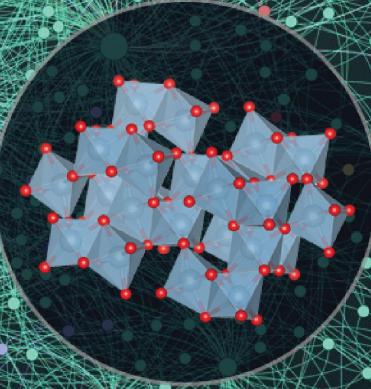
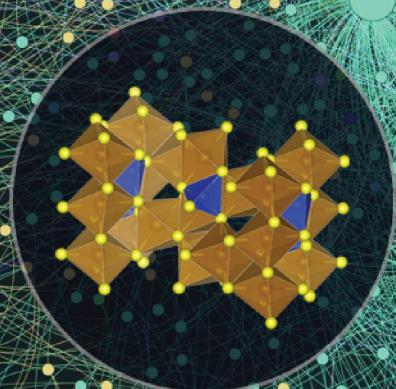


# Composition-based feature vector



How do we featurize a chemical formula for machine learning?



# How do we featurize a chemical formula for machine learning?



ceramic, hard, high melting point, strong



metal, alloy, intermediate melting point, ductile, corrosion resistant

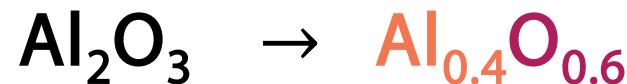


polymer, non-stick, low melting point, low friction

We use element properties to generate a Composition-Based Feature Vector (CBFV)

Element	Number	Atomic_Weight	Period	Group	Families	Metal	Nonmetal	Metalliod	Atomic_Radius
H	1	1.00794	1	1	7	0	1	0	0.53
He	2	4.002602	1	18	9	0	1	0	0.31
Li	3	6.941	2	1	1	1	0	0	1.67
Be	4	9.01218	2	2	2	1	0	0	1.12
B	5	10.811	2	13	6	0	0	1	0.87
C	6	12.011	2	14	7	0	1	0	0.67
N	7	14.00674	2	15	7	0	1	0	0.56
O	8	15.9994	2	16	7	0	1	0	0.48
F	9	18.998403	2	17	8	0	1	0	0.42
Ne	10	20.1797	2	18	9	0	1	0	0.38
Na	11	22.989768	3	1	1	1	0	0	1.9
Mg	12	24.305	3	2	2	1	0	0	1.45
Al	13	26.981539	3	13	5	1	0	0	1.18
Si	14	28.0855	3	14	6	0	1	0	1.11

Unique formulae end up with unique vector representations

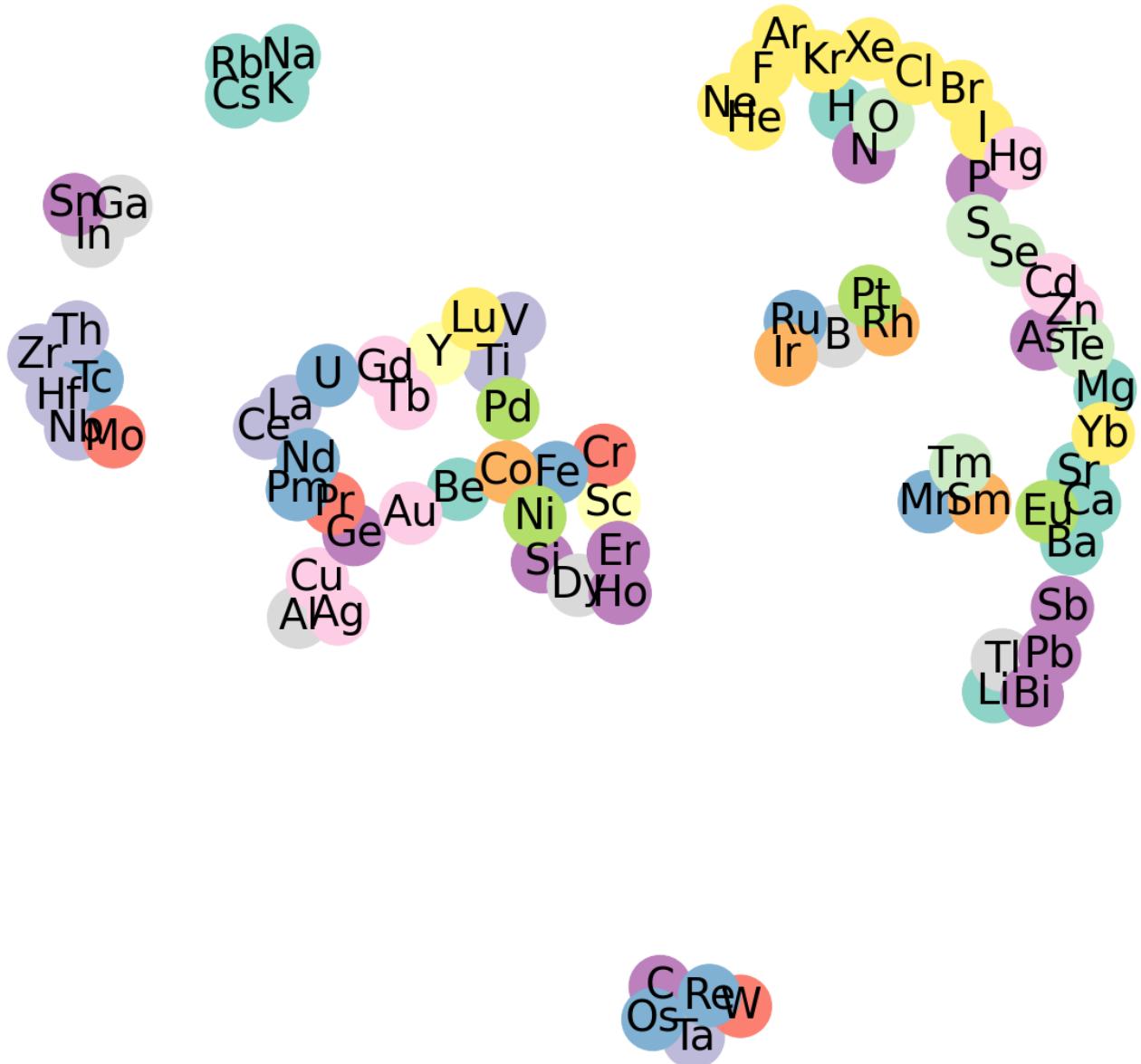


C	6	12.011	2	14	7	0	1	0	0.67
N	7	14.00674	2	15	7	0	1	0	0.56
O	8	15.9994	2	16	7	0	1	0	0.48
F	9	18.998403	2	17	8	0	1	0	0.42
Ne	10	20.1797	2	18	9	0	1	0	0.38
Na	11	22.989768	3	1	1	1	0	0	1.9
Mg	12	24.305	3	2	2	1	0	0	1.45
Al	13	26.981539	3	13	5	1	0	0	1.18

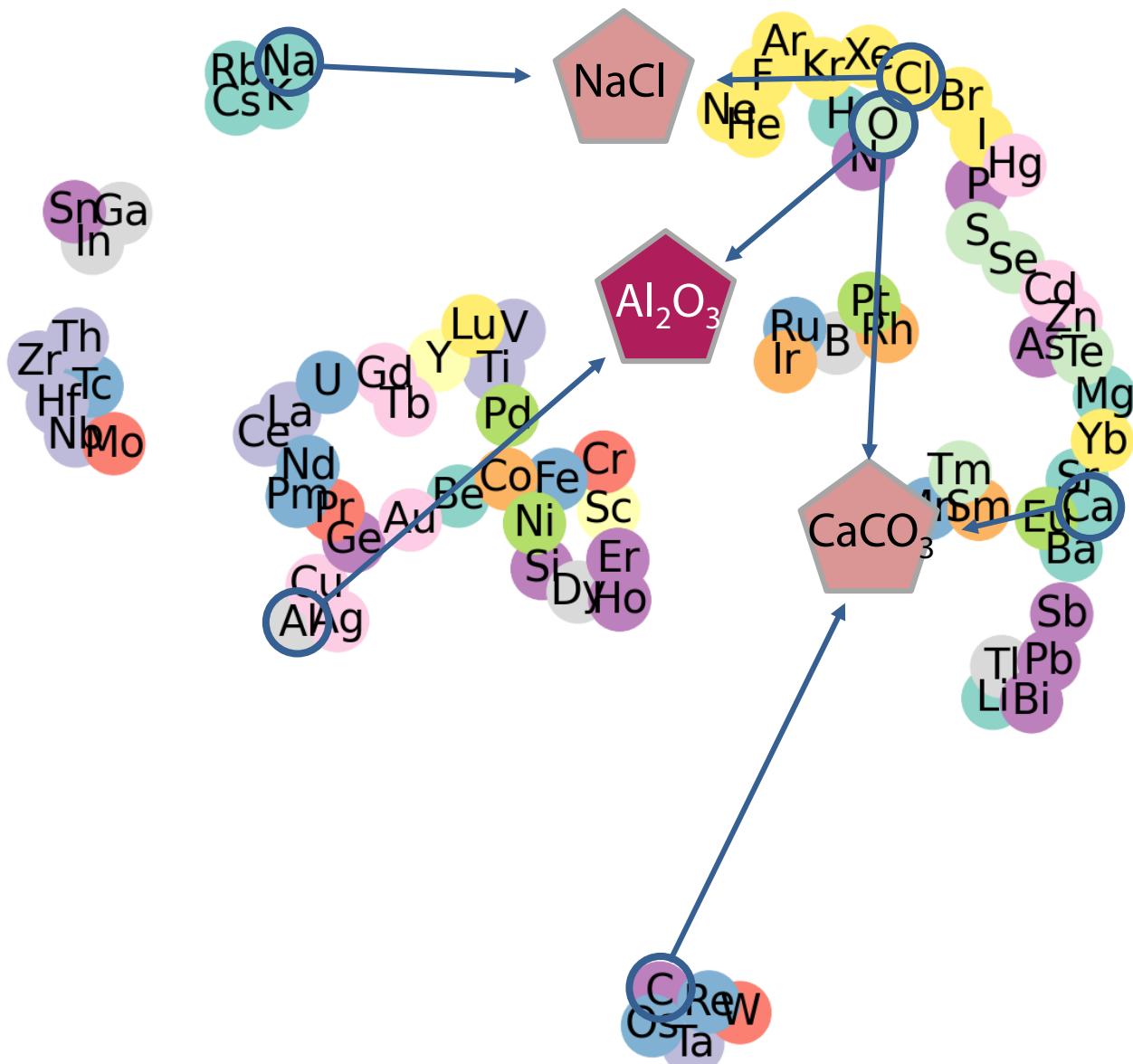
$$\begin{aligned}\text{Al}_{0.4}\text{O}_{0.6} = & [8 \times 0.6 + 13 \times 0.4, \\ & 16 \times 0.6 + 27 \times 0.4, \\ & \dots, \\ & 0.48 \times 0.6 + 1.18 \times 0.4]\end{aligned}$$

$$\begin{aligned}\text{Al}_{0.4}\text{O}_{0.6} = & [10, \\ & 20.4, \\ & \dots, \\ & 3.35]\end{aligned}$$

Projecting this CBFV into reduced dimensions we see cluster of chemically similar materials

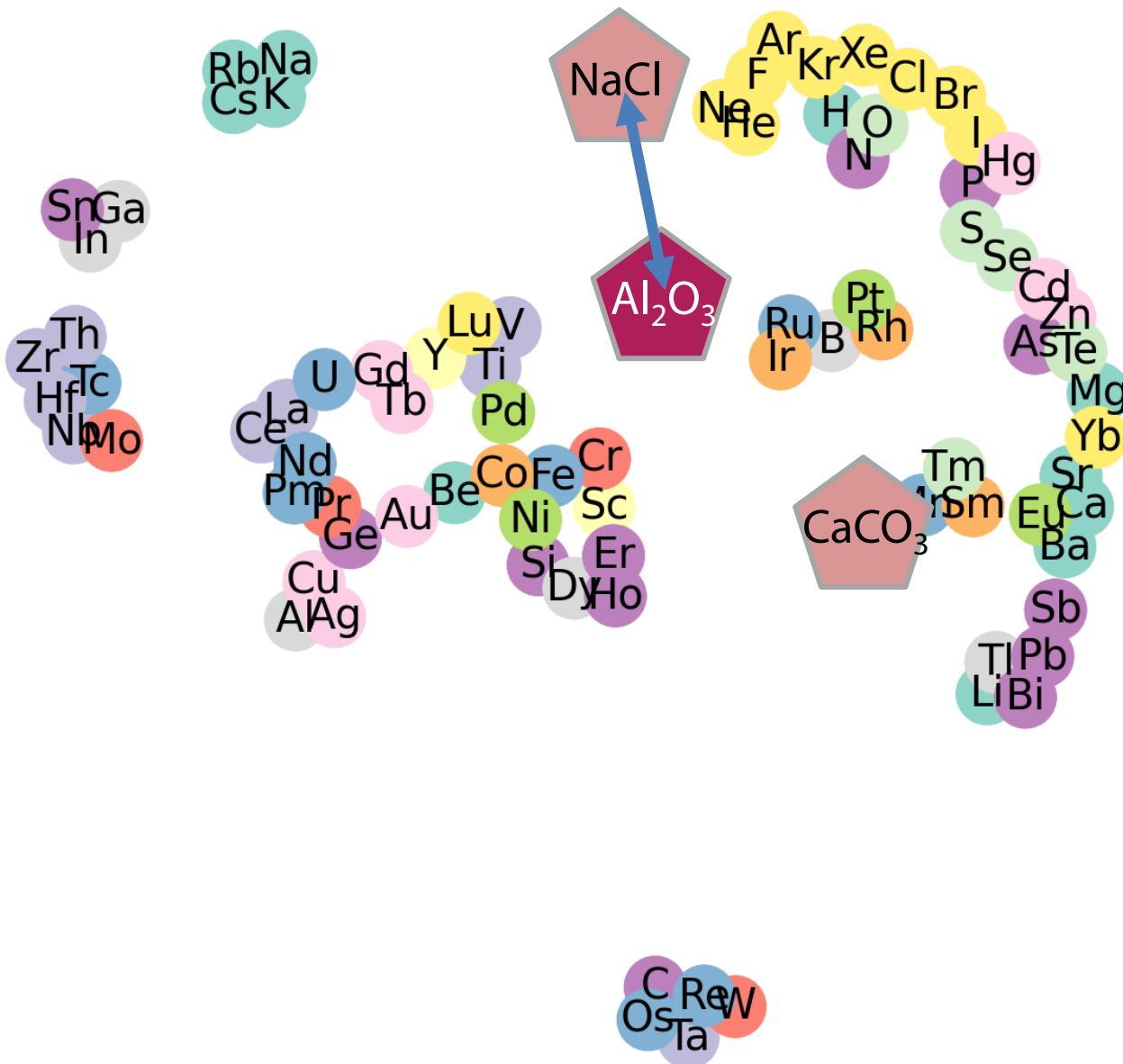


Projecting this CBFV into reduced dimensions we see cluster of chemically similar materials



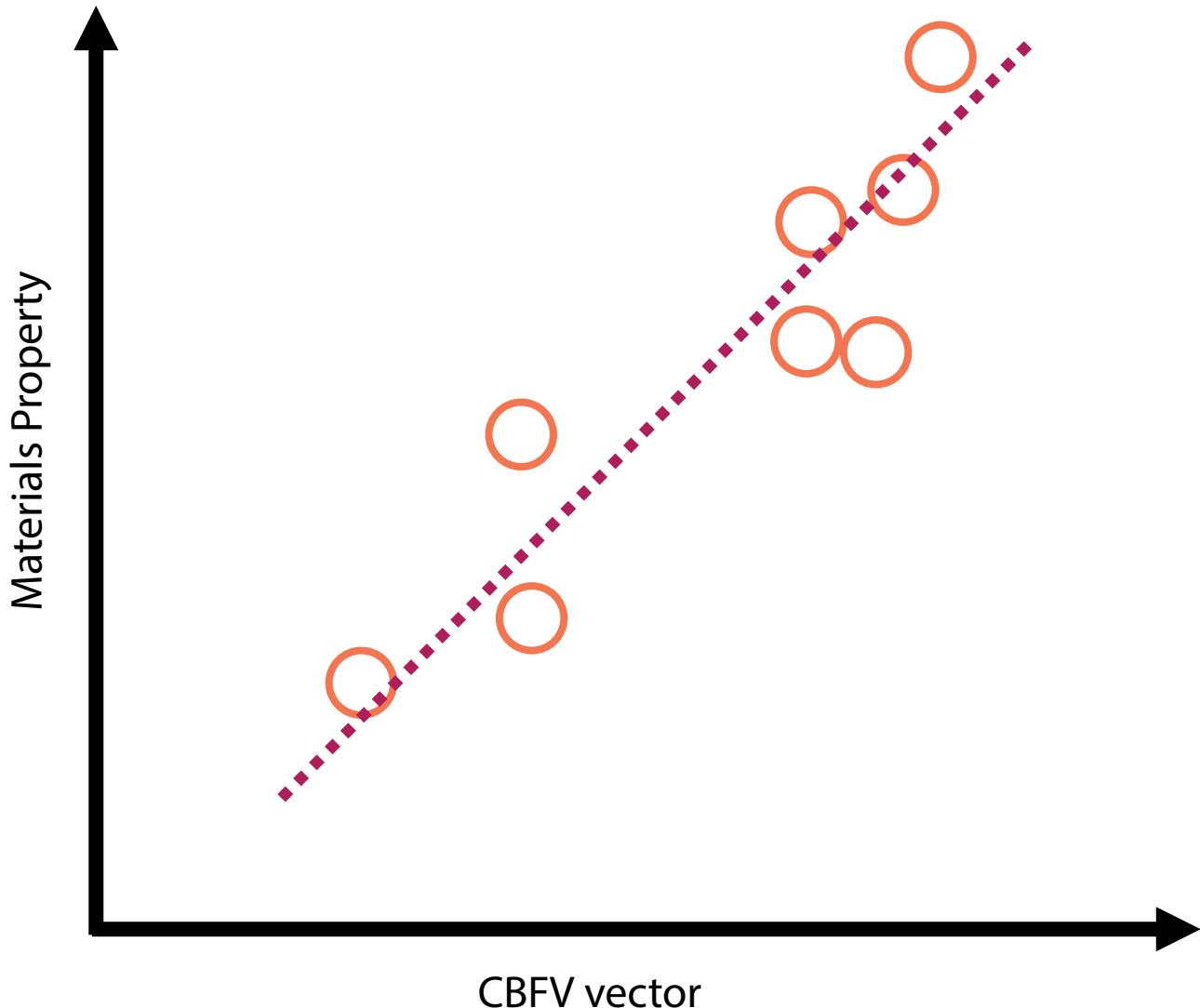
Compounds appear at intermediate positions based on pure elemental constituents

Compounds near to one another should share similar properties



Nearest Neighbor:  
Find closest datapoint  
that already has **label**

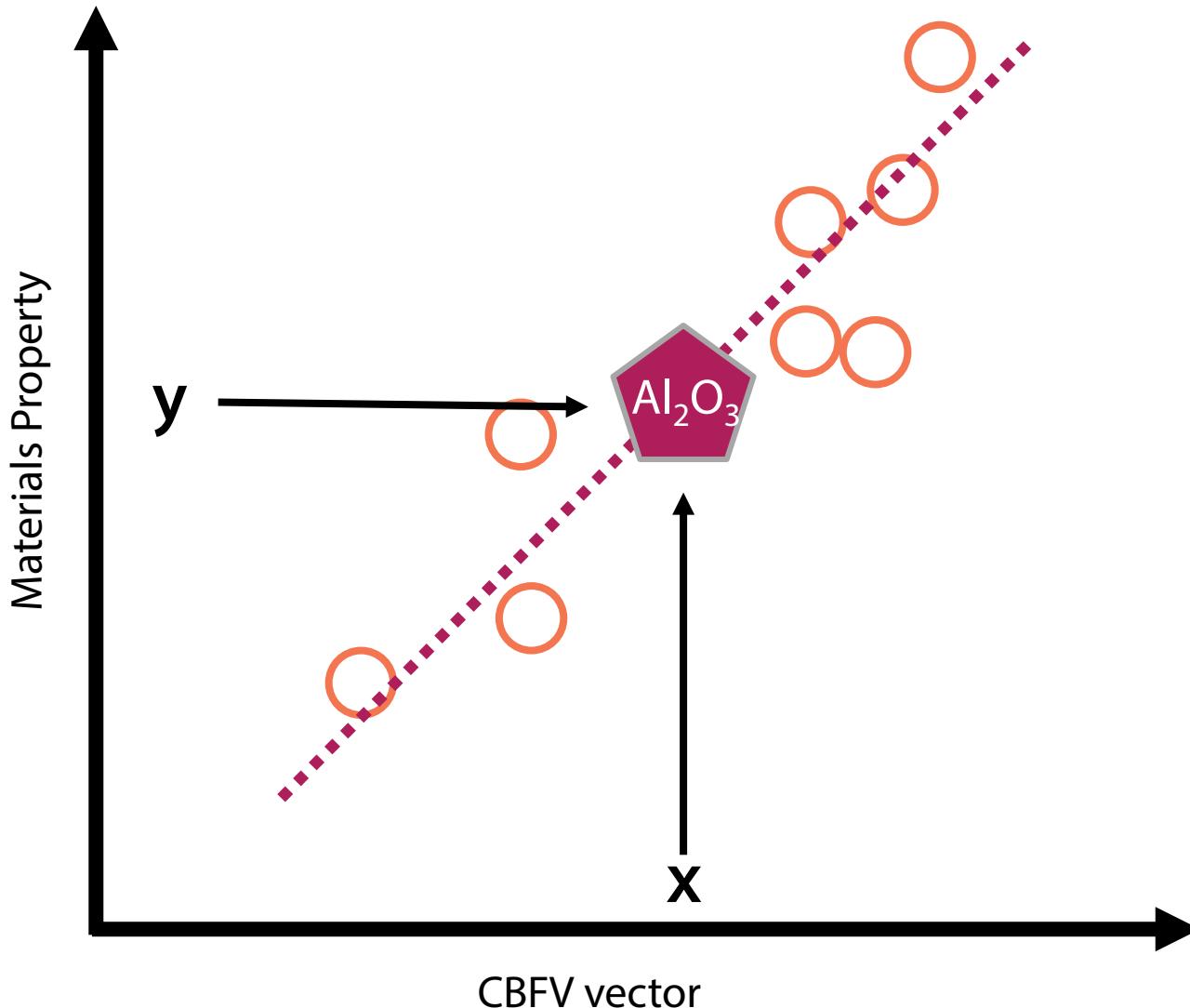
Many different models and algorithms available for prediction



$$y = mx + b$$

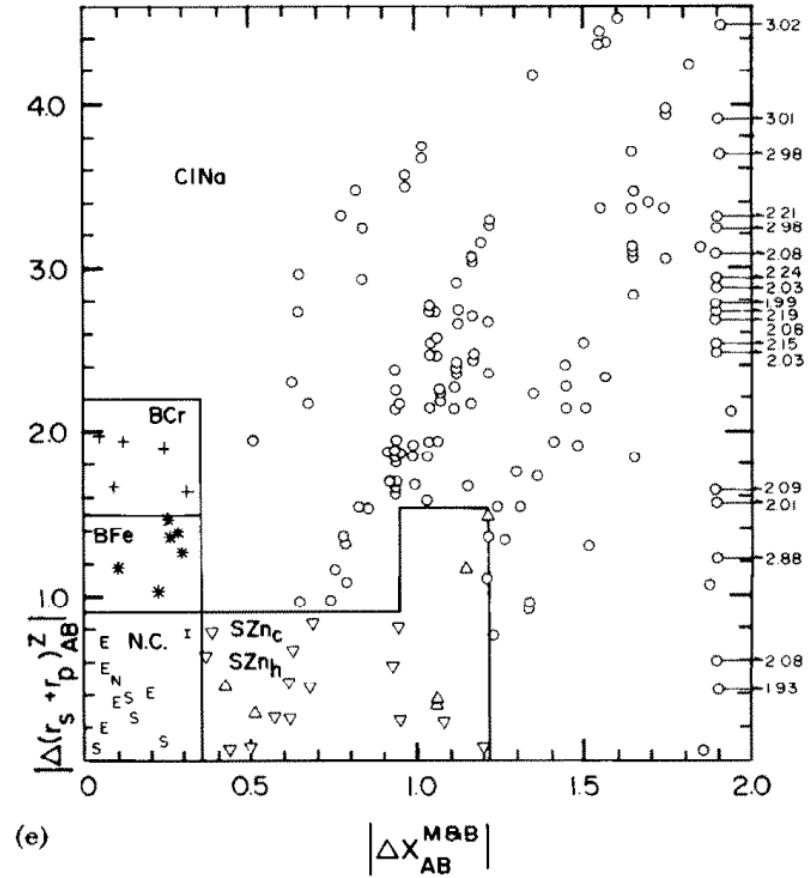
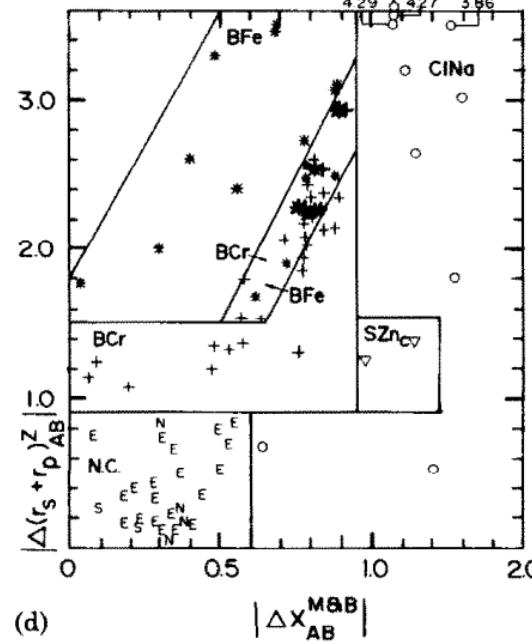
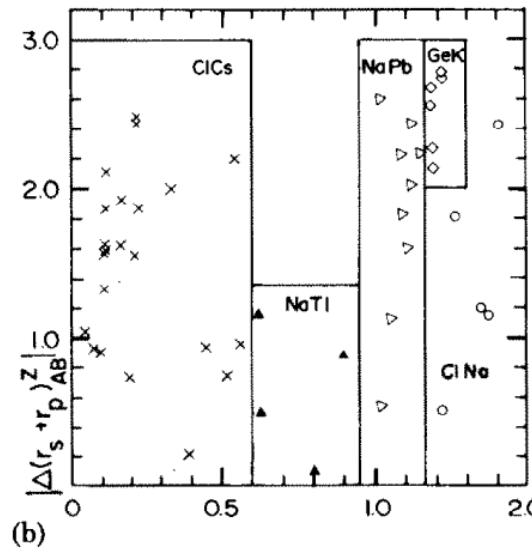
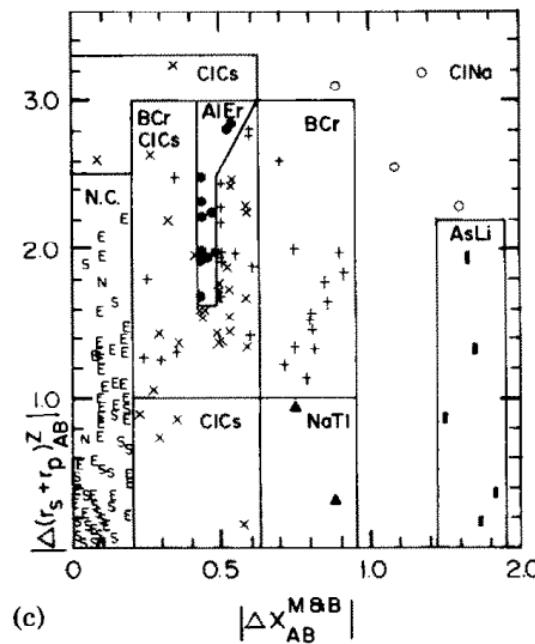
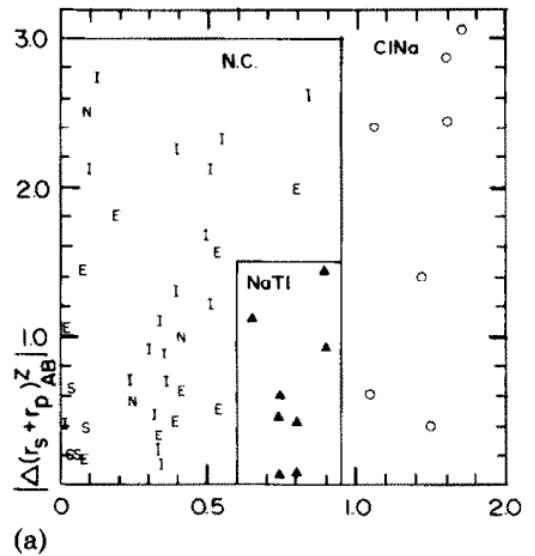
Fit line to  
known data

Many different models and algorithms available for prediction



- $$y = mx + b$$
1. Given a new  $x$
  2. Generate a prediction  $y$

# Villars developed chemical descriptors long before materials informatics!



P. Villars, 1983, Journal of  
Less Common Metals

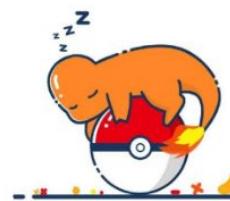
There are many different elemental property databases used for CBFV



# Oliynyk

## idocx/Atom2Vec

A python implement of Atom2Vec: a simple way to describe atoms for machine learning



materialsintelligence/  
**mat2vec**

Supplementary Materials for Tshitoyan et al.  
"Unsupervised word embeddings capture latent knowledge from materials science literature",  
Nature (2019).



# MAGPIE

Materials Agnostic Platform for Informatics and Exploration

JARVIS, MAGPIE, and Oliynyk are all based off of elemental property databases

# Periodic Table of Elements

[Cite](#)

[Download](#)



TABLE

LIST W/PROPERTIES

GAME

1 H Hydrogen Nonmetal	3 Li Lithium Alkali Metal	4 Be Beryllium Alkaline Earth M	1 H Hydrogen Nonmetal
11 Na Sodium Alkali Metal	12 Mg Magnesium Alkaline Earth M		
19 K Potassium Alkali Metal	20 Ca Calcium Alkaline Earth M	21 Sc Scandium Transition Metal	22 Ti Titanium Transition Metal
37 Rb Rubidium Alkali Metal	38 Sr Strontium Alkaline Earth M	39 Y Yttrium Transition Metal	40 Zr Zirconium Transition Metal
55 Cs Cesium Alkali Metal	56 Ba Barium Alkaline Earth M	*	72 Hf Hafnium Transition Metal
87 Fr Francium Alkali Metal	88 Ra Radium Alkaline Earth M	**	73 Ta Tantalum Transition Metal
			74 W Tungsten Transition Metal
		*	104 Rf Rutherfordium Transition Metal
		**	105 Db Dubnium Transition Metal
			106 Sg Seaborgium Transition Metal
		*	57 La Lanthanum Lanthanide
		**	58 Ce Cerium Lanthanide
			59 Pr Praseodymium Lanthanide
			89 Ac Actinium Actinide
			90 Th Thorium Actinide
			91 Pa Protactinium Actinide
			92 U Uranium Actinide
			93 Np Neptunium Actinide
			94 Pu Plutonium Actinide
			95 Am Americium Actinide
			96 Cm Curium Actinide
			97 Bk Berkelium Actinide
			98 Cf Californium Actinide
			99 Es Einsteinium Actinide
			100 Fm Fermium Actinide
			101 Md Mendelevium Actinide
			102 No Nobelium Actinide
			103 Lr Lawrencium Actinide

## ELEMENT PROPERTIES

29  
**Cu**  
Copper  
Transition Metal

[Copper Element Page](#)

## DISPLAY PROPERTY/TREND

### Chemical Group Block

5 B Boron Metalloid	6 C Carbon Nonmetal	7 N Nitrogen Nonmetal	8 O Oxygen Nonmetal	9 F Fluorine Halogen	2 He Helium Noble Gas
13 Al Aluminum Transition I	14 Si Silicon Metalloid	15 P Phosphorus Nonmetal	16 S Sulfur Nonmetal	17 Cl Chlorine Halogen	18 Ar Argon Noble Gas
31 Ga Gallium Transition I	32 Ge Germanium Metalloid	33 As Arsenic Metalloid	34 Se Selenium Nonmetal	35 Br Bromine Halogen	36 Kr Krypton Noble Gas
49 Tn Thorium Transition I	50 Sn Tin Post-Transition I	51 Sb Antimony Metalloid	52 Te Tellurium Metalloid	53 I Iodine Halogen	54 Xe Xenon Noble Gas
81 Tl Thallium Transition I	82 Pb Lead Post-Transition I	83 Bi Bismuth Post-Transition I	84 Po Polonium Metalloid	85 At Astatine Halogen	86 Rn Radon Noble Gas
113 Nh Nhrium Transition I	114 Fl Flerovium Post-Transition I	115 Mc Moscovium Post-Transition I	116 Lv Livermorium Post-Transition I	117 Ts Tennessine Halogen	118 Og Oganesson Noble Gas
56 Dy Dysprosium Lanthanide	67 Ho Holmium Lanthanide	68 Er Erbium Lanthanide	69 Tm Thulium Lanthanide	70 Yb Ytterbium Lanthanide	71 Lu Lutetium Lanthanide
98 Cf Californium Actinide	99 Es Einsteinium Actinide	100 Fm Fermium Actinide	101 Md Mendelevium Actinide	102 No Nobelium Actinide	103 Lr Lawrencium Actinide

# mat2vec extracts chemical domain knowledge from journal abstracts via word embeddings

## LETTER

<https://doi.org/10.1038/s41586-019-1335-8>

### Unsupervised word embeddings capture latent knowledge from materials science literature

Vahé Tshitoyan<sup>1,3\*</sup>, John Dagdelen<sup>1,2</sup>, Leigh Weston<sup>1</sup>, Alexander Dunn<sup>1,2</sup>, Ziqin Rong<sup>1</sup>, Olga Kononova<sup>2</sup>, Kristin A. Persson<sup>1,2</sup>, Gerbrand Ceder<sup>1,2,\*</sup> & Anubhav Jain<sup>1\*</sup>

The overwhelming majority of scientific knowledge is published as text, which is difficult to analyse by either traditional statistical analysis or modern machine learning methods. By contrast, the main source of machine-interpretable data for the materials research community has come from structured property databases<sup>1,2</sup>, which encompass only a small fraction of the knowledge present in the research literature. Beyond property values, publications contain valuable knowledge regarding the connections and relationships between data items as interpreted by the authors. To improve the identification and use of this knowledge, several studies have focused on the retrieval of information from scientific literature using supervised natural language processing<sup>3–10</sup>, which requires large hand-labelled datasets for training. Here we show that materials science knowledge present in the published literature can be efficiently encoded as information-dense word embeddings<sup>11–13</sup> (vector representations of words) without human labelling or supervision. Without any explicit insertion of chemical knowledge, these embeddings capture complex materials science concepts such as the underlying structure of the periodic table and structure–property relationships in materials. Furthermore, we demonstrate that an unsupervised method can recommend materials for functional applications several years before their discovery. This suggests that latent knowledge regarding future discoveries is to a large extent embedded in past publications. Our findings highlight the possibility of extracting knowledge and relationships from the massive body of scientific literature in a collective manner, and point towards a generalized approach to the mining of scientific literature.

Assignment of high-dimensional vectors (embeddings) to words in a text corpus in a way that preserves their syntactic and semantic relationships is one of the most fundamental techniques in natural language processing (NLP). Word embeddings are usually constructed using machine learning algorithms such as GloVe<sup>13</sup> or Word2vec<sup>11,12</sup>, which use information about the co-occurrences of words in a text corpus. For example, when trained on a suitable body of text, such methods should produce a vector representing the word 'iron' that is closer by cosine distance to the vector for 'steel' than to the vector for 'organic'. To train the embeddings, we collected and processed approximately 3.3 million scientific abstracts published between 1922 and 2018 in more than 1,000 journals deemed likely to contain materials-related research, resulting in a vocabulary of approximately 500,000 words. We then applied the skip-gram variation of Word2vec, which is trained to predict context words that appear in the proximity of the target word as a means to learn the 200-dimensional embedding of that target word, to our text corpus (Fig. 1a). The key idea is that, because words with similar meanings often appear in similar contexts, the corresponding embeddings will also be similar. More details about the model are included in the Methods and in Supplementary Information sections S1 and S2, where we also discuss alternative algorithm options such as GloVe. We find that, even though no chemical information or interpretation is added to the algorithm, the obtained word embeddings

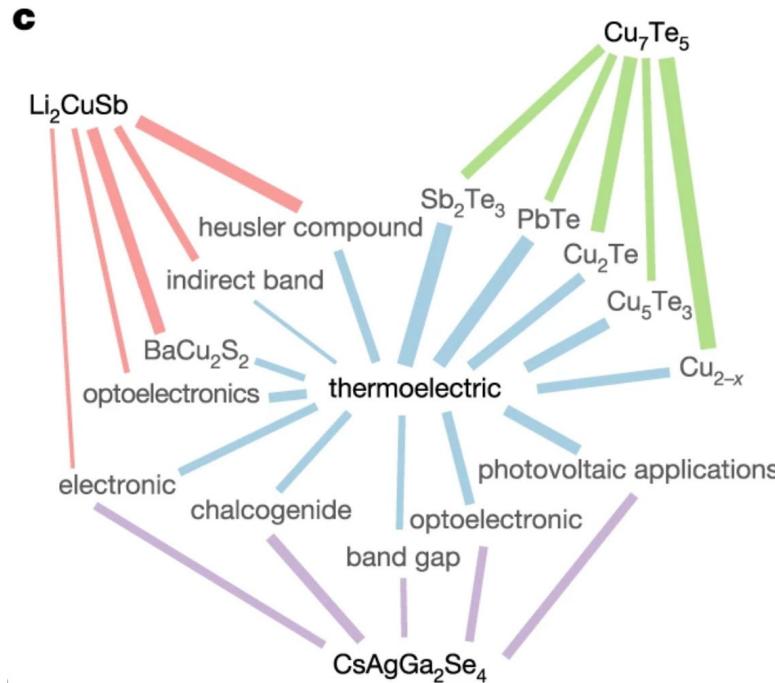
behave consistently with chemical intuition when they are combined using various vector operations (projection, addition, subtraction). For example, many words in our corpus represent chemical compositions of materials, and the five materials most similar to  $\text{LiCoO}_2$  (a well-known lithium-ion cathode compound) can be determined through a dot product (projection) of normalized word embeddings. According to our model, the compositions with the highest similarity to  $\text{LiCoO}_2$  are  $\text{LiMn}_2\text{O}_4$ ,  $\text{LiNi}_{0.5}\text{Mn}_{1.5}\text{O}_4$ ,  $\text{LiNi}_{0.8}\text{Co}_{0.2}\text{O}_2$ ,  $\text{LiNi}_{0.8}\text{Co}_{0.15}\text{Al}_{0.05}\text{O}_2$  and  $\text{LiNiO}_2$ —all of which are also lithium-ion cathode materials.

Similar to the observation made in the original Word2vec paper<sup>11</sup>, these embeddings also support analogies, which in our case can be domain-specific. For instance, 'NiFe' is to 'ferromagnetic' as 'IrMn' is to '?', where the most appropriate response is 'antiferromagnetic'. Such analogies are expressed and solved in the Word2vec model by finding the nearest word to the result of subtraction and addition operations between the embeddings. Hence, in our model,

ferromagnetic– $\text{NiFe} + \text{IrMn} \approx$  antiferromagnetic

To better visualize such embedded relationships, we projected the embeddings of  $\text{Zr}$ ,  $\text{Cr}$  and  $\text{Ni}$ , as well as their corresponding oxides and crystal structures, onto two dimensions using principal component analysis (Fig. 1b). Even in reduced dimensions, there is a consistent operation in vector space for the concepts 'oxide of' ( $\text{Zr} - \text{ZrO}_2 \approx \text{Cr} - \text{Cr}_2\text{O}_3 \approx \text{Ni} - \text{NiO}$ ) and 'structure of' ( $\text{Zr} - \text{HCP} \approx \text{Cr} - \text{BCC} \approx \text{Ni} - \text{FCC}$ ). This suggests that the positions of the embeddings in space encode materials science knowledge such as the fact that zirconium has a hexagonal close packed (HCP) crystal structure under standard conditions and that its principal oxide is  $\text{ZrO}_2$ . Other types of materials analogies captured by the model, such as functional applications and crystal symmetries, are listed in Extended Data Table 1. The accuracies for each category are close to 50%—similar to the baseline set in the original Word2vec study<sup>12</sup>. We stress that Word2vec treats these entities simply as strings, and no chemical interpretation is explicitly provided to the model; rather, materials knowledge is captured through the positions of the words in scientific abstracts. Notably, we also found that embeddings of chemical elements are representative of their positions in the periodic table when projected onto two dimensions (Extended Data Fig. 1a, b; Supplementary Information sections S4 and S5) and can serve as effective feature vectors in quantitative machine learning models such as formation energy prediction—outperforming several previously reported curated feature vectors (Extended Data Fig. 1c, d; Supplementary Information section S6).

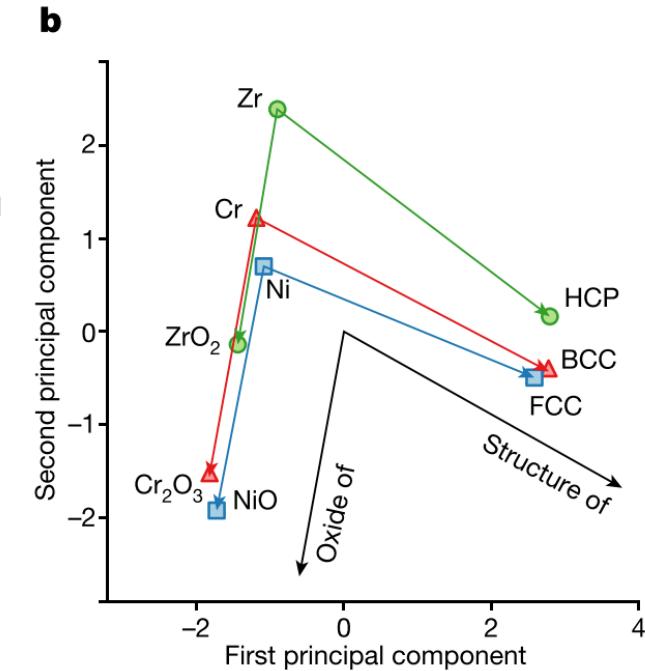
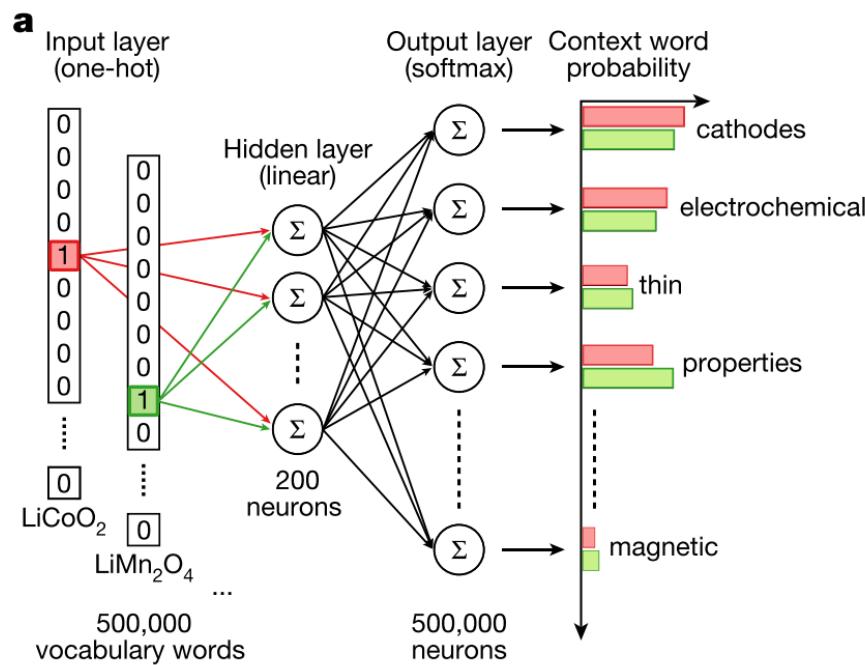
The main advantage and novelty of this representation, however, is that application keywords such as 'thermoelectric' have the same representation as material formulae such as ' $\text{Bi}_2\text{Te}_3$ '. When the cosine similarity of a material embedding and the embedding of 'thermoelectric' is high, one might expect that the text corpus necessarily includes abstracts reporting on the thermoelectric behaviour of this material<sup>14,15</sup>. However, we found that a number of materials that have relatively high cosine similarities to the word 'thermoelectric' never



**ABSTRACT**  $\text{Cu}_{9.1}\text{Te}_4\text{Cl}_3$  is a new polymorphic compound in the class of coinage metal polytelluride halides. Copper is highly mobile, which results in multiple order–disorder phase transitions in a limited temperature interval from 240 to 370 K. Mainly as a consequence of thermal transport properties, the compound's thermoelectric figure of merit reaches values up to  $ZT = 0.15$  in the temperature range between room temperature and 523 K. Its structure is closely related to that of  $\text{Ag}_{10}\text{Te}_4\text{Br}_3$ , another coinage metal polytelluride halide, which represents the first p–n–p-switchable semiconductor approachable by a simple temperature change. The title compound outperforms  $\text{Ag}_{10}\text{Te}_4\text{Br}_3$  in terms of thermoelectric properties by 1 order of magnitude and therefore acts as a link between the class of p–n–p compounds and thermoelectric

<sup>1</sup>Lawrence Berkeley National Laboratory, Berkeley, CA, USA. <sup>2</sup>Department of Materials Science and Engineering, University of California, Berkeley, CA, USA. <sup>3</sup>Present address: Google LLC, Mountain View, CA, USA. \*e-mail: vahé.tshitoyan@gmail.com; ceder@lbl.gov; ajain@lbl.gov

# mat2vec extracts chemical domain knowledge from journal abstracts via word embeddings



**Fig. 1 | Word2vec skip-gram and analogies.** **a**, Target words ‘LiCoO<sub>2</sub>’ and ‘LiMn<sub>2</sub>O<sub>4</sub>’ are represented as vectors with ones at their corresponding vocabulary indices (for example, 5 and 8 in the schematic) and zeros everywhere else (one-hot encoding). These one-hot encoded vectors are used as inputs for a neural network with a single linear hidden layer (for example, 200 neurons), which is trained to predict all words mentioned within a certain distance (context words) from the given target word. For similar battery cathode materials such as LiCoO<sub>2</sub> and LiMn<sub>2</sub>O<sub>4</sub>, the context words that occur in the text are mostly the same (for example,

‘cathodes’, ‘electrochemical’, and so on), which leads to similar hidden layer weights after the training is complete. These hidden layer weights are the actual word embeddings. The softmax function is used at the output layer to normalize the probabilities. **b**, Word embeddings for Zr, Cr and Ni, their principal oxides and crystal symmetries (at standard conditions) projected onto two dimensions using principal component analysis and represented as points in space. The relative positioning of the words encodes materials science relationships, such that there exist consistent vector operations between words that represent concepts such as ‘oxide of’ and ‘structure of’.

# Atom2Vec creates a vector based on known chemical ratios

## Learning atoms for materials discovery

Quan Zhou<sup>a</sup>, Peizhe Tang<sup>a</sup>, Shenxiu Liu<sup>a</sup>, Jinbo Pan<sup>b</sup>, Qimin Yan<sup>b</sup>, and Shou-Cheng Zhang<sup>a,c,1</sup>

<sup>a</sup>Department of Physics, Stanford University, Stanford, CA 94305-4045; <sup>b</sup>Department of Physics, Temple University, Philadelphia, PA 19122; and <sup>c</sup>Stanford Institute for Materials and Energy Sciences, SLAC National Accelerator Laboratory, Menlo Park, CA 94025

Contributed by Shou-Cheng Zhang, June 4, 2018 (sent for review February 2, 2018; reviewed by Xi Dai and Stuart P. Parkin)

Exciting advances have been made in artificial intelligence (AI) during recent decades. Among them, applications of machine learning (ML) and deep learning techniques brought human-competitive performances in various tasks of fields, including image recognition, speech recognition, and natural language understanding. Even in Go, the ancient game of profound complexity, the AI player has already beaten human world champions convincingly with and without learning from the human. In this work, we show that our unsupervised machines (Atom2Vec) can learn the basic properties of atoms by themselves from the extensive database of known compounds and materials. These learned properties are represented in terms of high-dimensional vectors, and clustering of atoms in vector space classifies them into meaningful groups consistent with human knowledge. We use the atom vectors as basic input units for neural networks and other ML models designed and trained to predict materials properties, which demonstrate significant accuracy.

atomism | machine learning | materials discovery

The past 20 y witnessed the accumulation of an unprecedentedly massive amount of data in materials science via both experimental explorations and numerical simulations (1–5). The huge datasets not only enable but also call for data-based statistical approaches. As a result, a new paradigm has emerged which aims to harness artificial intelligence (AI) and machine-learning (ML) techniques (6–10) to assist materials research and discovery. Several initial attempts have been made along this path (11–16). Most of them learned maps from materials information (input) to materials properties (output) based on known materials properties. The input or feature of materials involves descriptors of constituents: Certain physical or chemical attributes of atoms are taken, depending on the materials property under prediction (11, 14, 17). Despite the success so far, these works heavily rely on researchers' wise selection of relevant descriptors; thus the degree of intelligence is still very limited from a theoretical perspective. And practically, extra computations are usually unavoidable for machines to interpret such atom descriptors which are in the form of abstract human knowledge.

To create a higher level of AI and to overcome the practical limitation, we propose Atom2Vec in this paper, which lets machines learn their own knowledge about atoms from data. Atom2Vec considers only existence of compounds in a materials database, without reference to any specific property of materials. This massive dataset is leveraged for a learning feature in materials science in an unsupervised manner (18–23). Because of the absence of materials property labels, Atom2Vec is naturally prevented from being biased to one certain aspect. As a result, the learned knowledge can yield complete and universal descriptions of atoms in principle, as long as the dataset is sufficiently large and representative. Atom2Vec follows the core idea that properties of an atom can be inferred from the environments it lives in, which is similar to the distributional hypothesis in linguistics (24). In a compound, each atom can be selected as a target type, while the environment refers to all remaining atoms together with their positions relative to the target atom. Intuitively, similar atoms tend to appear in similar environments, which allows our Atom2Vec to extract knowledge from the associations between



PHYSICS

### Atom2Vec Workflow

We begin to illustrate the full workflow of Atom2Vec as shown in Fig. 1. To capture relations between atoms and environments, atom–environment pairs are generated for each compound in a materials dataset as the first step. Before pair generation, a more explicit definition of environment is needed, whereas atoms are represented by chemical symbols conveniently. Although a complete environment should involve both chemical composition and crystal structure as mentioned before, we take into account only the former here as a proof of concept. Under this simplification, environment covers two aspects: the count of the target atom in the compound and the counts of different atoms in the remains. As an example, let us consider the compound  $\text{Bi}_2\text{Se}_3$  from the mini-dataset of only seven samples given in Fig. 1. Two atom–environment pairs are generated from  $\text{Bi}_2\text{Se}_3$ : for atom  $\text{Bi}$ , the environment is represented as “(2) $\text{Se}_3$ ”; for atom  $\text{Se}$ , the environment is represented as “(3) $\text{Bi}_2$ .” Specifically, for the first pair, “(2)” in the environment (2) $\text{Se}_3$  means there are two target atoms ( $\text{Bi}$  here for the compound), while “ $\text{Se}_3$ ” indicates that three  $\text{Se}$  atoms exist in the environment. Following the notation, we collect all atom–environment pairs from the dataset and then record them in an atom–environment matrix  $X$ , where its entry  $X_{ij}$  gives the count of pairs with the  $i$ th atom and the  $j$ th environment. Such a matrix for the mini-dataset is also given in Fig. 1 for illustration purposes. Clearly, each row vector gives counts with different environments for one atom, and each column vector yields counts with different atoms for one environment. According to the previously mentioned intuition, two atoms behave similarly if their corresponding row vectors are close to each other in the vector space.

Although revealing similarity to some extent, descriptions of atoms in terms of row vectors of the atom–environment matrix are still very primitive and inefficient, since the vectors can be

### Significance

Motivated by the recent achievements of artificial intelligence (AI) in linguistics, we design AI to learn properties of atoms from materials data on its own. Our work realizes knowledge representation of atoms via computers and could serve as a foundational step toward materials discovery and design fully based on machine learning.

Author contributions: Q.Z., P.T., S.L., Q.Y., and S.-C.Z. designed research; Q.Z., P.T., and S.L. performed research; Q.Z., P.T., and S.L. contributed new reagents/analytic tools; Q.Z., P.T., S.L., J.P., and Q.Y. analyzed data; and Q.Z. and P.T. wrote the paper.

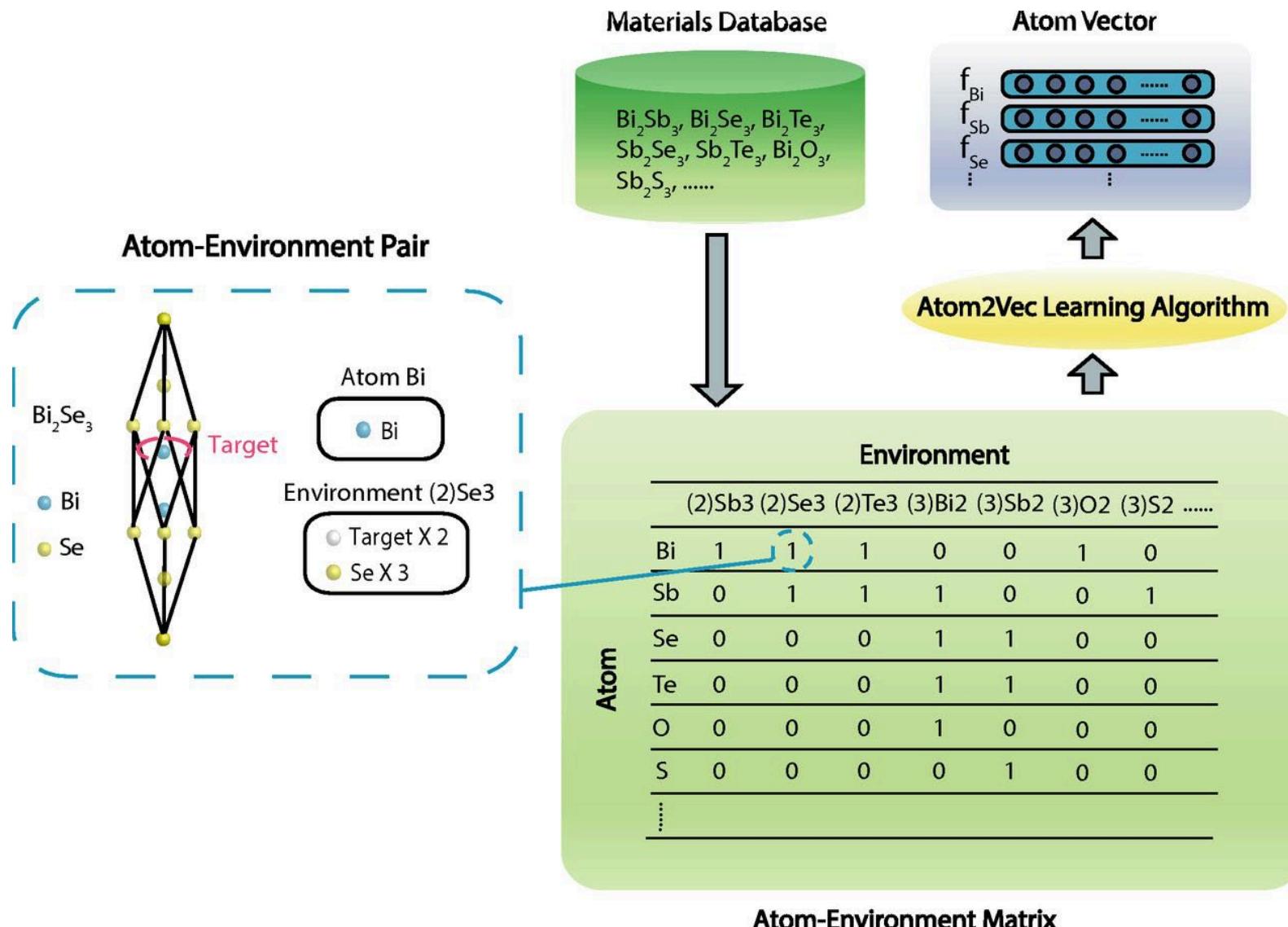
Reviewers: X.D., Institute of Physics, Chinese Academy of Sciences; and S.P.P., Max Planck Institute of Microstructure Physics in Halle, Martin-Luther-University Halle-Wittenberg. The authors declare no conflict of interest.

This open access article is distributed under Creative Commons Attribution-NonCommercial-NoDerivatives License 4.0 (CC BY-NC-ND).

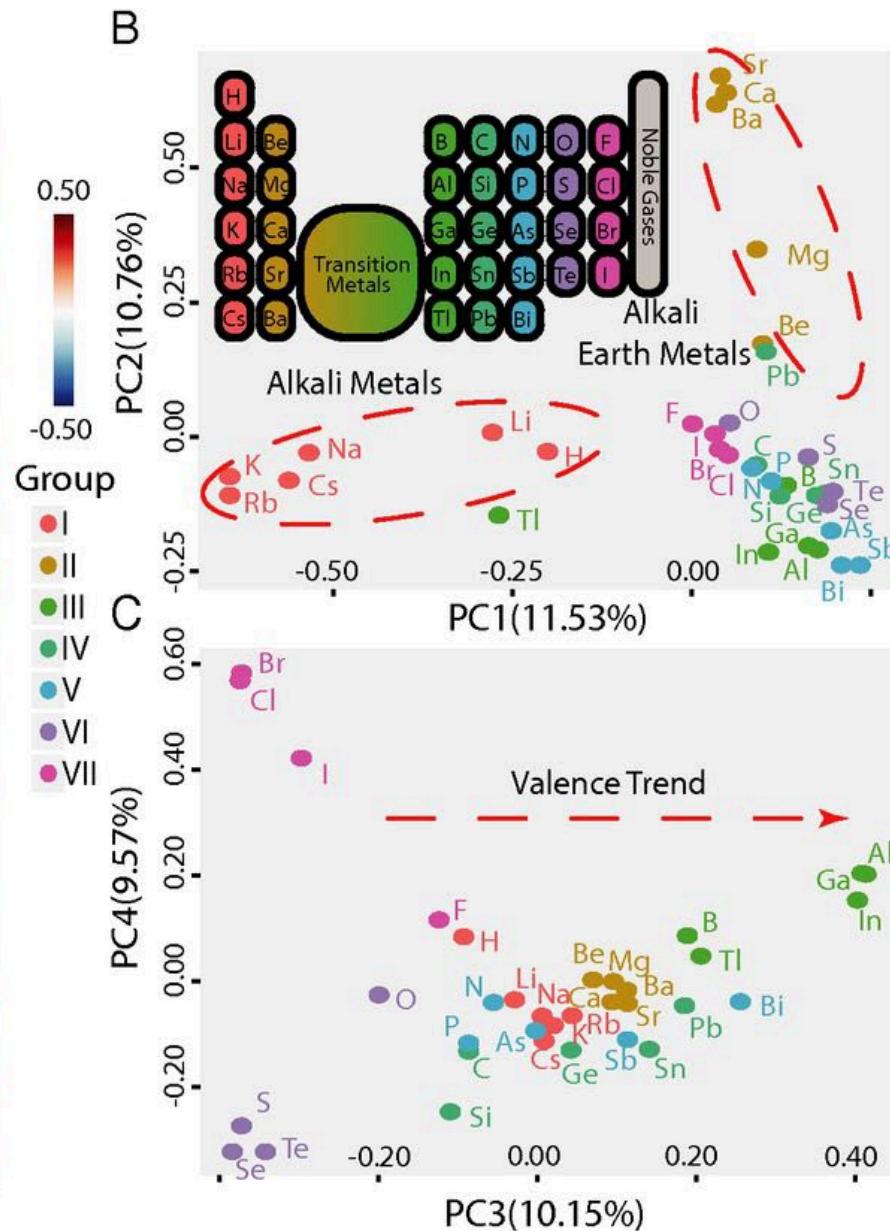
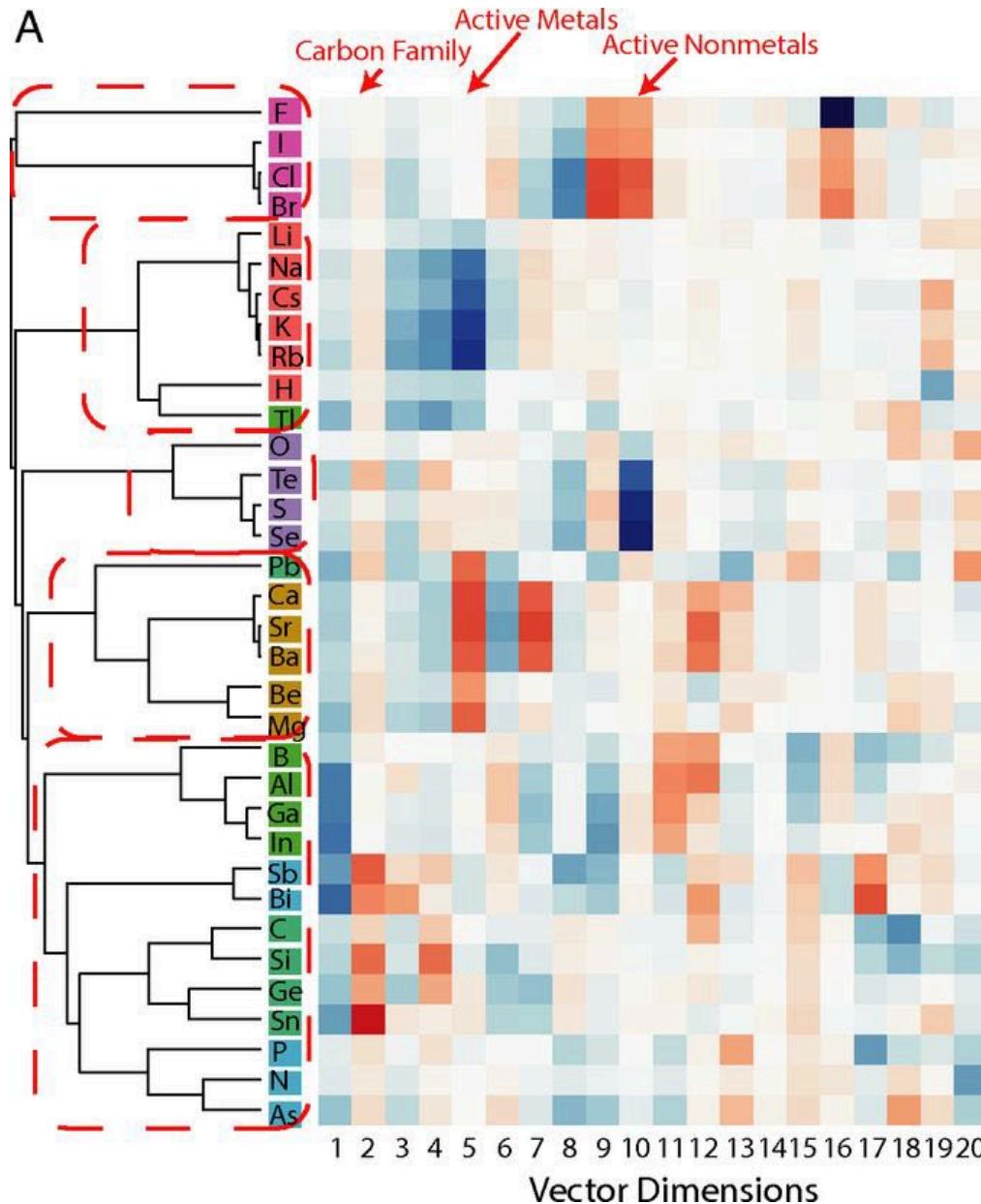
To whom correspondence should be addressed. Email: schang@stanford.edu.

This article contains supporting information online at [www.pnas.org/lookup/suppl/doi:10.1073/pnas.1801181115/DCSupplemental](http://www.pnas.org/lookup/suppl/doi:10.1073/pnas.1801181115/DCSupplemental).

Published online June 26, 2018.



# Atom2Vec creates a vector based on known chemical ratios



# ElemNet is a simple one-hot encoding feature paired with deep learning

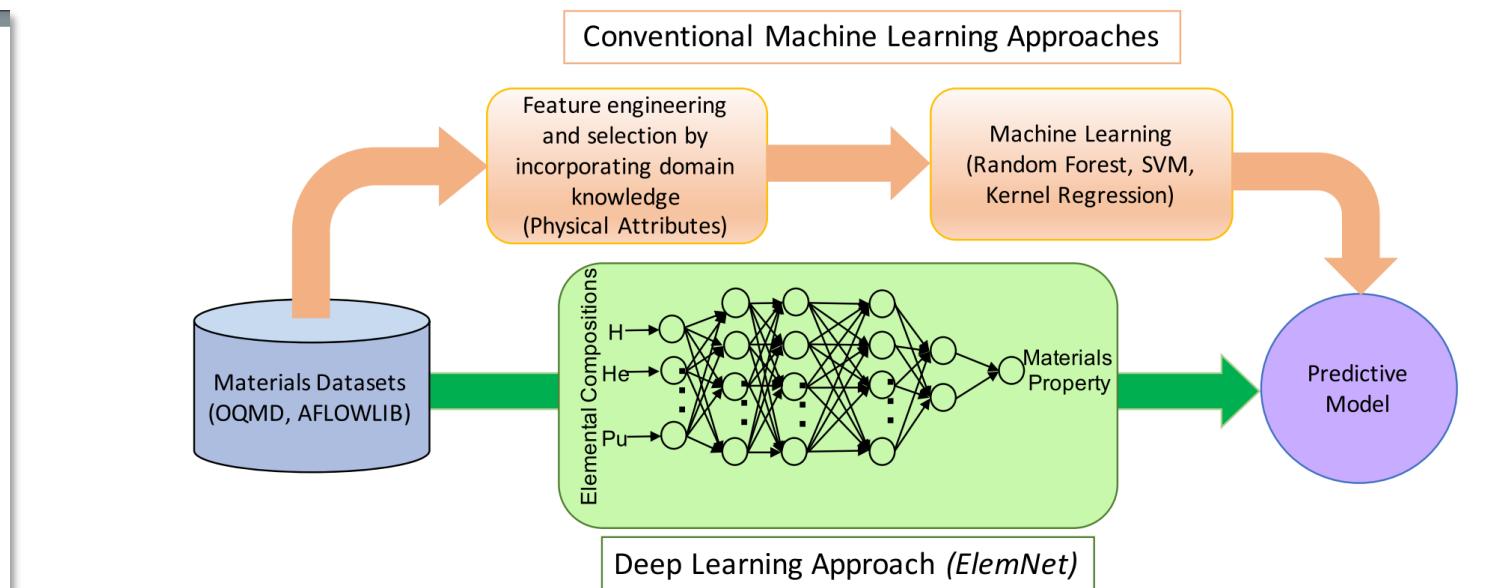
The image shows the cover page of a scientific publication. At the top right is the URL [www.nature.com/scientificreports/](http://www.nature.com/scientificreports/). Below it is the journal title "SCIENTIFIC REPORTS" in large, bold, black letters, with a red gear icon integrated into the letter "R". To the left of the title is the word "OPEN". The main title of the paper is "ElemNet: Deep Learning the Chemistry of Materials From Only Elemental Composition". Below the title, the authors are listed: Dipendra Jha<sup>1</sup>, Logan Ward<sup>2</sup>, Arindam Paul<sup>1</sup>, Wei-keng Liao<sup>1</sup>, Alok Choudhary<sup>1</sup>, Chris Wolverton<sup>3</sup> & Ankit Agrawal<sup>1</sup>. The text indicates the following publication details: Received: 1 August 2018, Accepted: 6 November 2018, Published online: 04 December 2018. The abstract begins with: "Conventional machine learning approaches for predicting material properties from elemental compositions have emphasized the importance of leveraging domain knowledge when designing model inputs. Here, we demonstrate that by using a deep learning approach, we can bypass such manual feature engineering requiring domain knowledge and achieve much better results, even with only a few thousand training samples. We present the design and implementation of a deep neural network model referred to as *ElemNet*; it automatically captures the physical and chemical interactions and similarities between different elements using artificial intelligence which allows it to predict the materials properties with better accuracy and speed. The speed and best-in-class accuracy of *ElemNet* enable us to perform a fast and robust screening for new material candidates in a huge combinatorial space; where we predict hundreds of thousands of chemical systems that could contain yet undiscovered compounds."

Materials scientists, condensed matter physicists and solid-state chemists rely on data generated by experiments and simulation-based models to discover new materials and understand their characteristics. For the major part of the history of materials science, experimental observations have been the primary means to know the various chemical and physical properties of materials<sup>1–6</sup>. Nevertheless, experimentation of all possible combinations of material composition and crystal structures is not feasible as that would be very expensive and time-consuming, and the composition space is practically infinite. Computational methods, such as Density Functional Theory (DFT)<sup>7</sup>, offer a less expensive means to predict many material properties and processes on the atomic level<sup>8</sup>. DFT calculations have offered opportunities for large-scale data collection such as the Open Quantum Materials Database (OQMD)<sup>9,10</sup>, the Automatic Flow of Materials Discovery Library (AFLOWLIB)<sup>11</sup>, the Materials Project<sup>12</sup>, and the Novel Materials Discovery (NoMaD)<sup>13</sup>; they contain DFT computed properties of  $\sim 10^6$ – $10^8$  of experimentally-observed and hypothetical materials. In the past few decades, such materials datasets have led to the new data-driven paradigm of materials informatics<sup>14–19</sup>. The availability of such large data resources has spurred the interest of researchers in applying advanced data-driven based machine learning (ML) techniques for accelerated discovery and design of new materials with select engineering properties<sup>19–39</sup>.

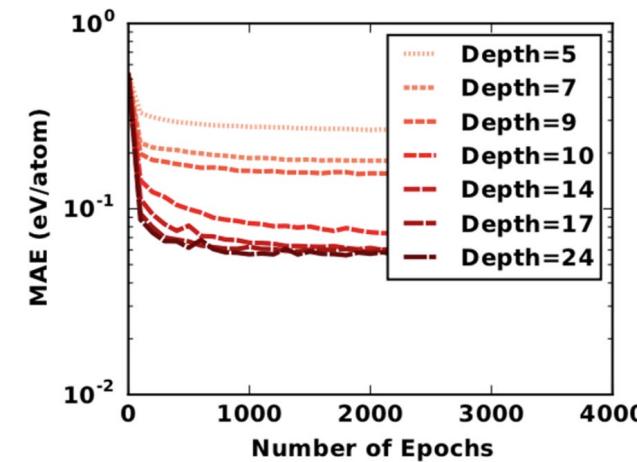
Conventionally, constructing an effective ML model requires first developing a suitable representation for the input data as shown in Fig. 1. As has been discussed in several recent works, the best representations are those that encode knowledge about the physics of the underlying problem. To that end, there have been many distinct approaches for encoding information regarding the composition<sup>21,32</sup> or crystal structure<sup>34,37,40,41</sup> of a material. For instance, Ward *et al.* developed a set of attributes based on the composition of a material that can be useful for problems including predicting formation enthalpies of crystalline materials and glass-forming ability of metal alloys<sup>32</sup>. Ghiringhelli *et al.*<sup>42</sup> analyzed the tendency for materials to form different crystal structures using thousands of descriptors. Developing ML models based on intuitive representations is evidently successful given the large number and growing rate of ML models constructed over the past several years using this approach<sup>8,18,43</sup>. However, the prediction accuracy for these problems is limited by our ability to feature engineer the materials representation to incorporate all the domain knowledge required to make correct predictions. Given that one of the major use cases of ML is for problems where the physics driving behavior is yet to be understood<sup>19</sup>, this limit

<sup>1</sup>Department of Electrical Engineering and Computer Science, Northwestern University, Evanston, USA.

<sup>2</sup>Computation Institute, University of Chicago, Chicago, USA. <sup>3</sup>Department of Materials Science and Engineering, Northwestern University, Evanston, USA. Correspondence and requests for materials should be addressed to A.A. (email: [ankitag@eeecs.northwestern.edu](mailto:ankitag@eeecs.northwestern.edu))



Layer Types	No. of units	Activation	Layer Positions
Fully-connected Layer	1024	ReLU	First to 4th
Drop-out (0.8)	1024		After 4th
Fully-connected Layer	512	ReLU	5th to 7th
Drop-out (0.9)	512		After 7th
Fully-connected Layer	256	ReLU	8th to 10th
Drop-out (0.7)	256		After 10th
Fully-connected Layer	128	ReLU	11th to 13th
Drop-out (0.8)	128		After 13th
Fully-connected Layer	64	ReLU	14th to 15th
Fully-connected Layer	32	ReLU	16th
Fully-connected Layer	1	Linear	17th



# Are CBFV features based on domain knowledge better?

Integrating Materials and Manufacturing Innovation  
https://doi.org/10.1007/s40192-020-00179-z

TECHNICAL ARTICLE

## Is Domain Knowledge Necessary for Machine Learning Materials Properties?

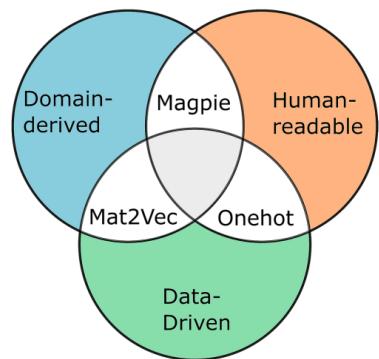
Ryan J. Murdock<sup>1</sup> · Steven K. Kauwe<sup>1</sup> · Anthony Yu-Tung Wang<sup>2</sup> · Taylor D. Sparks<sup>1</sup> 

Received: 4 June 2020 / Accepted: 9 July 2020  
© The Minerals, Metals & Materials Society 2020

### Abstract

New featurization schemes for describing materials as composition vectors in order to predict their properties using machine learning are common in the field of Materials Informatics. However, little is known about the comparative efficacy of these methods. This work sets out to make clear which featurization methods should be used across various circumstances. Our findings include, surprisingly, that simple fractional and random-noise representations of elements can be as effective as traditional and new descriptors when using large amounts of data. However, in the absence of large datasets or for data that is not fully representative, we show that the integration of domain knowledge offers advantages in predictive ability.

### Graphical abstract



**Keywords** Materials informatics · Machine learning · Featurization · Descriptors · Neural networks

Taylor D. Sparks  
sparks@eng.utah.edu

<sup>1</sup> Materials Science and Engineering Department, University of Utah, Salt Lake City, UT 84109, USA

<sup>2</sup> Technische Universität Berlin, Fachgebiet Keramische Werkstoffe/Chair of Advanced Ceramic Materials, 10623 Berlin, Germany

Published online: 27 August 2020



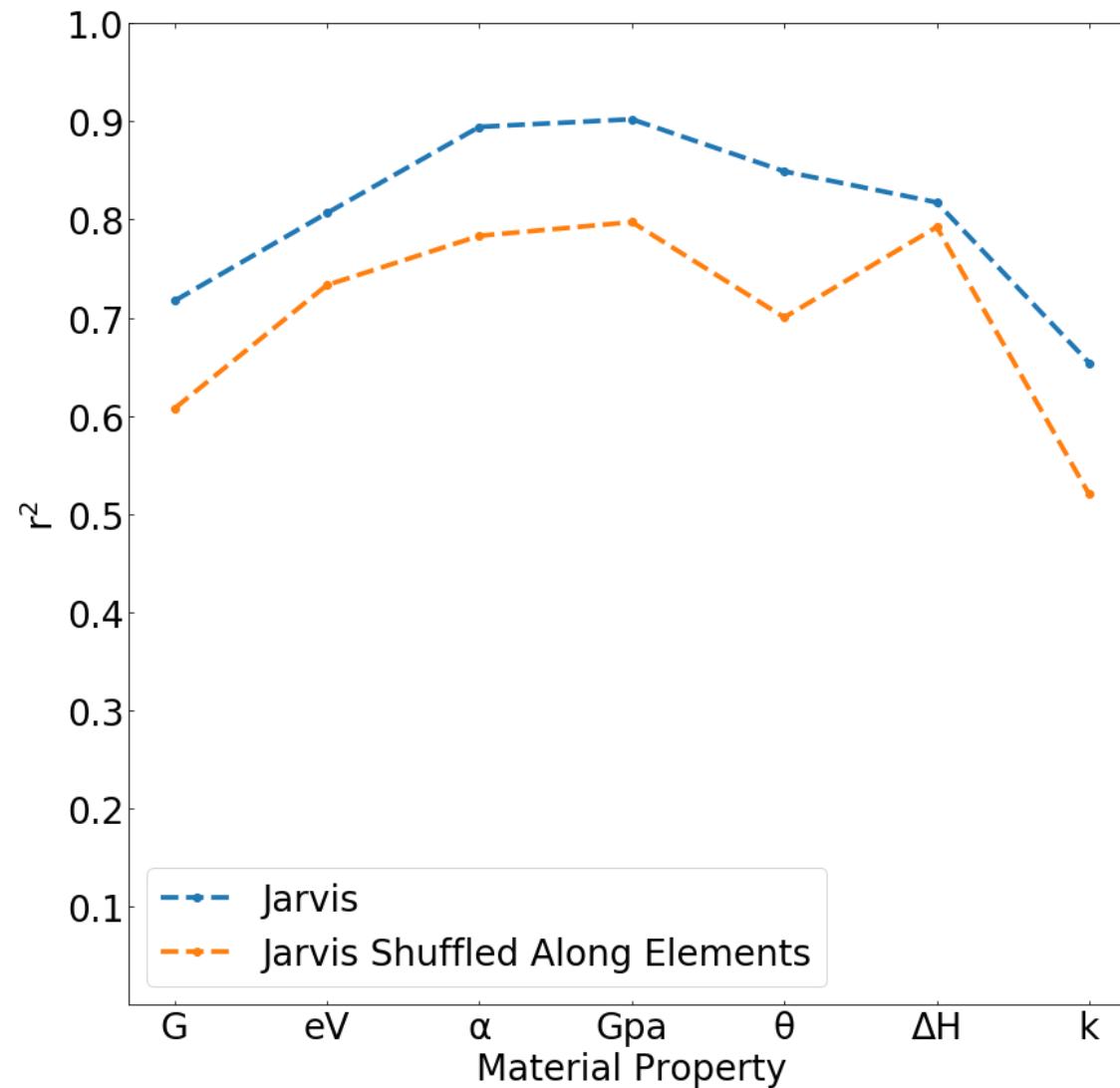
## Consider JARVIS featurization code from NIST

1	element	Number	MendeleevNumber	AtomicWeight	MeltingT	Column	Row	CovalentRadius	Electronegativity
11	Ne	10	99	20.1791	24.56	18	2	58	1.63
12	Na	11	2	22.98976928	370.87	1	3	166	0.93
13	Mg	12	68	24.305	923	2	3	141	1.31
14	Al	13	73	26.9815386	933.47	13	3	121	1.61
15	Si	14	78	28.0855	1687	14	3	111	1.9
16	P	15	83	30.973762	317.3	15	3	107	2.19

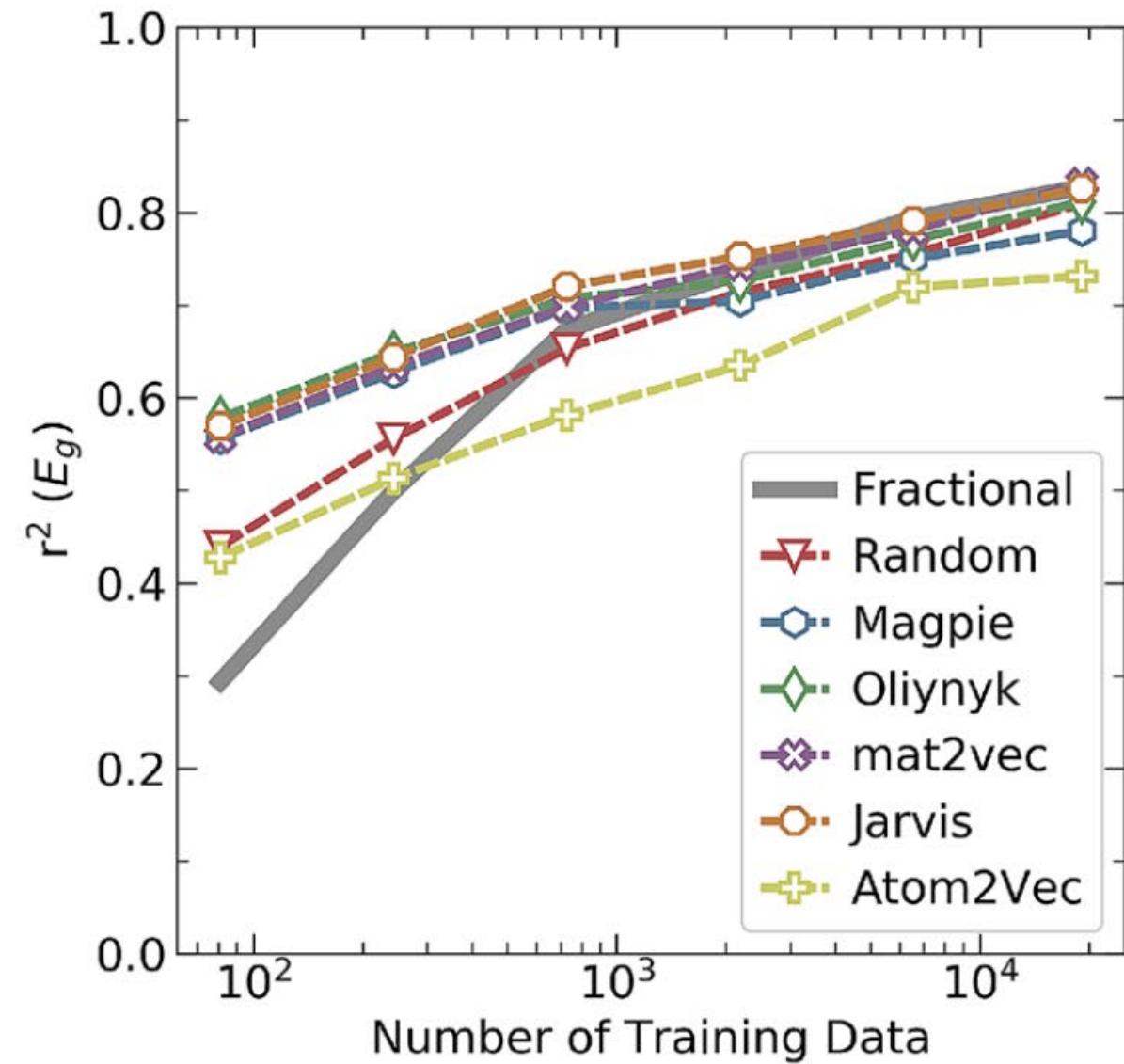
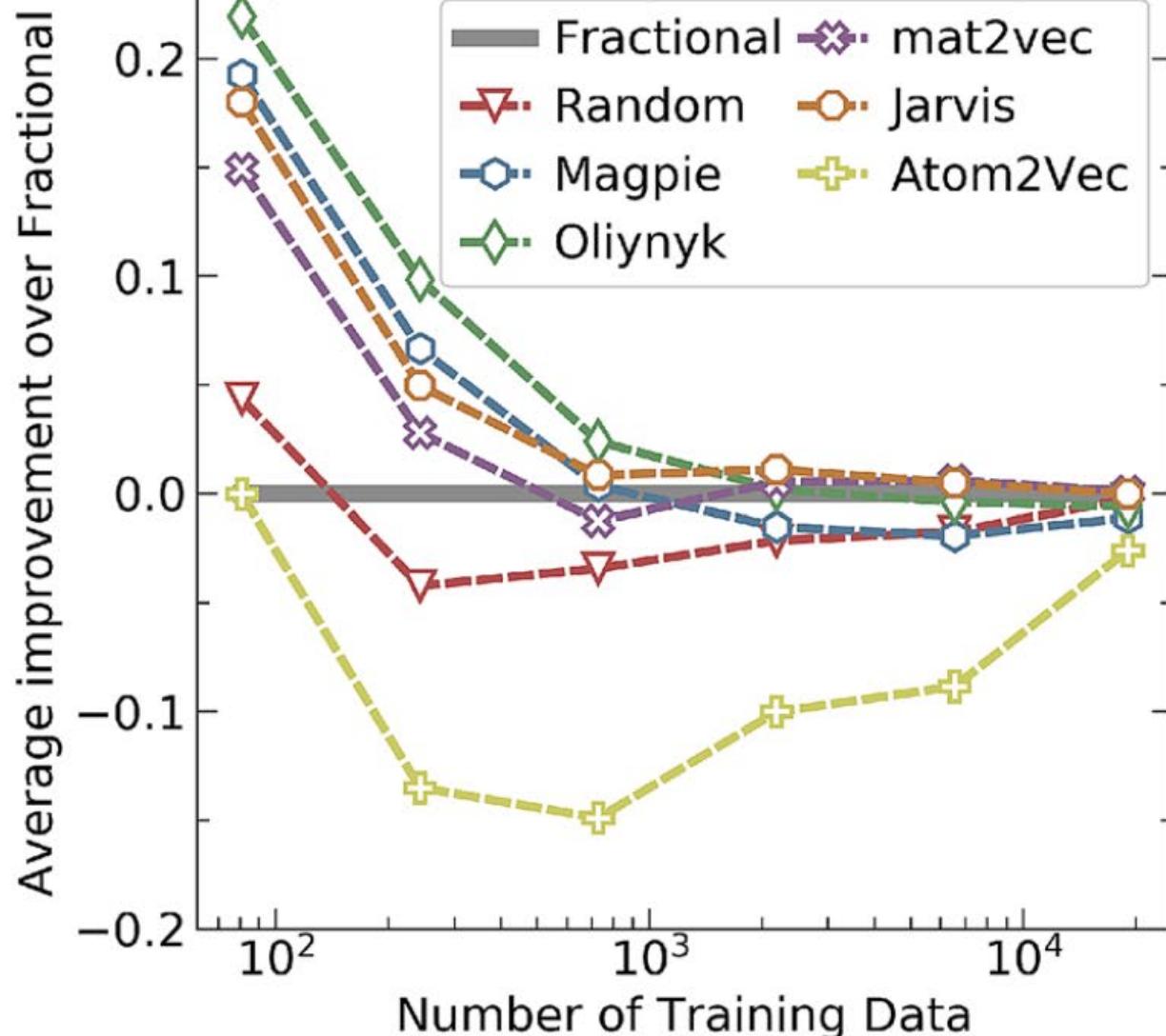
1	element	Number	MendeleevNumber	AtomicWeight	MeltingT	Column	Row	CovalentRadius	Electronegativity
11	Ne	10	99	20.1791	24.56	18	2	58	1.63
12	Na	11	2	22.98976928	370.87	1	3	166	0.93
13	Mg	12	68	24.305	923	2	3	141	1.31
14	Al	13	73	26.9815386	933.47	13	3	121	1.61
15	Si	14	78	28.0855	1687	14	3	111	1.9
16	P	15	83	30.973762	317.3	15	3	107	2.19

Will JARVIS break if we swap some rows?

# JARVIS with random shuffling does worse (sanity check!)

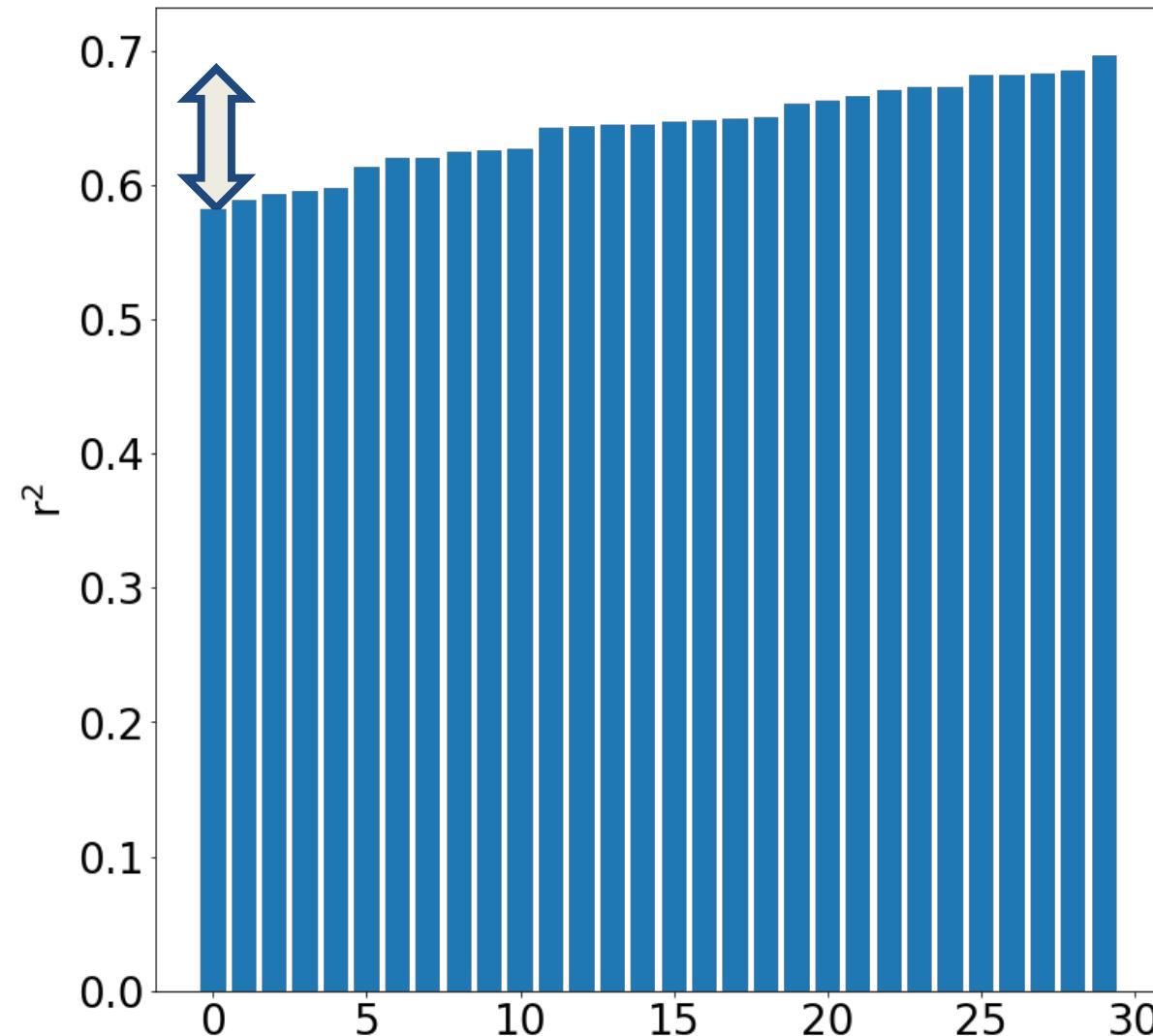


Domain knowledge helps a bit, but not much and primarily in low data limit

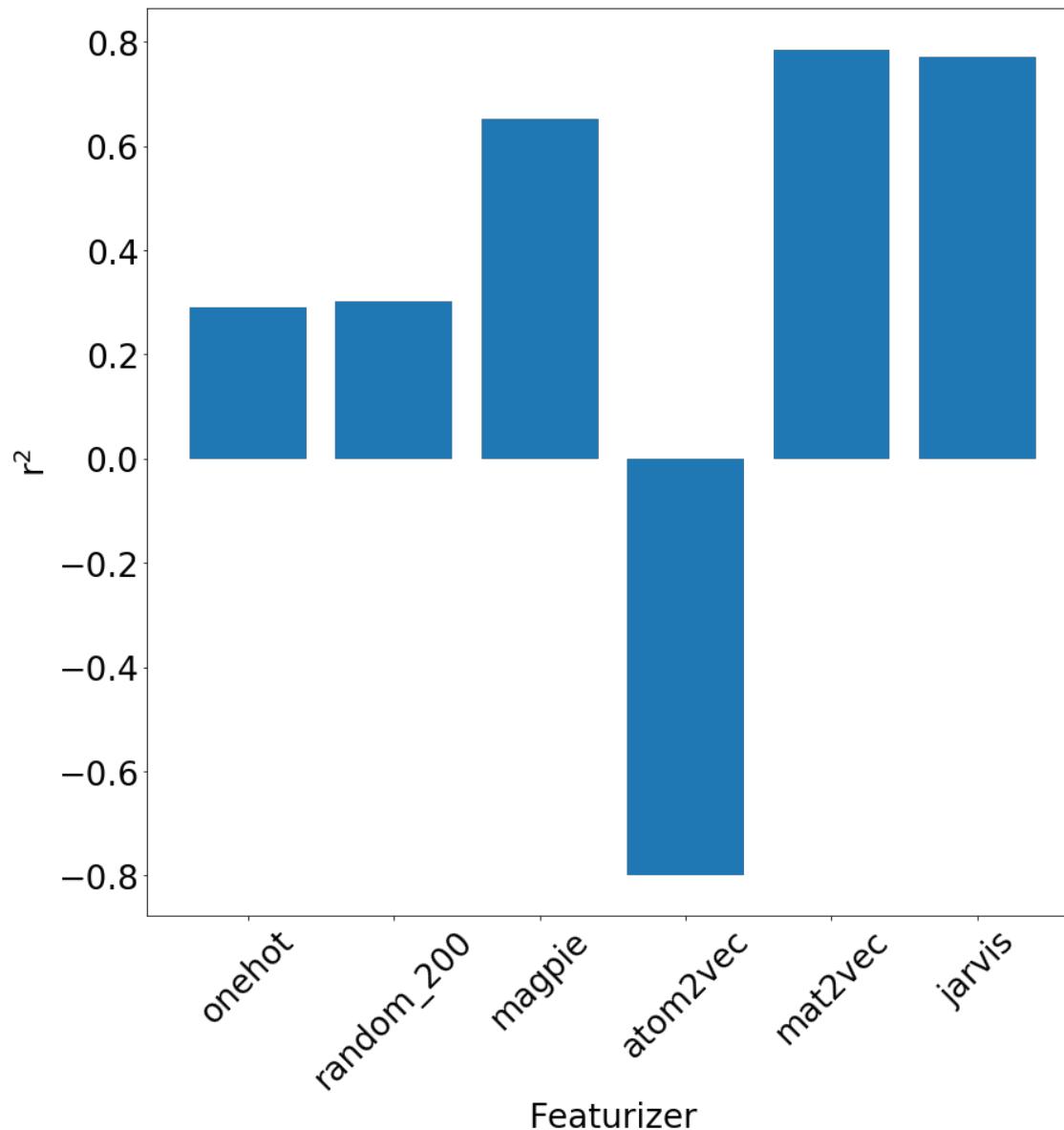


Most feature comparisons are not very meaningful!

Over 15% variation just from selecting from your randomly splitting cross-validation set!



## Feature engineering may help us extrapolate



# Our CBFV package makes it very easy to compare many different CBFV techniques in python

## CBFV 1.0.1



Latest version

pip install CBFV



Released: about 17 hours ago

Tool for quickly creating a composition-based feature vector

### Navigation

Project description

Release history

Download files

### Project links

Homepage

### Statistics

GitHub statistics:

Stars: 3

Forks: 1

Open issues/PRs: 2

### Project description

#### CBFV Package

Tool to quickly create a composition-based feature vectors from materials datafiles.

### Installation

The source code is currently hosted on GitHub at: <https://github.com/kaaiian/CBFV>

Binary installers for the latest released version are available at the [Python Package Index \(PyPI\)](#)

```
# PyPI  
pip install CBFV
```

### Making the composition-based feature vector

The CBFV package assumes your data is stored in a pandas dataframe of the following structure:

# Structure-based feature vector

