

Unsupervised word embeddings capture latent knowledge from materials science literature

Vahe Tshitoyan^{1,3*}, John Dagdelen^{1,2}, Leigh Weston¹, Alexander Dunn^{1,2}, Ziqin Rong¹, Olga Kononova², Kristin A. Persson^{1,2}, Gerbrand Ceder^{1,2*} & Anubhav Jain^{1*}

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^{1,2}, 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^{3–10}, 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^{11–13} (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 GloVe¹³ or Word2vec^{11,12}, 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 LiCoO₂ (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 LiCoO₂ are LiMn₂O₄, LiNi_{0.5}Mn_{1.5}O₄, LiNi_{0.8}Co_{0.2}O₂, LiNi_{0.8}Co_{0.15}Al_{0.05}O₂ and LiNiO₂—all of which are also lithium-ion cathode materials.

Similar to the observation made in the original Word2vec paper¹¹, 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,

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

To better visualize such embedded relationships, we projected the embeddings of Zr, Cr and 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’ (Zr – ZrO₂ ≈ Cr – Cr₂O₃ ≈ Ni – NiO) and ‘structure of’ (Zr – HCP ≈ Cr – BCC ≈ Ni – 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 ZrO₂. 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 study¹². 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 Bi₂Te₃. 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^{14,15}. However, we found that a number of materials that have relatively high cosine similarities to the word ‘thermoelectric’ never

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

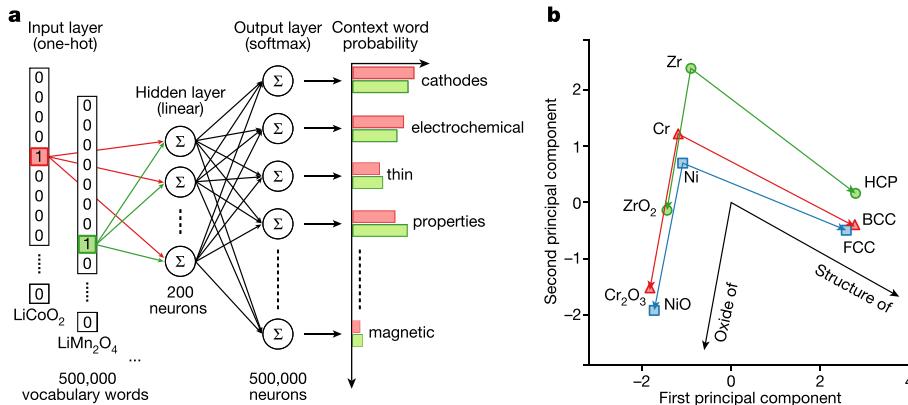


Fig. 1 | Word2vec skip-gram and analogies. **a**, Target words ‘ LiCoO_2 ’ and ‘ LiMn_2O_4 ’ are represented as vectors with ones at their corresponding vocabulary indices (for example, 5 and 8 in the schematic) and zeros everywhere else (one-hot encoding). These one-hot encoded vectors are used as inputs 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_2 and LiMn_2O_4 , 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’.

appeared explicitly in the same abstract with this word, or any other words that unequivocally identify materials as thermoelectric (Fig. 2a). Rather than dismissing these instances as spurious, we investigated whether such cases could be usefully interpreted as predictions of novel thermoelectric materials.

As a first test, we compared our predicted thermoelectric compositions with available computational data. Specifically, we identified compounds mentioned in our text corpus more than three times that are also present in a dataset¹⁶ that reports the thermoelectric power factors (an important component of the overall thermoelectric figure of merit, zT) of approximately 48,000 compounds calculated using density functional theory (DFT)^{17,18} (see Methods). A total of 9,483 compounds overlap between the two datasets, of which 7,663 were never mentioned alongside thermoelectric keywords in our text corpus and can be considered candidates for prediction. To obtain the predictions, we ranked each of these 7,663 compounds by the dot product of their normalized output embedding with the word embedding of ‘thermoelectric’ (see Supplementary Information sections S1 and S3 regarding the use of output versus word embeddings). This ranking can be interpreted as the likelihood that that material will co-occur with the word ‘thermoelectric’ in a scientific abstract despite this never occurring explicitly in the text corpus. The distributions of DFT maximum power factor values for all 9,483 materials (separated into known thermoelectrics and candidates) are plotted in Fig. 2b, and the values of the 10 highest ranked candidates from the word embedding approach are indicated with dashed lines. We find that the top ten predictions all exhibit computed power factors significantly greater than the average of candidate materials (green), and even slightly higher than the average of known thermoelectrics (purple). The average maximum power factor of $40.8 \mu\text{W K}^{-2} \text{cm}^{-1}$ for these top ten predictions is 3.6 times larger than the average of candidate materials ($11.5 \mu\text{W K}^{-2} \text{cm}^{-1}$) and 2.4 times larger than the average of known thermoelectrics ($17.0 \mu\text{W K}^{-2} \text{cm}^{-1}$). Moreover, the three highest power factors from the top ten predictions are at the 99.6th, 96.5th and 95.3rd percentiles of known thermoelectrics. We note that in contrast to supervised methods, our embeddings are based only on the text corpus and are not trained or modified in any manner using the DFT data.

Next, we compared the same model directly against experimentally measured power factors and zTs ¹⁹. Because our approach does not provide numerical estimations of these quantities, we compared the relative ranking of candidates through the Spearman rank correlation²⁰ for the 83 materials that appear both in our text corpus and the experimental

set. We obtained a 59% and 52% rank correlation of experimental results with the embedding-based ranking for maximum power factor and maximum zT , respectively. Unexpectedly, our model outperformed the DFT dataset of power factors used in the previous paragraph, which exhibits only a 31% rank correlation with the experimental maximum power factors.

Finally, we tested whether our model—if trained at various points in the past—would have correctly predicted thermoelectric materials reported later in the literature. Specifically, we generated 18 different ‘historical’ text corpora consisting only of abstracts published before cutoff years between 2001 and 2018. We trained separate word embeddings for each historical dataset, and used these embeddings to predict the top 50 thermoelectrics that were likely to be reported in future (test) years. For every year past the date of prediction, we tabulated the cumulative percentage of predicted thermoelectric compositions that were reported in the literature alongside a thermoelectric keyword. Figure 3a depicts the result from each such ‘historical’ dataset as a thin grey line. For example, the light grey line labelled ‘2015’ depicts the percentage of the top 50 predictions made using the model trained only on scientific abstracts published before 1 January 2015, and that were subsequently reported in the literature alongside a thermoelectric keyword after one, two, three or four years (that is, the years 2015–2018). Overall, our results indicate that materials from the top 50 word embedding-based predictions (red line) were on average eight times more likely to have been studied as thermoelectrics within the next five years as compared to a randomly chosen unstudied material from our corpus at that time (blue) and also three times more likely than a random material with a non-zero DFT bandgap (green). The use of larger corpora that incorporate data from more recent years improved the rate of successful predictions, as indicated by the steeper gradients for later years in Fig. 3a.

To examine these results in more detail, we focus on the fate of the top five predictions determined using only abstracts published before the year 2009. Figure 3b plots the evolution of the prediction rank of these top five compounds as more abstracts are added in subsequent years. One of these compounds, CuGaTe_2 , represents one of the best present-day thermoelectrics and would have been predicted as a top five compound four years before its publication in 2012²¹. Two of the other predictions, ReS_2 and CdIn_2Te_4 , were suggested in the literature to be good thermoelectrics^{22,23} only approximately 8–9 years after the point at which they would have first appeared in the top five list from our algorithm. We note that the sharp increase in the rank of layered ReS_2 in 2015 coincides with the discovery of a record zT for

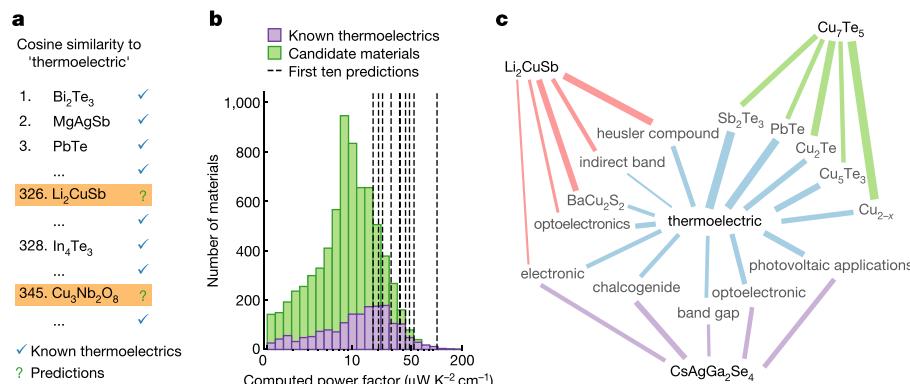


Fig. 2 | Prediction of new thermoelectric materials. **a**, A ranking of thermoelectric materials can be produced using cosine similarities of material embeddings with the embedding of the word ‘thermoelectric’. Highly ranked materials that have not yet been studied for thermoelectric applications (do not appear in the same abstracts as words ‘ZT’, ‘zT’, ‘seebeck’, ‘thermoelectric’, ‘thermoelectrics’, ‘thermoelectrical’, ‘thermoelectricity’, ‘thermoelectrically’ or ‘thermopower’) are considered to be predictions that can be tested in the future. **b**, Distributions of the power factors computed using density functional theory (see Methods) for 1,820 known thermoelectrics in the literature (purple) and 7,663 candidate materials not yet studied as thermoelectric (green). Power factors of the first ten predictions not studied as thermoelectrics in our text corpus and for which computational data are available (Li₂CuSb, CuBi₂S₂, CdIn₂Te₄, CsGeI₃, PdSe₂, KAg₂SB₄, LuRhO₃, MgB₂C₂, Li₃Sb and TlSbSe₂) are shown with black dashed lines. **c**, A graph showing how the context words of materials predicted to be thermoelectrics connect to the

SnSe²⁴—also a layered material. The final two predictions, HgZnTe and SmInO₃, contain expensive (Sm, In) or toxic (Hg) elements and have not been studied yet, and SmInO₃ has dropped appreciably in ranking with the addition of more data. The top 10 predictions for each year between 2001 and 2018 are available in Supplementary Table S3.

To illustrate how materials never mentioned next to the word ‘thermoelectric’ are identified as thermoelectrics with high expected probability, we investigated the series of connections that can lead to a prediction. In Fig. 2c, we present three materials from our top five predictions (Extended Data Table 2) alongside some of the key context words that connect these materials to ‘thermoelectric’. For instance, CsAgGa₂Se₄ has high likelihood of appearing next to ‘chalcogenide’, ‘band gap’, ‘optoelectronic’ and ‘photovoltaic applications’: many good thermoelectrics are chalcogenides, the existence of a bandgap is

word thermoelectric. The width of the edges between ‘thermoelectric’ and the context words (blue) is proportional to the cosine similarity between the word embeddings of the nodes, whereas the width of the edges between the materials and the context words (red, green and purple) is proportional to the cosine similarity between the word embeddings of context words and the output embedding of the material. The materials are the first (Li₂CuSb), third (CsAgGa₂Se₄) and fourth (Cu₇Te₅) predictions. The context words are top context words according to the sum of the edge weights between the material and the word ‘thermoelectric’. Wider paths are expected to make larger contributions to the predictions. Examination of the context words demonstrates that the algorithm is making predictions on the basis of crystal structure associations, co-mentions with other materials for the same application, associations between different applications, and key phrases that describe the material’s known properties.

crucial for the majority of thermoelectrics, and there is a large overlap between optoelectronic, photovoltaic and thermoelectric materials (see Supplementary Information section S8). Consequently, the correlations between these keywords and CsAgGa₂Se₄ led to the prediction. This direct interpretability is a major advantage over many other machine learning methods for materials discovery. We also note that several predictions were found to exhibit promising properties despite not being in any well known thermoelectric material classes (see Supplementary Information section S10). This demonstrates that word embeddings go beyond trivial compositional or structural similarity and have the potential to unlock latent knowledge not directly accessible to human scientists.

As a final step, we verified the generalizability of our approach by performing historical validation of predictions for three additional keywords—‘photovoltaics’, ‘topological insulator’ and ‘ferroelectric’. We

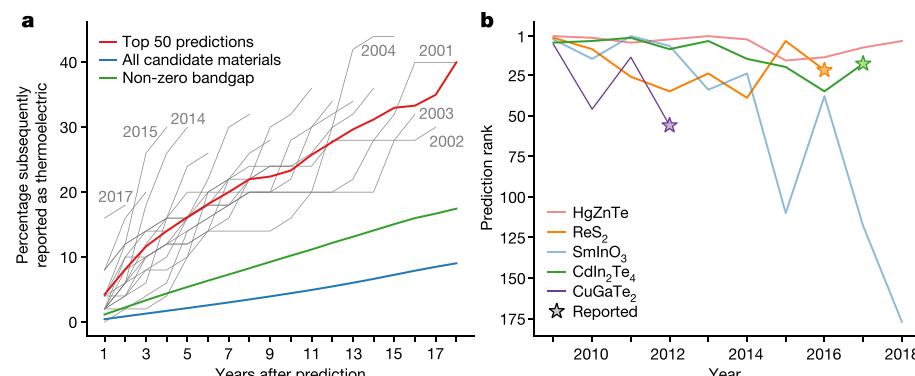


Fig. 3 | Validation of the predictions. **a**, Results of prediction of thermoelectric materials using word embeddings obtained from various historical datasets. Each grey line uses only abstracts published before that year to make predictions (for example, predictions for 2001 are performed using abstracts from 2000 and earlier). The lines plot the cumulative percentage of predicted materials subsequently reported as thermoelectrics in the years following their predictions; earlier predictions

can be analysed over longer test periods, resulting in longer grey lines. The results are averaged (red) and compared to baseline percentages from either all materials (blue) or non-zero DFT bandgap²⁷ materials (green). **b**, The top five predictions from the year 2009 dataset, and evolution of their prediction ranks as more data are collected. The marker indicates the year of first published report of one of the initial top five predictions as a thermoelectric.

emphasize that the word embeddings used for these predictions are the same as those for thermoelectrics predictions; we have simply modified the dot product to be with a different target word. Notably, with almost no change in procedure, we find trends similar to the ones in Fig. 3a for all three functional applications, with the results summarized in Extended Data Fig. 2 and Extended Data Table 3.

The success of our unsupervised approach can partly be attributed to the choice of the training corpus. The main purpose of abstracts is to communicate information in a concise and straightforward manner, avoiding unnecessary words that may increase noise in embeddings during training. The importance of corpus selection is demonstrated in Extended Data Table 4, where we show that discarding abstracts unrelated to inorganic materials science improves performance, and models trained on the set of all Wikipedia articles (about ten times more text than our corpus) perform substantially worse on materials science analogies. Contrary to what might seem like the conventional machine learning mantra, throwing more data at the problem is not always the solution. Instead, the quality and domain-specificity of the corpus determine the utility of the embeddings for domain-specific tasks.

We suggest that the methodology described here can also be generalized to other language models, such that the probability of an entity (such as a material or molecule) co-occurring with words that represent a target application or property can be treated as an indicator of performance. Such language-based inference methods can become an entirely new field of research at the intersection between natural language processing and science, going beyond simply extracting entities and numerical values from text and leveraging the collective associations present in the research literature. Substitution of Word2vec with context-aware embeddings such as BERT²⁵ or ELMo²⁶ could lead to improvements for functional material predictions, as these models are able to change the embedding of the word based on its context. They substantially outperform context-independent embeddings such as Word2vec or GloVe across all conventional NLP tasks. Also, in addition to co-occurrences, these models can capture more complex relationships between words in the sentence, such as negation. In the current study, the effects of negation are somewhat mitigated because scientific abstracts often emphasize positive relationships. However, a natural extension of this work is to parse the full texts of articles. We expect the full texts will contain more negative relationships and in general more variable and complex sentences, and will therefore require more powerful methods.

Scientific progress relies on the efficient assimilation of existing knowledge in order to choose the most promising way forward and to minimize re-invention. As the amount of scientific literature grows, this is becoming increasingly difficult, if not impossible, for an individual scientist. We hope that this work will pave the way towards making the vast amount of information found in scientific literature accessible to individuals in ways that enable a new paradigm of machine-assisted scientific breakthroughs.

Online content

Any methods, additional references, Nature Research reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41586-019-1335-8>.

Received: 19 December 2018; Accepted: 8 May 2019;
Published online 3 July 2019.

- Hill, J. et al. Materials science with large-scale data and informatics: unlocking new opportunities. *MRS Bull.* **41**, 399–409 (2016).
- Butler, K. T., Davies, D. W., Cartwright, H., Isayev, O. & Walsh, A. Machine learning for molecular and materials science. *Nature* **559**, 547–555 (2018).
- Friedman, C., Kra, P., Yu, H., Krauthammer, M. & Rzhetsky, A. GENIES: a natural-language processing system for the extraction of molecular pathways from journal articles. *Bioinformatics* **17**, S74–S82 (2001).
- Müller, H. M., Kenny, E. E. & Sternberg, P. W. Textpresso: an ontology-based information retrieval and extraction system for biological literature. *PLoS Biol.* **2**, e309 (2004).
- Swain, M. C. & Cole, J. M. ChemDataExtractor: a toolkit for automated extraction of chemical information from the scientific literature. *J. Chem. Inf. Model.* **56**, 1894–1904 (2016).
- Eltyeb, S. & Salim, N. Chemical named entities recognition: a review on approaches and applications. *J. Cheminform.* **6**, 17 (2014).
- Kim, E. et al. Materials synthesis insights from scientific literature via text extraction and machine learning. *Chem. Mater.* **29**, 9436–9444 (2017).
- Leaman, R., Wei, C. H. & Lu, Z. TmChem: a high performance approach for chemical named entity recognition and normalization. *J. Cheminform.* **7**, S3 (2015).
- Krallinger, M., Rabal, O., Lourenço, A., Oyarzabal, J. & Valencia, A. Information retrieval and text mining technologies for chemistry. *Chem. Rev.* **117**, 7673–7761 (2017).
- Spangler, S. et al. Automated hypothesis generation based on mining scientific literature. In Proc. 20th ACM SIGKDD Intl Conf. Knowledge Discovery and Data Mining 1877–1886 (ACM, 2014).
- Mikolov, T., Corrado, G., Chen, K. & Dean, J. Efficient estimation of word representations in vector space. Preprint at <https://arxiv.org/abs/1301.3781> (2013).
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. & Dean, J. Distributed representations of words and phrases and their compositionality. Preprint at <https://arxiv.org/abs/1310.4546> (2013).
- Pennington, J., Socher, R. & Manning, C. GloVe: global vectors for word representation. Proc. 2014 Conf. Empirical Methods in Natural Language Processing (EMNLP) 1532–1543 (Association for Computational Linguistics, 2014).
- Liu, W. et al. New trends, strategies and opportunities in thermoelectric materials: a perspective. *Materials Today Physics* **1**, 50–60 (2017).
- He, J. & Tritt, T. M. Advances in thermoelectric materials research: looking back and moving forward. *Science* **357**, eaak9997 (2017).
- Ricci, F. et al. An ab initio electronic transport database for inorganic materials. *Sci. Data* **4**, 170085 (2017).
- Hohenberg, P. & Kohn, W. Inhomogeneous electron gas. *Phys. Rev.* **136**, B864–B871 (1964).
- Kohn, W. & Sham, L. J. Self-consistent equations including exchange and correlation effects. *Phys. Rev.* **140**, A1133–A1138 (1965).
- Gaultois, M. W. et al. Data-driven review of thermoelectric materials: performance and resource considerations. *Chem. Mater.* **25**, 2911–2920 (2013).
- Spearman, C. The proof and measurement of association between two things. *Am. J. Psychol.* **15**, 72–101 (1904).
- Plirdpring, T. et al. Chalcopyrite CuGaTe₂: a high-efficiency bulk thermoelectric material. *Adv. Mater.* **24**, 3622–3626 (2012).
- Tian, H. et al. Low-symmetry two-dimensional materials for electronic and photonic applications. *Nano Today* **11**, 763–777 (2016).
- Pandey, C., Sharma, R. & Sharma, Y. Thermoelectric properties of defect chalcopyrites. *AIP Conf. Proc.* **1832**, 110009 (2017).
- Zhao, L.-D. et al. Ultralow thermal conductivity and high thermoelectric figure of merit in SnSe crystals. *Nature* **508**, 373–377 (2014).
- Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. BERT: pre-training of deep bidirectional transformers for language understanding. Preprint at <https://arxiv.org/abs/1810.04805> (2018).
- Peters, M. E. et al. Deep contextualized word representations. Preprint at <https://arxiv.org/abs/1802.05365> (2018).
- Jain, A. et al. The materials project: a materials genome approach to accelerating materials innovation. *APL Mater.* **1**, 011002 (2013).

Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

© The Author(s), under exclusive licence to Springer Nature Limited 2019

METHODS

Data collection and processing. We obtained approximately 3.3 million abstracts, primarily focused on materials science, physics, and chemistry, through a combination of Elsevier's Scopus and Science Direct application programming interfaces (APIs) (<https://dev.elsevier.com/>), the Springer Nature API (<https://dev.springernature.com/>), and web scraping. Parts of abstracts (or full abstracts) that were in foreign languages were removed using text search and regular expression matching, as were articles with metadata types corresponding to 'Announcement', 'BookReview', 'Erratum', 'EditorialNotes', 'News', 'Events' and 'Acknowledgement'. Abstracts with titles containing keywords 'Foreword', 'Prelude', 'Commentary', 'Workshop', 'Conference', 'Symposium', 'Comment', 'Retract', 'Correction', 'Erratum' and 'Memorial' were also selectively removed from the corpus. Some abstracts contained leading or trailing copyright information, which was removed using regular expression matching and heuristic rules. Leading words and phrases such as 'Abstract:' were also removed using similar methods. We further retained only abstracts related to inorganic materials according to a binary classifier (see 'Abstract classification' below). We tuned the classifier for high recall to guarantee the presence of the majority of relevant abstracts at the expense of retaining some irrelevant ones. Removing irrelevant abstracts substantially improved the performance of our algorithm, as discussed in more detail in Supplementary Information section S2. The 1.5 million abstracts that were classified as relevant were tokenized using ChemDataExtractor⁵ to produce the individual words. The tokens that were identified as valid chemical formulae using pymatgen²⁸ combined with regular expression and rule-based techniques were normalized such that the order of elements and common multipliers did not matter (NiFe is the same as Fe₅₀Ni₅₀). Valence states of elements were split into separate tokens (for example, Fe(III) becomes two separate tokens, Fe and (III)). We also performed selective lower-casing and deaccenting. If the token was not a chemical formula or an element symbol, and if only the first letter was uppercase, we lower-cased the word. Thus, chemical formulae and abbreviations stayed in their common form, whereas words at the beginning of sentences and proper nouns were lower-cased. Numbers with units were often not tokenized correctly by ChemDataExtractor. We addressed this in the processing step by splitting the common units from numbers and converting all numbers to a special token <nUm>. This reduced the vocabulary size by approximately 20,000 words. We found that correct preprocessing, especially the choice of phrases to include as individual tokens, substantially improved the results. The code used for preprocessing is available at <https://github.com/materialsintelligence/mat2vec>.

Abstract classification. This work focuses on inorganic materials science. However, our corpus contained some abstracts that fell outside this scope (for example, articles on polymer science). We removed articles outside our targeted area of research literature by training a binary classifier that could label abstracts as 'relevant' or 'not relevant'. We annotated 1,094 randomly selected abstracts; of these, 588 were labelled as 'relevant' and 494 were labelled 'not relevant'. The labelled abstracts were used as data to train a classifier; we used a linear classifier based on logistic regression, where each document is described by a term frequency-inverse document frequency (tf-idf) vector. The classifier achieved an accuracy (f1-score) of 89% using fivefold cross-validation.

Word2vec training. We used the Word2vec implementation in gensim (<https://radimrehurek.com/gensim/>) with a few modifications. We found that skip-gram with negative sampling loss ($n = 15$) performed best (see Supplementary Information section S2 for comparison between models). The vocabulary consisted of all words that occurred more than five times as well as normalized chemical formulae, independent of the number of mentions. The phrases were generated using a minimum phrase count of 10, score threshold of 15 (ref. ¹²) and phrase depth of 2. The latter meant that we repeated the process twice, allowing generation of up to four grams. We also included common terms such as '·', 'of', 'to', 'a' and 'the', which in exceptional cases led to phrases with more tokens. For example, 'state-of-the-art thermoelectric' is one of the five 8-token phrases in our vocabulary. At the end of each phrase generation cycle, we removed phrases that contained punctuation and numbers. The size of the vocabulary approximately doubled after phrase generation. The rest of the hyperparameters were as follows: we used 200-dimensional embeddings, a learning rate of 0.01 decreasing to 0.0001 in 30 epochs, a context window of 8 and subsampling with a 10^{-4} threshold, which subsamples approximately the 400 most common words. Hyperparameters were optimized for performance on approximately 15,000 grammatical and 15,000 materials science analogies, with the score defined as the percentage of correctly 'solved' analogies from the two sets. Hyperparameter optimization and the choice of the corpus are also discussed in more detail in Supplementary Information section S2. The code used for the training and the full list of analogies used in this study are available at <https://github.com/materialsintelligence/mat2vec>.

Thermoelectric power factors. Each materials structure optimization and band structure calculation was performed with density functional theory (DFT)

using the projector augmented wave (PAW)²⁹ pseudopotentials and the Perdew-Burke-Ernzerhof (PBE)³⁰ generalized-gradient approximation (GGA), implemented in the Vienna Ab initio Simulation Package (VASP)^{31,32}. A +U correction was applied to transition metal oxides¹⁶. Seebeck coefficient (S) and electrical conductivity (σ) were calculated using the BoltzTrap package³³ using a constant relaxation time of 10^{-14} s at simulated temperatures between 300 K and 1,300 K and for carrier concentrations (doping) between 10^{16} cm^{-3} and 10^{22} cm^{-3} . A 48,770-material subset of the calculations was taken from a previous work¹⁶, the remaining calculations were performed in this work using the software atomate³⁴. All calculations used the pymatgen²⁸ Python library within the FireWorks³⁵ workflow management framework. To more realistically evaluate the thermoelectric potential of a candidate material, we devised a simple strategy to condense the complex behaviour of the S and σ tensors into a single power factor metric. For each semiconductor type $\eta \in \{n, p\}$, temperature T , and doping level c , the S and σ tensors were averaged over the three crystallographic directions, and the average power factor, PF_{avg} , was computed. PF_{avg} is a crude estimation of the polycrystalline power factor from the power factor of a perfect single crystal. To account for the complex behaviour of S and σ with T , c , and η , we then took the maximum average power factor over T , c , and η constrained to a maximum cutoff temperature T_{cut} and maximum cutoff doping c_{cut} . Formally, this is $\text{PF}_{\text{avg}, \text{max}}^{T_{\text{cut}}, c_{\text{cut}}} = \max \text{PF}(\eta, T, c)$ such that $T \leq T_{\text{cut}}, c \leq c_{\text{cut}}$. We chose $T_{\text{cut}} = 600$ K and $c_{\text{cut}} = 10^{20} \text{ cm}^{-3}$ because these values resulted in better correspondence with the experimental dataset than more optimistic values, owing to the limitations of the constant relaxation time approximation. The resulting power factor, $\text{PF}_{\text{avg}, \text{max}}^{600 \text{ K}, 10^{20}}$, is equated with 'computed power factor' in this study. To rank materials according to experimental power factors (or zT), we used the maximum value for a given stoichiometry across all experimental conditions present in the dataset from Gaulois et al.¹⁹.

Data availability

The scientific abstracts used in this study are available via Elsevier's Scopus and Science Direct API's (<https://dev.elsevier.com/>) and the Springer Nature API (<https://dev.springernature.com/>). The list of DOIs used in this study, the pre-trained word embeddings and the analogies used for validation of the embeddings are available at <https://github.com/materialsintelligence/mat2vec>. All other data generated and analysed during the current study are available from the corresponding authors on reasonable request.

Code availability

The code used for text preprocessing and Word2vec training are available at <https://github.com/materialsintelligence/mat2vec>.

28. Ong, S. P. et al. Python Materials Genomics (pymatgen): a robust, open-source Python library for materials analysis. *Comput. Mater. Sci.* **68**, 314–319 (2013).
29. Kresse, G. & Joubert, D. From ultrasoft pseudopotentials to the projector augmented-wave method. *Phys. Rev. B Condens. Matter Mater. Phys.* **59**, 1758–1775 (1999).
30. Perdew, J. P., Burke, K. & Ernzerhof, M. Generalized gradient approximation made simple. *Phys. Rev. Lett.* **77**, 3865–3868 (1996).
31. Kresse, G. & Furthmüller, J. Efficient iterative schemes for ab initio total-energy calculations using a plane-wave basis set. *Phys. Rev. B Condens. Matter* **54**, 11169–11186 (1996).
32. Kresse, G. & Furthmüller, J. Efficiency of ab-initio total energy calculations for metals and semiconductors using a plane-wave basis set. *Comput. Mater. Sci.* **6**, 15–50 (1996).
33. Madsen, G. K. & Singh, D. J. Boltztrap. A code for calculating band-structure dependent quantities. *Comput. Phys. Commun.* **175**, 67–71 (2006).
34. Mathew, K. et al. Atomate: a high-level interface to generate, execute, and analyze computational materials science workflows. *Comput. Mater. Sci.* **139**, 140–152 (2017).
35. Jain, A. et al. Fireworks: a dynamic workflow system designed for high-throughput applications. *Concurr. Comput.* **27**, 5037–5059 (2013).
36. Yang, X., Dai, Z., Zhao, Y., Liu, J. & Meng, S. Low lattice thermal conductivity and excellent thermoelectric behavior in Li₃Sb and Li₃Bi. *J. Phys. Condens. Matter* **30**, 425401 (2018).
37. Wang, Y., Gao, Z. & Zhou, J. Ultralow lattice thermal conductivity and electronic properties of monolayer 1T phase semimetal SiTe₂ and SnTe₂. *Physica E* **108**, 53–59 (2019).
38. Mukherjee, M., Yumnam, G. & Singh, A. K. High thermoelectric figure of merit via tunable valley convergence coupled low thermal conductivity in Al^IB^{IV}C^V chalcopyrites. *J. Phys. Chem. C* **122**, 29150–29157 (2018).
39. Kim, E. et al. Machine-learned and codified synthesis parameters of oxide materials. *Sci. Data* **4**, 170127 (2017).
40. Faber, F. A., Lindmaa, A., Von Lilienfeld, O. A. & Armiento, R. Machine learning energies of 2 million elpasolite (ABC₂D₆) crystals. *Phys. Rev. Lett.* **117**, 135502 (2016).
41. Zhou, Q. et al. Learning atoms for materials discovery. *Proc. Natl. Acad. Sci. USA* **115**, E6411–E6417 (2018).

Acknowledgements This work was supported by Toyota Research Institute through the Accelerated Materials Design and Discovery program. We thank T. Botari, M. Horton, D. Mrdjenovich, N. Mingione and A. Faghaninia for discussions.

Author contributions All authors contributed to the conception and design of the study, as well as writing of the manuscript. V.T. developed the data processing pipeline, trained and optimized the Word2vec embeddings, trained the machine learning models for property predictions and generated the thermoelectric predictions. V.T., J.D. and L.W. analysed the results and developed the software infrastructure for the project. J.D. trained and optimized the GloVe embeddings and developed the data acquisition infrastructure. L.W. performed the abstract classification. A.D. performed the DFT calculation of thermoelectric power factors. Z.R. contributed to data acquisition. O.K. developed the code

for normalization of material formulae. A.D., Z.R. and O.K. contributed to the analysis of the results. K.A.P., G.C. and A.J. supervised the work.

Competing interests The authors declare no competing interests.

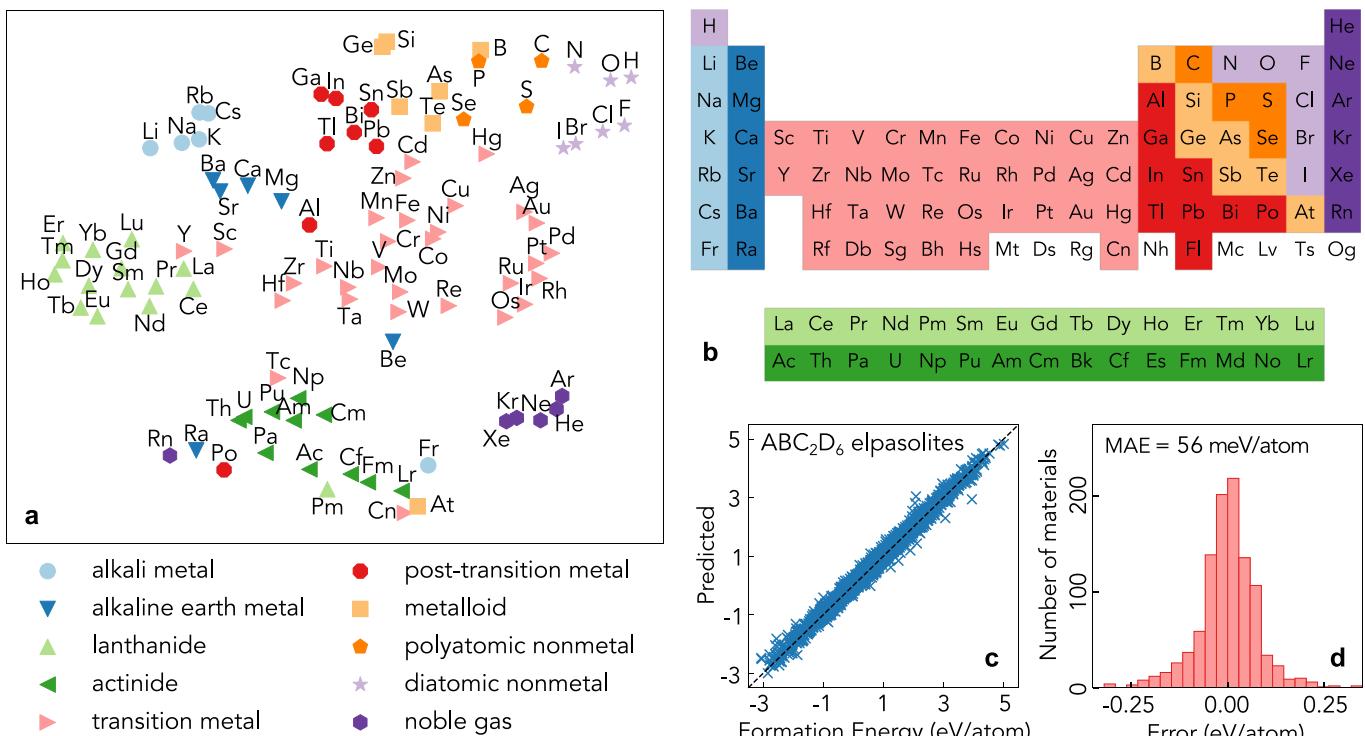
Additional information

Extended data is available for this paper at <https://doi.org/10.1038/s41586-019-1335-8>.

Supplementary information is available for this paper at <https://doi.org/10.1038/s41586-019-1335-8>.

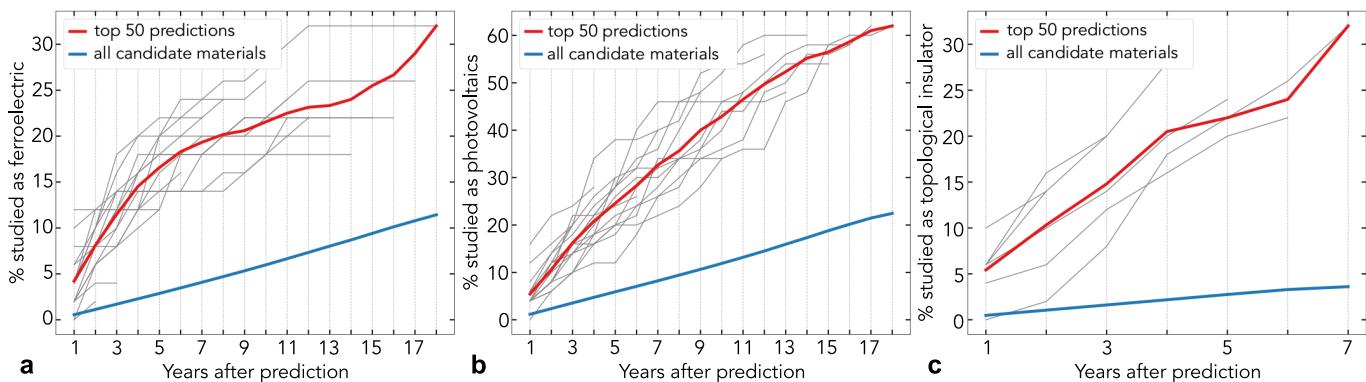
Correspondence and requests for materials should be addressed to V.T., G.C. and A.J.

Reprints and permissions information is available at <http://www.nature.com/reprints>.



Extended Data Fig. 1 | Chemistry is captured by word embeddings.
a, Two-dimensional *t*-distributed stochastic neighbour embedding (*t*-SNE) projection of the word embeddings of 100 chemical element names (for example, ‘hydrogen’) labelled with the corresponding element symbols and grouped according to their classification. Chemically similar elements are seen to cluster together and the overall distribution exhibits a topology reminiscent of the periodic table itself (compare to **b**). Arranged from top left to bottom right are the alkali metals, alkaline earth metals, transition metals, and noble gases while the trend from top right to bottom left generally follows increasing atomic number (see Supplementary Information section S4 for a more detailed discussion). **b**, The periodic table coloured according to the classification shown in **a**. **c**, Predicted

versus actual (DFT) values of formation energies of approximately 10,000 ABC₂D₆ elpasolite compounds⁴⁰ using a simple neural network model with word embeddings of elements as features (see Supplementary Information section S6 for the details of the model). The data points in the plot use fivefold cross-validation. **d**, Error distribution for the 10% test set of elpasolite formation energies. With no extensive optimization, the word embeddings achieve a mean absolute error (MAE) of 0.056 eV per atom, which is substantially smaller than the 0.1 eV per atom error reported for the same task in the original study using hand-crafted features⁴⁰ and the 0.15 eV per atom achieved in a recent study using element features automatically learned from crystal structures of more than 60,000 compounds⁴¹.



Target word / phrase	Indicator words of a potentially existing study
ferroelectric	ferroelectric, antiferroelectric, ferroelectrics, ferroelectricity, ferro-electricity, relaxor, paraelectric, para-electric, multiferroics, multiferroic, anti-ferroelectric, paraelectricity, ferroelectric, para-electric, ferro-electric, piezoelectric, PZT, pyroelectric, piezo-electric, pyro-electric, magnetoelectric, magnetoelectricity
photovoltaics	solar, photovoltaic, PV, photodevices, photoelectronics, optoelectronic, optoelectronics, nano-optoelectronics, nano-optoelectronic, opto-electronic, opto-electronics, photodiodes, photodiode, photodetectors, photodetector, photosensor, photosensors, photosensing, LED, LEDs
d topological insulator	topological, topologically

Extended Data Fig. 2 | Historical validations of functional material predictions. **a-c,** Ferroelectric (**a**), photovoltaic (**b**) and topological insulator predictions (**c**) using word embeddings obtained from various historical datasets, similar to Fig. 3a. For ferroelectrics and photovoltaics, the range of prediction years is 2001–2018. The phrase ‘topological insulator’ obtained its own embedding in our corpus only in 2011 (owing to count and vocabulary size limits), so it is possible to analyse the results only over a shorter time period (2011–2018). Each grey line uses only

abstracts published before a certain year to make predictions. The lines show the cumulative percentage of predicted materials studied in the years following their predictions; earlier predictions can be analysed over longer test periods. The results are averaged in red and compared to baseline percentages from all materials. **d**, The target word or phrase used to rank materials for each application (based on cosine similarity), and the corresponding words used as indicators for a potentially existing study.

Extended Data Table 1 | Materials science analogies

Relationship	Example vector operation	Answer	Validation pairs	Accuracy (%)
Chemical element names	helium - He + Fe	= iron	8372	71.4
Crystal symmetries	cubic - GaAs + CdSe	= hexagonal	2034	35.4
Crystal structure names	zincblende - GaP + GaN	= wurtzite	556	18.7
Elemental crystal structures	dhcp - La + Cr	= bcc	1198	48.6
Principal oxides	Al ₂ O ₃ - Al + Si	= SiO ₂	650	48.8
Units	pressure - Pa + Hz	= frequency	452	35.4
Magnetic properties	ferromagnetic - NiCo + IrMn	= antiferromagnetic	622	41.0
Applications	thermoelectric - PbTe + LiFePO ₄	= cathode materials	-	-
Grammar	structures - structure + energy	= energies	15162	61.6
Total			29046	60.1

Examples of verified word analogies corresponding to various materials science concepts. The first column lists the types of tested analogies. The second column is an example vector operation for the corresponding analogy type, with the observed answer listed in the third column. The fourth column gives the number of pairs used for scoring the corresponding analogy task, with the resulting score of our model shown in the fifth column. Application analogies were not tested quantitatively and the example is for demonstration purposes only. The full list of tested analogies is available at <https://github.com/materialsintelligence/mat2vec>.

Extended Data Table 2 | Top 50 thermoelectric predictions

Top 50 thermoelectric predictions				
1. Li ₂ CuSb	11. MgB ₂ C ₂	21. CuIn ₅ S ₈	31. Ag ₃ SbS ₃	41. Eu ₂ CuSi ₃
2. Cu ₃ Nb ₂ O ₈	12. AlGaSb	22. AlFe ₂ B ₂	32. (CH ₃ NH ₃) ₃ Bi ₂ I ₉	42. Cu ₂ ZnSiS ₄
3. CsAgGa ₂ Se ₄	13. Li ₃ Sb*	23. CeTe	33. Ba ₄ Ga ₄ SnSe ₁₂	43. Bi ₄ Br ₄
4. Cu ₇ Te ₅	14. Ba ₂₄ Si ₁₀₀	24. Pb _{0.902} Sn _{0.098} Se	34. In ₃ Se ₂	44. (YbS) _{1.25} CrS ₂
5. Ge ₁₅ Sb ₄₇ Te ₃₈	15. GaNAsP	25. Bi _{0.95} La _{0.05} FeO ₃	35. Ag ₂ PbGeS ₄	45. KCu ₂ SbS ₃
6. CsGeI ₃	16. ZnGa ₂ Te ₄	26. CdSnP ₂ *	36. AgCrO ₂	46. Cu ₂ GeTe ₃
7. KAg ₂ SbS ₄	17. Cu ₃ TaS ₄	27. PdTe	37. TlCu ₂ Se ₂	47. NaLaS ₂
8. SnTe ₂ *	18. HgMnTe	28. HgZnTe	38. AgGa	48. Hg _{0.78} Cd _{0.22} Te
9. Ni ₂ Te ₃	19. MnBi ₂ Se ₄	29. Pr _{0.7} Ca _{0.3} Mn _{0.95} Co _{0.05} O ₃	39. BSb	49. InSn
10. Yb ₁₁ AlSb ₉	20. Ag ₆ Si ₂ O ₇	30. Cd ₄ GeSe ₆	40. CdIn ₂ S ₂ Se ₂	50. ReSSe

The top 50 thermoelectric predictions using the full text corpus available at the time of writing. Some of these have practical limitations (for example, the presence of air-sensitive species or toxic and expensive elements), but others appear to be experimentally testable candidates. An exhaustive manual literature search revealed that, from the first 150 predictions using the full corpus of collected abstracts published through 2018, 48 materials (32%) had already been studied as thermoelectrics in papers that were not represented in our corpus, many of which were published within the last two years. In the top 50 listed here we have excluded any predictions for which we could find thermoelectric reports outside our corpus.

*Materials reported as good thermoelectrics while this manuscript was being prepared and reviewed^{36–38}.

Extended Data Table 3 | Top five functional material predictions and context words

Prediction	Top 10 most contributing context words
Topological Insulator	
1. Sc ₂ CF ₂	armchair direction, zigzag direction, phosphorene, Sc ₂ C(OH) ₂ , MXene, semiconducting, armchair, semiconductor, Sc ₂ AIC, strongly anisotropic
2. LaCuOSe	layered oxychalcogenides, oxychalcogenides, oxychalcogenide, degenerate semiconductor, semiconductor, (LaO)CuS, LaAgSeO, p - type, ZrCuSiAs, LaCuOTe
3. Co ₂ FeAl	heusler alloy, Co ₂ Cr _{0.6} Fe _{0.4} Al, full - heusler, Co ₂ FeAl _{1-x} Si _x , heusler, half - metallic, spin polarized, heusler compound, MFTJs, ferromagnet
4. Ca ₅ In ₂ Sb ₆	Ca ₅ Al ₂ Sb ₆ , Ca ₅ Ga ₂ Sb ₆ , zintl compound, zintl compounds, Sr ₅ In ₂ Sb ₆ , thermoelectric, thermoelectric properties, carrier concentration, carrier mobility, effective mass
5. AgBiP ₂ Se ₆	atomically thin, ferroelectricity, semiconductor, two - dimensional, ferroelectric, monolayer, devices, band edge, ground state, plane
Photovoltaics	
1. MoOHCF	electrochromic window, electrochromic, electrochromic device, vivid color, counter electrode, visible wavelengths, prussian blue, transmittance, optical transmittance, thin film
2. MoN ₂	renewable energy, appealing, applications, NIBs, great potential, realization, ion batteries, dinitride, promising, MoN ₃
3. Ni _{0.4} Co _{0.6} (OH) ₂	flexible, supercapacitors, peony - like, fabrication, step hydrothermal, strategy, CFC, carbon fiber, cloth, superior
4. NiFeS	NiVS, highly efficient, low cost, FeNiS ₂ , [NiFe], NiFeVS, efficient electrocatalyst, Ni-Fe-V, OER, promising
5. Si ₂ BN	graphenelike, have attracted, NB ₂ Si, anode material, much attention, nanostructures, hydrogen storage, battery anode, Si ₃ B ₃ N ₇ , buckled
Ferroelectric	
1. BaTiSi ₂ O ₇	spontaneous polarization, dielectric, lead - free, fresnoite, ceramics, [TiO ₅], difficult to prepare, glass, phase, properties
2. GdTiO ₃	BaTiO ₃ , BTO, SmTiO ₃ , SrTiO ₃ , gate dielectric, ferrimagnetic, mott insulator, perovskite, pyrochlore, magnetic ordering
3. TlGaTe ₂	dielectric, dependences of the permittivity, dielectric constant, TlInSe ₂ , TlInTe ₂ , relaxors, permittivity, x(TlGaTe ₂)x, ε(T), dielectric relaxation
4. Ba ₅ NdTi ₃ Ta ₇ O ₃₀	dielectric, Bi ₄ Ti ₃ O ₁₂ , tungsten – bronze, dielectrics, Ba ₄ Nd ₂ Ti ₄ Ta ₆ O ₃₀ , co-substitution, ceramics, tetragonal, x, co-modification
5. Pb ₃ Mn ₇ O ₁₅	dielectric, Pb ₂ Te ₃ O ₈ , magnetic ordering, Zn ₂ PbO ₄ , antiferromagnetic, Mn ₂ Te ₃ O ₈ , Pb ₃ Rh ₇ O ₁₅ , Pb ₂ ZnTeO ₆ , orthorhombic, manganite

The top five predictions and top ten most important context words leading to the prediction for **topological insulators, photovoltaics and ferroelectrics** using the full text corpus. A list of context words that could indicate prior study in the target domain have already been excluded in the process of making the predictions, as mentioned in Extended Data Fig. 2d. Furthermore, we have excluded any predictions for which we could find reports outside our corpus for the target application.

Extended Data Table 4 | Importance of the text corpus

Text corpus	Materials	Grammar	All	Corpus size
Wikipedia	2.6	72.8	51.0	2.81B words
Wikipedia elements	2.7	72.1	41.4	1.08B words
Wikipedia materials	2.2	72.8	41.3	781M words
All abstracts	43.3	58.3	51.0	643M words
Relevant abstracts	48.9	54.9	52.0	290M words
Pre-trained model from Kim et al ³⁹	10.4	47.1	30.8	640k papers

The top analogy scores in per cent for materials science and grammatical analogy tasks for **different corpora**. All models except Kim et al.³⁹ were trained using CBOW—continuous bag of words, the other variant of Word2vec, alongside skip-gram—with the same hyper-parameters (negative sampling loss with 15 samples, 10^{-4} downsampling, window 8, size 200, initial learning rate 0.01, 30 training epochs, minimum word count 5) and no phrases. We used the English Wikipedia dump from March 1, 2018. ‘Wikipedia elements’ corresponds to a subset of articles that mention a chemical element name (for example, ‘gold’), whereas ‘Wikipedia materials’ corresponds to a subset that mention at least one material formula. The smallest corpus on which we train our model has the best performance on materials-related analogies, whereas the largest corpus has the best performance for grammar. We believe this is due to the highly specialized nature of the relevant abstracts, suitable for the tested analogy pairs. We used the ‘Relevant abstracts’ corpus throughout this study.

In the format provided by the authors and unedited.

Unsupervised word embeddings capture latent knowledge from materials science literature

Vahe Tshitoyan^{1,3*}, John Dagdelen^{1,2}, Leigh Weston¹, Alexander Dunn^{1,2}, Ziqin Rong¹, Olga Kononova², Kristin A. Persson^{1,2}, Gerbrand Ceder^{1,2*} & Anubhav Jain^{1*}

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

Unsupervised word embeddings capture latent knowledge from materials science literature

Supplementary Information

Vahe Tshitoyan¹, John Dagdelen^{1,2}, Leigh Weston¹, Alexander Dunn^{1,2}, Ziqin Rong¹,
Olga Kononova², Kristin A. Persson^{1,2}, Gerbrand Ceder^{1,2} & Anubhav Jain¹

¹*Lawrence Berkeley National Laboratory, Berkeley 94720, California, USA*

²*Department of Materials Science and Engineering, University of California, Berkeley,
California 94720, USA*

S1 Word2vec Skip-gram

Skip-gram, one of the two variants of Word2vec, is explained schematically in fig. S1a. Assume we have $V = 500,000$ unique words in the vocabulary with each word assigned an arbitrarily index, so that it can be represented as a V -dimensional vector with zeros everywhere except that index. This representation is called one-hot encoding. Word2vec skip-gram loops through all words in the training text and uses its one-hot encoding as an input for a neural network. The task of the network is to predict all words within a certain

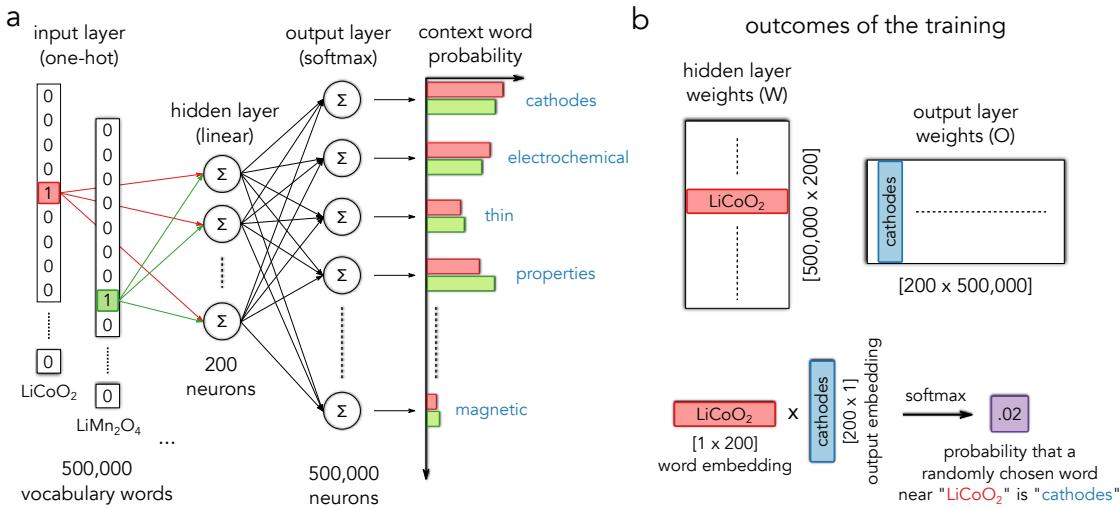


Figure S1: Word2vec skip-gram. **a.** A neural network with a single linear hidden layer learns to predict context words for every word in the vocabulary. For battery cathode materials LiCoO_2 and LiMn_2O_4 the network has to predict mostly the same context words. This results in similar hidden layer weights and therefore similar word embeddings. The softmax function is used at the output to produce normalized probabilities. **b.** Matrices W and O are the outcomes of the training, corresponding to the weights of the hidden and the output layers. Rows of W are called word embeddings, whereas columns of O are called output embeddings. The product of the two types of embeddings is the probability of the corresponding words to be used in close proximity in the text.

distance from this center word (usually ranging from 2 - 10 words away)*. While there is no single correct answer - every word occurs alongside 100s or 1000s of other words - the

*Larger word window often captures semantic relationships better, whereas smaller windows capture the syntactic relationships. We chose a relatively large window size of 8 to focus on semantic relationships, since this is more relevant for materials science relationships such as oxides of materials or common crystal structures.

end goal is not to correctly predict all neighbours but to learn compressed representations for the words. This representation is encoded in the weights of the single linear hidden layer of the neural network at the end of the training. The weights of the hidden layer are given by a $[V \times n]$ dimensional matrix W (fig. S1b), where n is the size of the space we set to “embed” the words in (200 in our case). When the one-hot encoded vector of the center word is fed into the network, all it does is select the corresponding row from matrix W . Then the output layer uses this row as an input for the softmax classifier to predict one of the neighbouring words. The classifier has to predict the same words for the words that occur in the same context, therefore, the network will adjust the corresponding rows of the matrix W to optimize this task. These row vectors are referred to as word vectors or word embeddings. Similarly, columns of the $[n \times V]$ matrix O of output weights are called output embeddings. In this notation, the task of the neural network is reduced to multiplying the row w of matrix W with the columns of matrix O and applying a softmax function, producing the probabilities of every word in the vocabulary to be next to the word w (fig. S1b).

The other variation of Word2vec is called continuous bag of words (CBOW). The neural network architecture is very similar, except instead of using the center word to predict the context words it uses the average embedding of the context words (hence, bag of words) to predict the center word. In the next section we demonstrate that Skip-

gram generally works better than CBOW for our application, therefore, we use Skip-gram throughout this work.

S2 Word2vec optimization

We tuned hyper-parameters of Word2vec to optimize its performance on the combined materials science and grammatical analogies. The full list of categorized analogies is available with supplementary materials. We found that including phrases as described in the Methods section of the main text improves the performance by approximately 4% for both CBOW and Skip-gram architectures, as shown in Table S1. We also find that Skip-gram performs approximately 4% better than CBOW both with and without phrases. We used negative sampling loss since it is faster to train. The rest of the hyperparameter optimization is summarized in table S2. We also trained GloVe embeddings¹ resulting in slightly worse performance compared to Word2vec (Table S1).

We check if analogy-based optimization leads to better performance for materials predictions using two additional metrics - one to quantify the quality of the predictions and the other for the quality of the ranking. For predictions, we use the average power factor of the first 10 predicted thermoelectrics. For the ranking, we compute the Spearman rank correlation⁵ of our ranking versus approximately 80 experimental thermoelectric figures

Algorithm	Materials	Grammar	All
Default	38.0	50.4	44.4
CBOW	48.9	54.9	52.0
CBOW + phrases	54.2	58.0	56.2
Skip-gram	54.7	58.2	56.5
Skip-gram + phrases	<u>58.9</u>	<u>61.6</u>	<u>60.3</u>
GloVe + phrases	53.8	56.0	55.0

Table S1: Algorithm choice. Top 1 analogy scores in % for materials science and grammatical analogy tasks. Each task consists of approximately 15,000 analogy pairs. The answer is considered correct only if the first nearest word matches the expected analogy. The default algorithm uses the original hyperparameters of the Word2vec code⁷, whereas the other four Word2vec algorithms use the optimized hyper-parameters. The GloVe algorithm uses the recommended parameters from the original paper¹, found to perform the best after trying to optimize the context window and the parameter alpha.

initial learning rate:	0.001	0.003	<u>0.01</u>	0.03	0.1
	50.6	54.6	<u>56.8</u>	55.1	52.6
downsampling:	10^{-3}	<u>10^{-4}</u>	10^{-5}	10^{-6}	
	56.8	<u>58.2</u>	56.5	50.6	
dimension:	100	<u>200</u>	300	400	
	54.7	<u>60.4</u>	<u>60.5</u>	59.0	
negative samples:	5	8	10	12	<u>15</u>
	59.3	59.5	59.8	59.8	<u>60.3</u>

Table S2: Hyper-parameter optimization. Top 1 analogy score in % for various hyper-parameter choices.

Only one parameter is varied while the rest are kept the same.

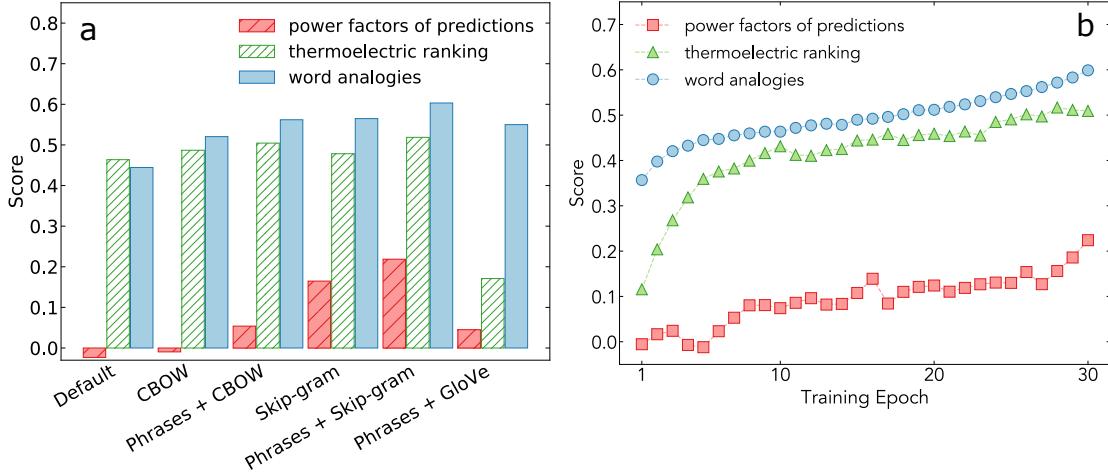


Figure S2: Accuracy of predictions. **a.** Performance metrics for different algorithms and parameters. Word analogies (blue) are analogy scores based on materials science and grammatical analogy tasks^{2–4}. Thermoelectric ranking score (green) is the Spearman rank correlation coefficient⁵ between the rank of our predictions and the experimentally measured thermoelectric figures of merit for approximately 80 materials⁶. For comparison, the correlation between the DFT and the experimental power factors from the same dataset is 0.31. The power factor score (red) is defined as $\frac{PF_{pred10} - PF_{mean}}{PF_{best10} - PF_{mean}}$, where PF_{mean} is the average power factor of all candidates, PF_{pred10} is the average power factor of the first 10 predictions and PF_{best10} is the average of the 10 highest power factors. The default algorithm uses the original hyperparameters of the Word2vec code⁷. The CBOW and Skip-gram use optimized hyperparameters with or without the common phrases. The GloVe model uses the recommended hyperparameters from the original paper¹. We found that hyper-parameter tuning changed the analogy scores for GloVe by less than a percent, however, we did not perform an extensive optimization similar to Word2vec. **b.** Evolution of the scores in a. for the “Phrases + Skip-gram” model over 30 training epochs. The learning rate decreases linearly from 10^{-2} to 10^{-4} .

of merit⁶. Fig. S2a shows the scores after 30 training epochs[†] for different models and parameters. Similar to the analogy scores, we see that Skip-gram performs better than CBOW, and that the inclusion of phrases results in performance gains. Additionally, all of the Word2vec models outperform GloVe at ranking the thermoelectrics. We attribute this to the predictive nature of Word2vec and the use of output embeddings for ranking and predictions (see the next section). GloVe is count-based and does not provide an additional set of output embeddings. Fig. S2b shows these evaluation metrics as functions of training epochs for the Skip-gram model with phrases. Until after 5 training epochs the predictions are not better than a random guess (power factor score of 0). The scores begin to improve following this initialization, and a substantial gain is made during the last few epochs of fine-tuning the embeddings. A similar trend is seen for all the metrics.

S3 Word versus output embeddings for predictions

The ranking (and consequently predictions) are performed by multiplying the embedding of the application keyword (e.g. “thermoelectric”) with the embeddings of all materials (with some count threshold, more than 3 in our case). For the application keyword we always use the normalized word embedding. However, for the materials we attempt to use either the word or the output embedding (fig. S1b). If we use word embeddings, the

[†]An epoch corresponds to a single full pass over the corpus.

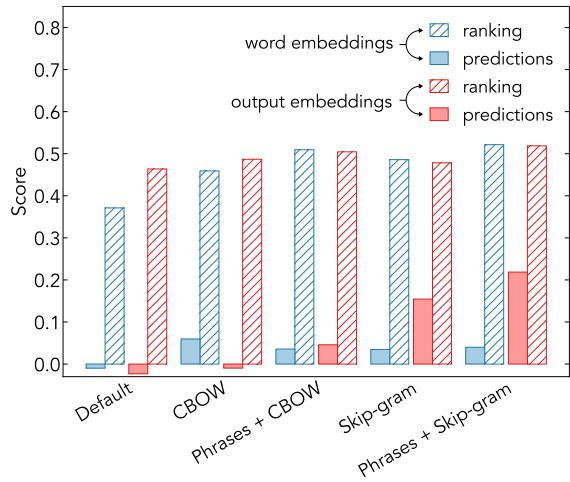


Figure S3: Word vs output embeddings.. Word embeddings corresponds to using word embeddings both for the application keyword and the material formula, whereas output embeddings corresponds to using the output embedding of the formula and the word embedding of the application keyword. The definitions of the scores are the same as in fig. S2a.

ranking is based on similarity of the application keyword and the material word. One can think of this as their interchangeability in text. If instead we use the normalized output embedding of the material, the predictions are based on the likelihood of the application keyword and the material formula being mentioned next to each other, if all materials were mentioned equal number of times in the text[‡]. This second approach generally yields better results as shown in fig. S3 and is used throughout this work.

[‡]The norm of the embedding was shown to depend on the number of mentions - with more common words usually having longer embedding vectors⁸.

S4 Word2vec element clustering versus periodic table

It is remarkable that using only relative positions of words in scientific text the algorithm learns a high dimensional representation for elements that is very similar to the periodic table when projected onto a plane. However, not all of the structure of our t-SNE projected word embeddings match well with the periodic table. Given that this is a context-based representation, it is unsurprising that the inert noble gases are far removed from the rest of the elements whereas post-transition metals, metalloids, and alkali metals, which are often used with each other in various applications, group closer together. The astute reader may observe that hydrogen is clustered with oxygen, nitrogen, and carbon; we attribute this to the fact that these elements are the main components of organic compounds. Similarly, Radon (Rn), radium (Ra) and polonium (Po), all radioactive elements, are found in closer proximity to uranium (U) and thorium (Th) in the plot than to their neighbors in the periodic table. Some elements, nevertheless, are completely out of place compared to the periodic table for what we believe to be non-physical reasons. We note that these elements' symbols overlap with common words that have the same spelling, such as “be” for beryllium, “at” for astatine or “Tc” for technetium which is also used to denote critical temperatures. Despite this, the high dimensionality of the embeddings enables relationships such as “being” - “Be” + “measure” \approx “measuring” and “BeO” - “Be” + “Mg” \approx “MgO” to be captured simultaneously, therefore, preserving both the chemical and the

syntactic relationships.

S5 Linear regression for elemental properties

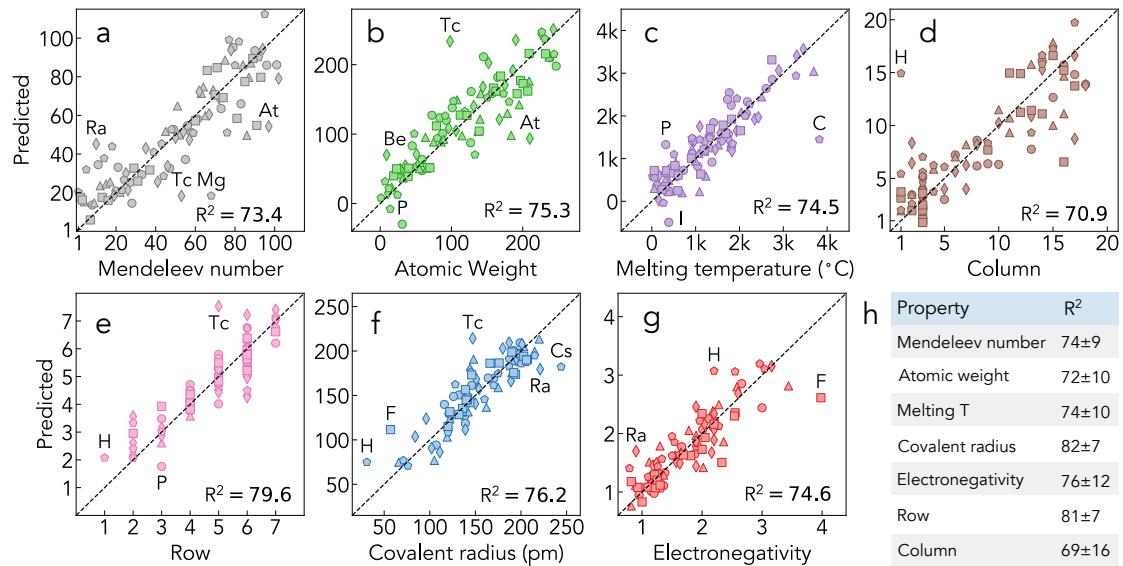


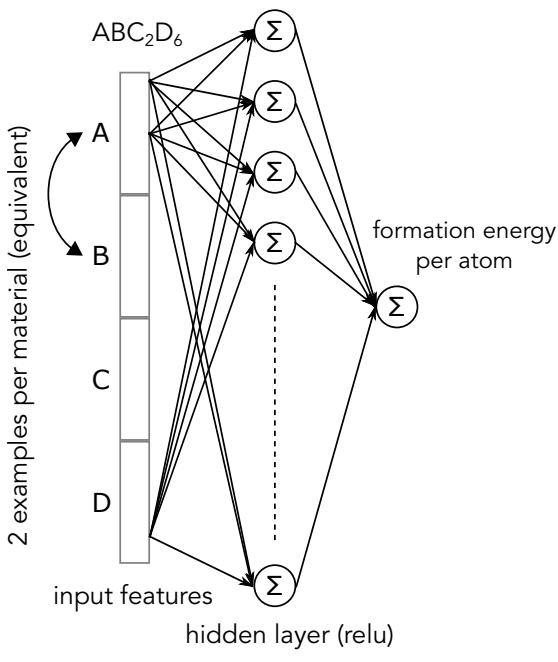
Figure S4: Predictions of Elemental Properties. **a-g.** 5-fold cross-validated predictions of 7 elemental properties using linear regression. The first 15 principal components of word embeddings of element names (e.g. “hydrogen”) were used as features. The 5 different shapes indicate the exact cross-validation splitting, such that each shape (e.g. square) represents a set of validation elements predicted using the training elements represented by the 4 other shapes (e.g. triangles, diamonds, circles, pentagons). The splitting was determined randomly. **h.** Means and standard deviations of validation R^2 scores (in percent) from 20 random 80% (training) / 20% (validation) splits.

We can determine whether there exist directions in the embedding space that correlate with elemental properties by fitting a linear regression to predict each property using

embeddings as features. We test on the following 7 elemental properties: Mendeleev number, atomic weight, melting temperature, covalent radius, electronegativity, as well as row and column in the periodic table. Since there are 200 features but only around 100 elements, even a model as simple as linear regression will overfit. To avoid this, we reduce the dimensionality to 15 by applying principal component analysis (PCA) to the normalized word embeddings. The new features are linear combinations of the original 200 and explain 65% of the total variance. Sample plots of predicted versus actual values using 5-fold cross-validation are shown in fig. S4a-g. The mean and standard deviations of R^2 for all tested properties are shown in fig. S4h. We do not perform model selection and there are no hyper-parameters to optimize, therefore, there is no need for a test set outside of the cross-validation.

S6 Formation energies of ABC_2D_6 elpasolites

We were able to predict formation energies of elpasolites with mean absolute errors as low as 55.7 meV/atom using only word embeddings (both Word2vec and Glove were tested) of their constituent elements as features. We use a dataset with approximately 10,000 ABC_2D_6 materials available from reference [9]. We use one of the simplest neural network architectures - a single fully connected hidden layer with ReLU (rectified linear unit) activation and a single output neuron (fig. S5a) - the same as reference [10]. For



a

Features	Hidden neurons	MAE (meV/atom)	
		validation	test
Word2vec embeddings	10	153.1 ± 6.7	156.3
	50	94.1 ± 3.1	87.6
	200	74.3 ± 2.1	69.8
	800	66.4 ± 1.7	62.0
	1600	67.0 ± 4.6	59.7
	3200	67.5 ± 3.6	55.7
GloVe embeddings	10	162.3 ± 38.2	145.4
	50	86.4 ± 4.7	85.8
	200	76.0 ± 6.1	65.0
	800	65.0 ± 1.9	64.0
	1600	64.9 ± 3.8	59.6
	3200	63.3 ± 3.6	56.4
Fundamental properties	10	373.3 ± 24.3	337.6
	50	304.3 ± 8.9	297.0
	200	280.4 ± 10.2	274.0
	800	268.5 ± 3.8	254.3
	1600	254.2 ± 10.0	288.7
	3200	267.3 ± 9.6	255.5
One-hot encoding	10	646.3 ± 10.8	606.0
	50	818.7 ± 16.8	760.4
	200	864.4 ± 15.3	836.4
	800	825.3 ± 19.6	795.8
	1600	808.7 ± 12.8	775.1
	3200	790.9 ± 18.1	761.2

b

Figure S5: Formation energies of ABC_2D_6 elpasolites. **a.** The architecture of the neural network used for predictions. **b.** Validation and test scores for 4 different feature choices as well as different hidden layer sizes. The performance for word embeddings does not improve much above 800 hidden neurons.

the input we concatenate embeddings of A, B, C and D elements, and also augment the data by creating 2 training examples for each material because A and B are equivalent. It is important to perform this data augmentation after splitting the data into training, validation and test sets to make sure every distinct material occurs only in one of the

sets. The mean absolute error on the test set of the best performing model decreases from 69.2 meV/atom to 55.7 meV/atom if we use this augmentation scheme. We also test alternative feature vectors with the same neural network architecture, such as one-hot encoding of elements and min-max scaled (all feature values between 0 and 1) vectors composed of the 7 elemental properties from the previous section. The performances for different features as well as different sizes of the hidden layer are summarized in fig. S5b, with word embeddings clearly performing the best. The displayed validation scores are the mean absolute errors (MAE) for 5-fold cross-validation. The test score is reported for a 10% test set separated before the training.

S7 Zero versus non-zero band gap classification

There are 1544 materials in our text corpus that have experimental band gaps in reference [11], with 603 materials having zero band gap and 941 materials having a non-zero band gap. Using 200-dimensional word embeddings of materials normalized to unit length as features, we trained a support vector machines (SVM) classifier with radial basis function (RBF) kernel to differentiate between zero vs non-zero band gap materials. Hyperparameter optimization for parameters C (regularization) and γ (inverse of the standard deviation of the kernel) was performed using grid search. An average f1-score over 20 random train / validation splits of 80% / 20% was used for scoring. The highest f1-score

of $90.8 \pm 1.0\%$ was obtained for $\gamma = 2.34$ and $C = 1.83$. In fig. S6a we plot a confusion matrix corresponding to a single 5-fold cross-validation applied to a re-shuffled (test) dataset using the optimal hyper-parameters. In fig. S6b we plot the distribution of the decision functions for these predictions, showing a good separation.

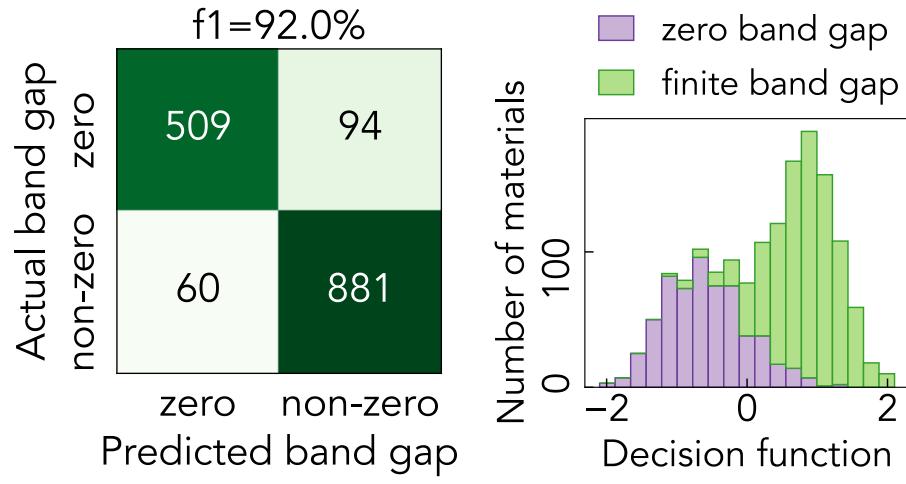


Figure S6: Prediction of zero vs non-zero band gaps. **a.** The confusion matrix of 5-fold cross-validation using hyper-parameters optimized on 20 other randomized train / validation splits. **b.** Distribution of decision function values of the 5-fold cross-validation shown in a. Values below 0 are classified as zero band gap, whereas above 0 as non-zero band gap.

S8 Material maps

Similar to chemical elements, one can visualize word embeddings of material formulas in 2D as shown in fig. S7a. We highlight a few large clusters using an unsupervised clustering

algorithm called DBSCAN¹². We also mark the most “connected” material within each cluster using PageRank¹³, an algorithm often used by search engines to rank web pages. Implementation details of DBSCAN and PageRank are discussed in the next section. If we zoom into the cluster with PbTe (fig. S7b), we see that these are all thermoelectric chalcogenides. Similarly, the cluster with LiFePO₄ contains mostly lithium-ion battery materials, the cluster with CdS is made up of materials used predominantly for solar cells, etc. We summarize the common elements in each cluster in fig. S7c and find that these elements correspond to those typically used within a particular functional application.

In addition to groups of similar materials, the long range order is also meaningful. Applications like thermoelectrics and photovoltaics merge into one another since both typically involve intermediate band semiconductors that need to be highly doped. Fig. S7c illustrates that the cluster marked as photovoltaics is composed mostly of sulfides and selenides - chalcogenides also used as thermoelectrics. Interestingly, the cluster with III-V semiconductors containing GaAs that can also be used for photovoltaics is far from CdS (II-VI semiconductor), since is not only the application that determines the position on the map but also the similarity of chemical compositions. This can be directly encoded when the name of the material is mentioned next to the formula in text, for example “gallium arsenide” next to “GaAs”. It turns out that many materials at the top of the map are oxides, bottom right are metallic, the ones stretching from the center to bottom are semiconductors

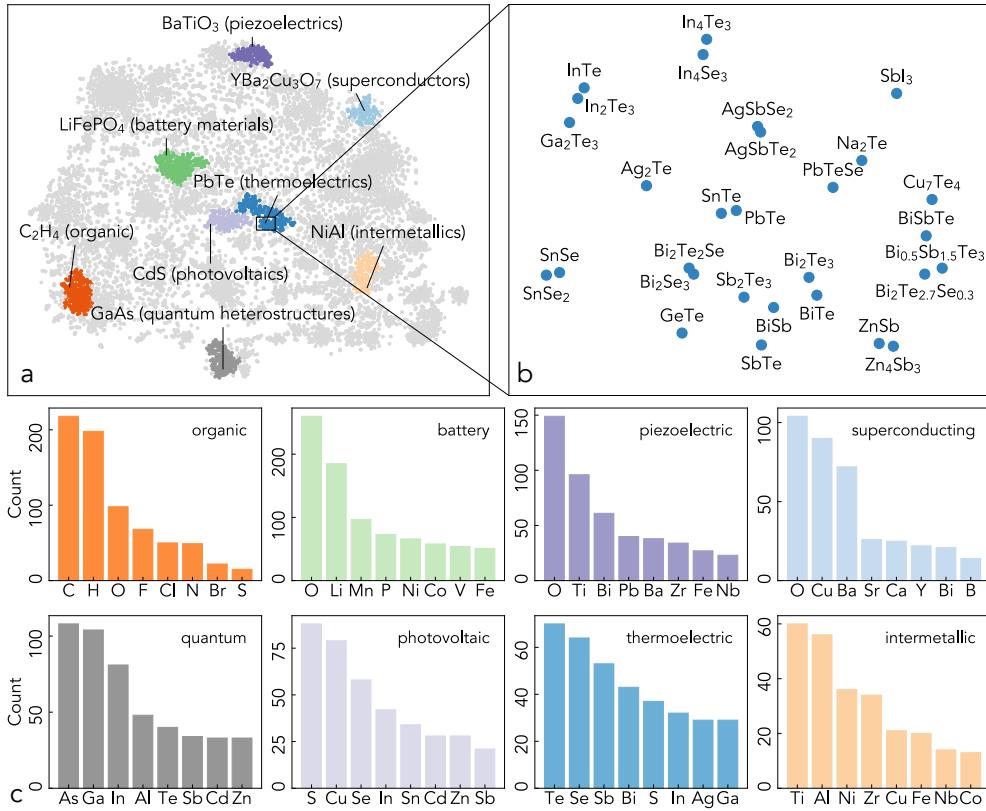


Figure S7: Material maps. **a.** t-SNE projection of 12,340 word embeddings corresponding to materials mentioned at least 10 times in the corpus. Each point represents a unique stoichiometry. The relative distance of materials can be interpreted as their context-based similarity. The materials are clustered in an unsupervised manner using DBSCAN¹², which groups together high density areas. The labeled material in each cluster corresponds to the “most connected” material within that cluster. This is determined using PageRank¹³ within each cluster, with weights corresponding to cosine similarities of word embeddings. An interactive version of the map can be found at reference [14]. **b.** A region of the map in a. in the vicinity of PbTe – one of the most common thermoelectric materials. **c.** Counts of the eight most common elements from each cluster in a., counted one per material independent on their stoichiometric ratios.

whereas the ones on the bottom left are organic. In fact, using only word embeddings of materials as features without any explicit knowledge of the compositions, we can predict if a material has a band-gap with 90.8% accuracy (f1-score), similar to a reported 91.4% score using a composition-based representation¹¹ (see Supplementary Information for the details).

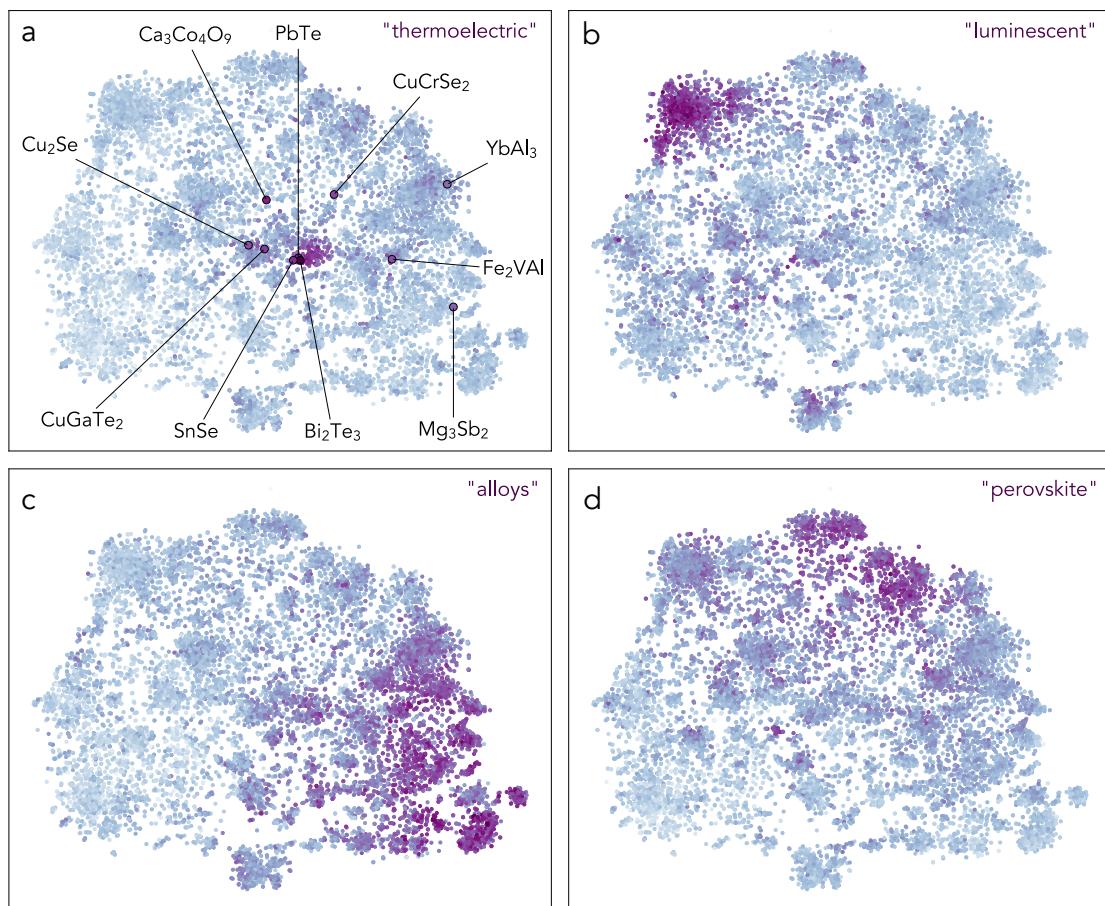


Figure S8: Dynamic material maps. Material maps highlighted according to various keywords. Darker colors correspond to more similarity.

We can create dynamic visualizations of material / keywords similarities by coloring each material on a 2D map according to its cosine similarity to that keyword - be that an application word (e.g. “thermoelectric”), a class of materials (e.g. “alloys”) or a crystal structure (e.g. “perovskite”). As an example, in fig. S8a we show the map highlighted according to the word “thermoelectric” - with darker colors corresponding to higher similarity. There are many types of thermoelectrics, hence, one should not expect all of them to cluster together. Some materials are used for other applications in other contexts, so they are further away from the main cluster containing conventional thermoelectrics such as Bi_2Te_3 and PbTe . SnSe , a recently discovered thermoelectric with a record power factor in 2014¹⁵ is also in this cluster. CuGaTe_2 is a well-known semiconductor also considered as a promising candidate for thin film solar cells¹⁶. Mg_3Sb_2 is a thermoelectric with Zintl structure¹⁷, Cu_2Se is an recently discovered ion-liquid like thermoelectric¹⁸, $\text{Ca}_3\text{Co}_4\text{O}_9$ is an oxide thermoelectric¹⁹, Fe_2VAl is a heusler-type nonmagnetic semimetal²⁰, CuCrSe_2 is a layered antiferromagnet and a superionic conductor²¹ whereas YbAl_3 is an intermetallic compound with a record power factor²². More extensive reviews of different types of thermoelectrics can be found at references [23] and [24]. Examples for a few other keywords are plotted in fig. S8b-d.

S9 Projection, clustering and ranking of materials embeddings.

Projection. In fig. S7a, the 200-dimensional embeddings of the 12,340 materials mentioned more than 10 times in our corpus were reduced to 2 dimensions using t-SNE^{25,26}. We used cosine distance between the embeddings as a metric, perplexity 30, learning rate 200, early exaggeration 12.0 and 10,000 iterations - with coordinates initialized using PCA.

Clustering. To group high density areas of the 2D projection for easier visualization, we used an unsupervised clustering technique called DBSCAN^{12,26}. We used neighbour distance cutoff $\epsilon = 2.75$ and a minimum count of 8, producing well separated clusters. Clusters with less than 120 materials were ignored. For the final visualization, we chose 8 clusters from the remaining 18.

Ranking. To find a representative material within each cluster, we use the implementation of PageRank¹³ available via igraph²⁷ software package. Each node of the undirected graph corresponds to a material, with the weights of the edges corresponding to cosine similarities between the materials. We used the default damping value of 0.85 to compute the ranks within each cluster, with the highest ranked materials labelled in fig. S7a. Globally, the five most connected materials in our corpus (excluding chemical elements) were TiO₂,

ZnO, SiO₂, Al₂O₃ and SiC, which are all used for a large spectrum of applications.

S10 Unconventional thermoelectric predictions

In addition to well known thermoelectric material classes, we observe predictions such as KAg₂SbS₄ (see Table 2 of extended data) that do not have strong similarity to known thermoelectrics. This particular compound has recently been suggested as a candidate photovoltaic material²⁸. Another example is BiOCl, an atypical oxychloride which was the top prediction in the 2010 historical corpus (Supplementary Table S5) and has a computed p-type power factor of $25.4 \mu\text{W}/\text{K}^2 \cdot \text{cm}$ (calculated using the constraints described in the Methods section of the main text) – ranking in the 93rd percentile of our dataset’s power factors. Potential issues with the oxychloride chemistry might be large band gaps and dopability. However, this material contains desirable band structure features, including two doubly-degenerate valence band peaks aligned at the off-symmetry points X and R²⁹ that are responsible for the high computed power factor. ZnSiP₂ was #3 in 2005 and has a similarly high p-type power factor of $33.2 \mu\text{W}/\text{K}^2 \cdot \text{cm}$ (95th percentile) and n-type power factor of $29.5 \mu\text{W}/\text{K}^2 \cdot \text{cm}$ (96th percentile among n-types power factors). Phosphides are typically thought to have high thermal conductivities not desirable for thermoelectric applications, however, this is not strictly true³⁰ and it is unclear if that is the case for this compound. This material is also interesting due to its band structure features, which include

a triply degenerate valence band peak and doubly-degenerate conduction band pocket at the gamma point³¹, as well as a degenerate conduction band pocket between Z and Σ_1 . Another notable example is $\text{Nd}_{0.5}\text{Sr}_{0.5}\text{MnO}_3$ (#5 in 2004). The computed power factor for this compound is missing from our dataset since it exhibits site disorder that is more difficult to model with density functional theory methods. However, it is known experimentally to have a relatively low thermal conductivity ($< 3 \text{ W/K} \cdot \text{m}$ at 300 K)³², high dopability, and high electrical conductivity ($> 300 \text{ S/cm}$ at 300 K)³³ – all promising indicators for a high zT material. Each of these examples may require additional synthesis and optimization work to overcome potential limitations in doping and thermal conductivity due to their unconventional chemistries, nevertheless, they are viable thermoelectric candidates that are not closely related to any mainstream thermoelectrics.

Year	Top 10 thermoelectric predictions	Total potential	Total abstracts
		predictions	in corpus
2001	HgMnTe, HgZnTe, EuLiH ₃ , CdGeP ₂ , La _{0.5} Sr _{0.5} MnO ₃ , VB ₂ , CoCr ₂ S ₄ , CdSeTe, Bi ₂ Sr ₂ CuO ₆ , AgInS ₂	13221	288178
2002	Mo ₃ Te ₄ , HgMnTe, ZrB ₂ , ZrSi ₂ , La _{0.5} Sr _{0.5} MnO ₃ , Mo ₅ Si ₃ , Ge ₂₂ Se ₇₈ , TmSb, BaLaCuO, Nd _{0.5} Sr _{0.5} MnO ₃	14181	331414
2003	EuB ₆ , CdGeP ₂ , HgMnTe, ReSe ₂ , Cd _{0.8} Zn _{0.2} Te, Yb ₄ As ₃ , HgZnTe, ReS ₂ , CoCr ₂ S ₄ , CuNb	15042	375079
2004	HgMnTe, V ₂ Ga ₅ , HgZnTe, Yb ₄ As ₃ , Nd _{0.5} Sr _{0.5} MnO ₃ , CoS ₂ , EuB ₆ , CdGeP ₂ , ReS ₂ , Ge ₂₂ Se ₇₈	15906	422439
2005	V ₂ Ga ₅ , BaSi ₂ , ZnSiP ₂ , HgZnTe, HgMnTe, CoCr ₂ S ₄ , EuB ₆ , Sb ₂ O ₅ , ReS ₂ , SbSI	16824	473567

2006	ReS ₂ , BaSi ₂ , TiSi, SmInO ₃ , ReSe ₂ , 17595	523433
	Na _{0.9} Mo ₆ O ₁₇ , HgMnTe, HgZnTe,	
	CeOs ₄ Sb ₁₂ , CoCr ₂ S ₄	
2007	ReS ₂ , HgZnTe, BaSi ₂ , SmInO ₃ , ReSe ₂ , 18510	580323
	EuB ₆ , LaOAgS, CeOs ₄ Sb ₁₂ , CdP ₂ , Sn ₄ P ₃	
2008	ReS ₂ , HgZnTe, ReSe ₂ , SbSI, GeI ₂ , SmInO ₃ , 19320	639825
	FeIn ₂ Se ₄ , Yb ₄ As ₃ , TeCl ₄ , CdIn ₂ Te ₄	
2009	HgZnTe, ReS ₂ , SmInO ₃ , CdIn ₂ Te ₄ , 20177	702186
	CuGaTe ₂ , ReSe ₂ , HgMnTe, TlSbSe ₂ ,	
	Co ₂ FeGa, (YbS) _{1.25} CrS ₂	
2010	BiOCl, HgZnTe, Co ₂ FeGa, CdIn ₂ Te ₄ , 21037	766690
	HgMnTe, (YbS) _{1.25} CrS ₂ , La _{0.9} Sr _{0.1} MnO ₃ ,	
	ReS ₂ , NiTe ₂ , NiP ₃	
2011	SmInO ₃ , CdIn ₂ Te ₄ , (YbS) _{1.25} CrS ₂ , 21679	831227
	FeIn ₂ Se ₄ , HgZnTe, NiP ₃ , CdGa ₂ O ₄ ,	
	AgInSe ₂ , La _{0.9} Sr _{0.1} MnO ₃ , ZrNCl	
2012	FeIn ₂ Se ₄ , (YbS) _{1.25} CrS ₂ , HgZnTe, 22446	904141
	CdGa ₂ O ₄ , YbTe, TlCrS ₂ , SmInO ₃ , HoCu ₂ ,	
	CdIn ₂ Te ₄ , SrNb _{0.01} Ti _{0.99} O ₃	

2013	HgZnTe, Zn _{0.7} Cd _{0.3} Se, CaYb ₂ S ₄ , CdIn ₂ Te ₄ , FeIn ₂ Se ₄ , CuSbS ₂ , CdGa ₂ O ₄ , TeCl ₄ , CeTe, (YbS) _{1.25} CrS ₂	23179	979040
2014	TlCrS ₂ , YbTe, HgZnTe, FeIn ₂ Se ₄ , SnSb ₂ Te ₄ , La _{0.7} Ca _{0.2} Sr _{0.1} MnO ₃ , CaYb ₂ S ₄ , HoCu ₂ , Bi _{0.95} La _{0.05} FeO ₃ , CdSnO ₃	24029	1067395
2015	TlCrS ₂ , YbTe, FeIn ₂ Se ₄ , ReS ₂ , Zn _{0.7} Cd _{0.3} Se, Co ₂ FeGa, ReSe ₂ , NiP ₃ , CoCrFeNi, LiGaSe ₂	24749	1159529
2016	YbTe, TeCl ₄ , FeIn ₂ Se ₄ , TlCrS ₂ , La _{0.7} Ca _{0.2} Sr _{0.1} MnO ₃ , Hf _{0.2} Zr _{0.8} O ₂ , CdSnP ₂ , MoSe, Bi _{0.95} La _{0.05} FeO ₃ , Pb _{0.902} Sn _{0.098} Se	25469	1257788
2017	YbTe, CdSnP ₂ , In ₃ Se ₂ , Hf _{0.2} Zr _{0.8} O ₂ , InFeZnO ₄ , TlSbSe ₂ , Sc ₂ CF ₂ , HgZnTe, Ag ₃ AuSe ₂ , TlCrS ₂	26184	1358468
2018	YbTe, In ₃ Se ₂ , ZnSnP ₂ , HgZnTe, TlSbSe ₂ , CdSnP ₂ , CuTe, TlCu ₂ Se ₂ , SbSI, MoSe	26804	1470230

Table S3: Top 10 thermoelectric predictions from each year. The list of materials is ordered from prediction #1 to prediction #10. Total candidates is the number of materials considered for the prediction, which includes all materials mentioned more than 3 times but not studied as thermoelectric before. Total abstracts is the number of (relevant) abstracts used to train the word embeddings.

1. Pennington, J., Socher, R. & Manning, C. Glove: Global Vectors for Word Representation. *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)* 1532–1543 (2014).
2. Mikolov, T., Corrado, G., Chen, K. & Dean, J. Efficient Estimation of Word Representations in Vector Space. In *International Conference on Learning Representations* (2013).
3. Mikolov, T., Sutskever, I., Chen, K., Corrado, G. & Dean, J. Distributed Representations of Words and Phrases and their Compositionality. In *International Conference on Neural Information Processing Systems*, 3111–3119 (2013).
4. <http://www.materialsintelligence.com/materials-analogies>.
5. Spearman, C. The proof and measurement of association between two things. *The American Journal of Psychology* **15**, 72–101 (1904).
6. Gaulois, M. W. *et al.* Data-driven review of thermoelectric materials: Performance and resource onsiderations. *Chemistry of Materials* **25**, 2911–2920 (2013).
7. <https://code.google.com/archive/p/word2vec/>.
8. Schakel, A. M. & Wilson, B. J. Measuring word significance using distributed representations of words. *arXiv preprint arXiv:1508.02297* (2015).

9. Faber, F. A., Lindmaa, A., Von Lilienfeld, O. A. & Armiento, R. Machine Learning Energies of 2 Million Elpasolite (ABC2D6) Crystals. *Physical Review Letters* **117**, 2–7 (2016). 1508.05315.
10. Zhou, Q. *et al.* Learning atoms for materials discovery. *Proceedings of the National Academy of Sciences* **115**, E6411–E6417 (2018).
11. Zhuo, Y., Mansouri Tehrani, A. & Brgoch, J. Predicting the Band Gaps of Inorganic Solids by Machine Learning. *Journal of Physical Chemistry Letters* **9**, 1668–1673 (2018).
12. Ester, M., Kriegel, H.-P., Sander, J. & Xu, X. A Density-based Algorithm for Discovering Clusters a Density-based Algorithm for Discovering Clusters in Large Spatial Databases with Noise. In *KDD'96 Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, 226–231 (1996).
13. Brin, S. & Page, L. The anatomy of a large-scale hypertextual web search engine. *Computer Networks and ISDN Systems* **30**, 107 – 117 (1998). Proceedings of the Seventh International World Wide Web Conference.
14. <http://www.materialsintelligence.com/materials-map>.
15. Zhao, L.-D. *et al.* Ultralow thermal conductivity and high thermoelectric figure of merit in SnSe crystals. *Nature* **508**, 373–377 (2014).

16. Plirdpring, T. *et al.* Chalcopyrite CuGaTe₂: A high-efficiency bulk thermoelectric material. *Advanced Materials* **24**, 3622–3626 (2012).
17. Condron, C. L., Kauzlarich, S. M., Gascoin, F. & Snyder, G. J. Thermoelectric properties and microstructure of Mg₃Sb₂. *Journal of Solid State Chemistry* **179**, 2252–2257 (2006).
18. Liu, H. *et al.* Copper ion liquid-like thermoelectrics. *Nature Materials* **11**, 422–425 (2012).
19. Funahashi, R. *et al.* An oxide single crystal with high thermoelectric performance in air. *Japanese Journal of Applied Physics* **39**, L1127–L1129 (2000).
20. Xu, B. *et al.* The structural, elastic and thermoelectric properties of Fe₂VAl at pressures. *Journal of Alloys and Compounds* **565**, 22–28 (2013).
21. Bhattacharya, S. *et al.* CuCrSe₂: A high performance phonon glass and electron crystal thermoelectric material. *Journal of Materials Chemistry A* **1**, 11289–11294 (2013).
22. Liang, J., Fan, D., Jiang, P., Liu, H. & Zhao, W. First-principles study of the thermoelectric properties of intermetallic compound YbAl₃. *Intermetallics* **87**, 27–30 (2017).

23. Liu, W. *et al.* New trends, strategies and opportunities in thermoelectric materials: A perspective. *Materials Today Physics* **1**, 50–60 (2017).
24. He, J. & Tritt, T. M. Advances in thermoelectric materials research: Looking back and moving forward. *Science* **357** (2017).
25. van der Maaten, L. & Hinton, G. Visualizing high-dimensional data using t-sne. *Journal of Machine Learning Research* **9**, 2579–2605 (2008).
26. <http://scikit-learn.org/stable/>.
27. <http://igraph.org>.
28. Nhalil, H. *et al.* Optoelectronic properties of candidate photovoltaic cu₂pbsis4, ag₂pbges4 and kag₂sbs4 semiconductors. *Journal of Alloys and Compounds* **746**, 405–412 (2018).
29. <https://www.materialsproject.org/materials/mp-22939/>.
30. Pöhls, J.-H. *et al.* Metal phosphides as potential thermoelectric materials. *Journal of Materials Chemistry C* **5**, 12441–12456 (2017).
31. <https://www.materialsproject.org/materials/mp-4763/>.

32. Fujishiro, H., Sugawara, S. & Ikebe, M. Anomalous phonon transport enhancement at first-order ferromagnetic transition in (gd, sm, nd) 0.55 sr0.45mno3. *Physica B: Condensed Matter* **316**, 331–334 (2002).
33. Cui, C., Tyson, T. A., Chen, Z. & Zhong, Z. Transport and structural study of pressure-induced magnetic states in nd 0.55 sr 0.45 mno 3 and nd 0.5 sr 0.5 mno 3. *Physical Review B* **68**, 214417 (2003).