy and chemical theory. *Journal of Open Source Software*, 4(38):1361, 2019.
- [24] Shyue Ping Ong, William Davidson Richards, Anubhav Jain, Geoffroy Hautier, Michael Kocher, Shreyas Cholia, Dan Gunter, Vincent L Chevrier, Kristin A Persson, and Gerbrand Ceder. Python materials genomics (pymatgen): A robust, open-source python library for materials analysis. *Computational Materials Science*, 68:314–319, 2013.
- [25] Anthony Yu-Tung Wang, Steven K Kauwe, Ryan J Murdock, and Taylor D Sparks. Compositionally restricted attention-based network for materials property predictions. *Npj Computational Materials*, 7(1):77, 2021.
- [26] Cameron J Hargreaves, Matthew S Dyer, Michael W Gaultois, Vitaliy A Kurlin, and Matthew J Rosseinsky. The earth mover’s distance as a metric for the space of inorganic compositions. *Chemistry of Materials*, 32(24):10610–10620, 2020.
- [27] Ashwin Narayan, Bonnie Berger, and Hyunghoon Cho. Density-preserving data visualization unveils dynamic patterns of single-cell transcriptomic variability. *bioRxiv*, pages 2020–05, 2020.
- [28] Leland McInnes, John Healy, and Steve Astels. hdbscan: Hierarchical density based clustering. *J. Open Source Softw.*, 2(11):205, 2017.
- [29] Bryce Meredig, Erin Antono, Carena Church, Maxwell Hutchinson, Julia Ling, Sean Paradiso, Ben Blaiszik, Ian Foster, Brenna Gibbons, Jason Hattrick-Simpers, et al. Can machine learning identify the next high-temperature superconductor? examining extrapolation performance for materials discovery. *Molecular Systems Design & Engineering*, 3(5):819–825, 2018.