


ARTICLE OPEN

Automated generation and ensemble-learned matching of X-ray absorption spectra

Chen Zheng¹, Kiran Mathew², Chi Chen¹, Yiming Chen¹, Hanmei Tang¹, Alan Dozier³, Joshua J. Kas⁴, Fernando D. Vila⁴, John J. Rehr⁴, Louis F. J. Piper^{5,6}, Kristin A. Persson² and Shyue Ping Ong¹ 

X-ray absorption spectroscopy (XAS) is a widely used materials characterization technique to determine oxidation states, coordination environment, and other local atomic structure information. Analysis of XAS relies on comparison of measured spectra to reliable reference spectra. However, existing databases of XAS spectra are highly limited both in terms of the number of reference spectra available as well as the breadth of chemistry coverage. In this work, we report the development of XASdb, a large database of computed reference XAS, and an Ensemble-Learned Spectra Identification (ELSIE) algorithm for the matching of spectra. XASdb currently hosts more than 800,000 K-edge X-ray absorption near-edge spectra (XANES) for over 40,000 materials from the open-science Materials Project database. We discuss a high-throughput automation framework for FEFF calculations, built on robust, rigorously benchmarked parameters. FEFF is a computer program uses a real-space Green's function approach to calculate X-ray absorption spectra. We will demonstrate that the ELSIE algorithm, which combines 33 weak "learners" comprising a set of preprocessing steps and a similarity metric, can achieve up to 84.2% accuracy in identifying the correct oxidation state and coordination environment of a test set of 19 K-edge XANES spectra encompassing a diverse range of chemistries and crystal structures. The XASdb with the ELSIE algorithm has been integrated into a web application in the Materials Project, providing an important new public resource for the analysis of XAS to all materials researchers. Finally, the ELSIE algorithm itself has been made available as part of *veidt*, an open source machine-learning library for materials science.

npj Computational Materials (2018)4:12; doi:10.1038/s41524-018-0067-x

INTRODUCTION

X-ray absorption spectroscopy (XAS) is a widely used technique in the study of the properties, physical states, and local environments of materials.^{1–3} When incident X-ray photons with energy greater than the binding energy are absorbed by an atom, a core-level electron is removed from its quantum level. In XAS, the absorption coefficient, $\mu(E)$ is measured as a function of X-ray energy E . Detailed descriptions of X-ray absorption theory and equation have been included in many excellent books and review papers.^{4,5}

The X-ray absorption fine structure (XAFS) is typically divided in to two regimes: X-ray absorption near-edge structure (XANES) and extended X-ray absorption fine structure (EXAFS).⁶ The XANES is a fingerprint of the oxidation states and coordination chemistries of the absorbing atom. Quantitative XANES analyses are typically difficult and are usually conducted in combination with principle component analysis or least-squares fitting. The EXAFS provides local atomic structure information, which can be extracted via coupling with theoretically calculated XAFS spectra using well-established software packages.⁷ One of the main challenges of interpreting XANES and EXAFS lies in obtaining reference spectra to fit the unknown spectra; measuring XAFS spectroscopy experimentally is laborious and time-consuming, requiring X-ray

beams of finely tunable energy that are accessible only through synchrotron radiation facilities.⁵ To the authors' knowledge, open reference databases usually contain at most hundreds of XAS spectra. For example, the electron energy-loss spectroscopy (EELS) database⁸ initiated in the 1990s contains 271 spectra, but only 21 of which are XAS spectra and 17 of which are K-edge spectra. EELS is theoretically equivalent to X-ray absorption⁹ under common acquisition conditions, but is of lower quality in terms of signal to noise ratio and energy resolution. Most XAS data are available only via publications in the literature, which cannot be extracted easily for comparison.

In recent years, theoretical calculations of XAFS have become more accurate and accessible due to the successful development of ab initio codes, such as the FEFF program,^{10,11} as well as advances in computing power. In this work, we will discuss the development of a high-throughput framework to generate a reference XAS database (XASdb) for all materials in the Materials Project¹² database. This framework combines the power of the Python Materials Genomics (pymatgen) materials analysis library¹³ with the FireWorks workflow management software¹⁴ to carry out hundreds of thousands of XAFS calculations using the FEFF9 code.¹⁰ This framework has been implemented in the Atomate package.¹⁵ More importantly, we have developed a novel automated XANES spectra matching algorithm that leverages

¹Department of NanoEngineering, University of California San Diego, 9500 Gilman Drive, Mail Code 0448, La Jolla, CA 92093-0448, USA; ²Department of Materials Science, University of California Berkeley, Berkeley, CA 94720, USA; ³Division of Applied Research and Technology, National Institute for Occupational Safety and Health, Centers for Disease Control, Cincinnati, OH 45226, USA; ⁴Department of Physics, University of Washington, Seattle, WA 98195, USA; ⁵Department of Physics, Applied Physics and Astronomy, Binghamton University, Binghamton, NY 13902, USA and ⁶Materials Science & Engineering, Binghamton University, Binghamton, NY 13902, USA

Correspondence: Kristin A. Persson (kapersson@lbl.gov) or Shyue Ping Ong (ongsp@eng.ucsd.edu)

These authors contributed equally: Chen Zheng, Kiran Mathew, Chi Chen

Received: 30 September 2017 Revised: 8 February 2018 Accepted: 12 February 2018

Published online: 20 March 2018

ensemble learning techniques to identify similar XANES spectra from our computed reference XASdb. We believe the combination of the XASdb with these machine-learned spectra matching tools will be an invaluable resource to the materials research community by greatly enhancing the efficiency at which experimental XAS spectra can be analyzed. It should be noted that this work primarily focuses on common K-edge XANES spectra; higher edge XANES and EXAFS computations and analysis are currently ongoing and will be discussed in future publications.

RESULTS AND DISCUSSION

We have selected the latest version (v9) of the popular FEFF program as our software of choice in this work. FEFF is a program for ab initio multiple-scattering calculations of XAFS and various other spectra for clusters of atoms. This choice is motivated by three factors: (i) FEFF-computed spectra has been shown to yield excellent agreement with experimentally measured spectra in a broad range of studies;^{16–18} (ii) FEFF calculations are relatively inexpensive compared to other approaches for computing XAS spectra (e.g., a typical FEFF calculation takes <1 h on a single node, while multi-day, multi-core calculations are necessary for DFT-based spectra calculations); and (iii) FEFF requires minimal adjustable parameters. These three advantages make FEFF an ideal candidate for automation to generate XAS spectra across a broad range of chemistries. A key step in any automation framework is benchmarking of computational parameters for convergence and accuracy. The benchmarking dataset and criteria are detailed in the Methods section. The Pearson correlation coefficient, as given by the following expression, is used as the benchmarking criterion.

$$S_{\text{Pearson}}(X, Y) = \frac{\sum_{i=1}^D (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\left(\sum_{i=1}^D (X_i - \bar{X})^2\right)\left(\sum_{i=1}^D (Y_i - \bar{Y})^2\right)}}, \quad (1)$$

where X_i and Y_i represent the absorption coefficients of two spectra on the same energy grid. The value of S_{Pearson} can range from -1 to 1 , with a value of 1 being a perfect match. Used in this context, the Pearson correlation coefficient is a similarity metric, i.e., it measures the degree of similarity between two spectra.

We have tested the convergence of the FEFF calculated spectra with respect to four parameters: the radius of the cluster considered in the full multiple-scattering (FMS) calculation (self-consistent field (SCF) rfms1), the total number of multiple-scattering paths considered (FMS rfms), the exchange-correlation potential (EXCHANGE), and the treatment of the core (COREHOLE) (see Methods for a detailed description of the FEFF input file).

The SCF rfms1 was varied from 2 to 8 \AA , and the spectrum at the highest value (8 \AA) was set as the reference for each material. Figure 1 shows the computed Pearson correlation coefficients between spectra computed at lower rfms1 and the reference. We find that the computed spectra are converged ($S_{\text{Pearson}} > 0.95$) at around $\text{rfms1} = 6 \text{ \AA}$ for all material, though the Al K-edge for aluminum nitride is converged only for $\text{rfms1} = 6.5 \text{ \AA}$. Given that the computational cost increases substantially for $\text{rfms1} > 7 \text{ \AA}$ (see Supplementary Fig. 1), we have chosen $\text{rfms1} = 7 \text{ \AA}$ as the default setting for SCF in the high-throughput XANES computations.

The rfms field in the FMS card was varied from 3.0 to 11.0 \AA at 1.0 \AA intervals, and the spectrum at the highest value (11 \AA) is set as the reference for each material. We find that the computed spectra are converged ($S_{\text{Pearson}} > 0.95$) around $\text{rfms} = 9 \text{ \AA}$ for all materials (see Supplementary Fig. 2(a)). Since the computational cost increases substantially for $\text{rfms} > 9 \text{ \AA}$ (see Supplementary Fig. 2(b)), we have chosen $\text{rfms} = 9 \text{ \AA}$ as the default setting for FMS in the high-throughput XANES computations.

In FEFF9, two approximations of the core-hole potentials have been implemented, i.e., a fully screened potential based on the

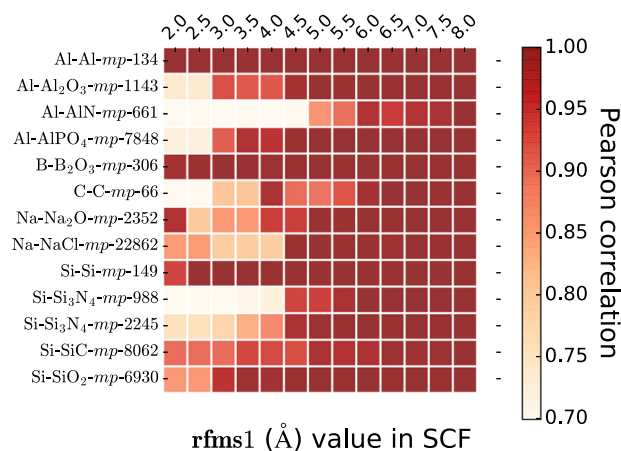


Fig. 1 Benchmarking results of rfms1 parameter in the SCF card for K-edge XANES of various materials. The rfms1 parameter specifies the radius of the cluster considered for the full multiple scattering during self-consistent potential calculations. Pearson correlation coefficients were calculated between spectra calculated at different rfms1 and the reference calculated at $\text{rfms1} = 8.0 \text{ \AA}$

final-state rule (FSR) and a linear random-phase-approximation (RPA) screening. Systematic reviews of these two approaches have been done by Rehr et al.¹⁹ We evaluated the performance of all three core-hole options in FEFF9 on the computed K-edge XANES. As shown in Supplementary Fig. 3(a), spectra obtained using both the FSR and RPA are in much better agreement with experimental results than ones without core-hole treatment. The spectra computed without a core-hole treatment lack the edge enhancement observed in the experiments. In general, spectra obtained using FSR and RPA are similar (Supplementary Fig. 3(b)). We have chosen RPA screening as the default setting for the high-throughput XANES computations as the FSR might breakdown for the L-shell metals.²⁰

Similar evaluations of the EXCHANGE card options reveal that the default Hedin-Lundquist model is the best option (see Supplementary Fig. 4).

Sensitivity of computed XAS spectra to lattice parameters

The FEFF code uses a self-consistent DFT calculation of the Fermi energy based on the real-space Green's function (RSGF) approach with muffin-tin potentials for a given lattice structure. Comparing to the full-potential calculations, we find that the FEFF calculation of the densities of states is typically in fairly good agreement with DFT for many materials. In the Materials Project, the Perdew-Berke-Ernzerhof (PBE)²¹ generalized gradient approximation functional was used as the default for all relaxation calculations. As it is well known that PBE leads to systematic errors of up to 5% in the lattice parameters (with a tendency to overestimate),^{22–25} we tested the sensitivity of computed XANES spectra to $\pm 5\%$ changes in the lattice parameters. The results are shown in Fig. 2.

We find that the Fermi energy level of the spectrum is sensitive to the lattice parameter variation (Fig. 2a). The Fermi energy level shifts towards lower energy as the lattice parameter increases, while the spacing of the spectral features contracts at the same time. An example for Na K-edge of Na_2O is shown in Fig. 2b, and additional examples are available in Supplementary Fig. 5.

A portion of the Fermi energy shift can be attributed to the artifacts of the FEFF's potential approximation model (see Supplementary Fig. 8). Nevertheless, the shape of the spectra remains unchanged. While different corrections to eliminate the artificial component of the dependence have been reported,²⁶ these approaches are not amenable to a high-throughput approach. Here, we note that due to the approximations used in

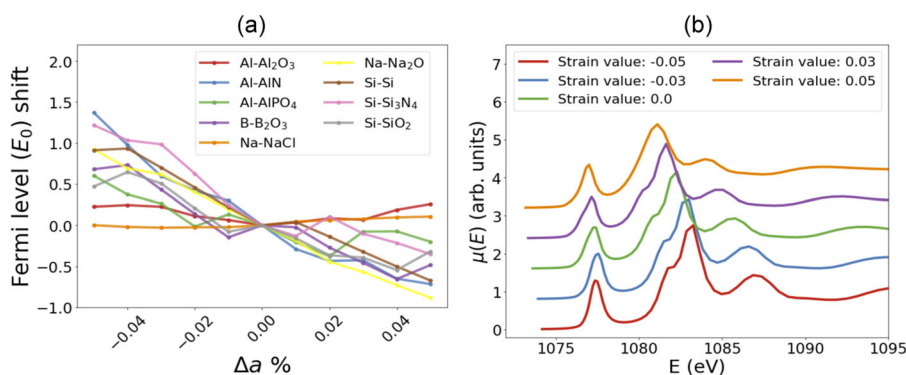


Fig. 2 **a** Relationship between the Fermi energy level of K-edge XANES and **a** lattice parameter changes. Fermi energy levels of the unstrained structures are used as references. **b** Visualization of Na K-edge XANES spectra in Na₂O (mp-2352) calculated with different applied strain values

FEFF, we need to calibrate the Fermi level with experimental spectra. Therefore, a pure energy shift only translates to an energy calibration value in the post processing.

In summary, the PBE-relaxed structures from the Materials Project can be used as the input for high-throughput XANES calculations, even though there are other functionals^{27,28} that may provide better lattice parameters estimates.^{29–32}

Workflow and database

Using the high-throughput parameters outlined above, we developed a high-throughput workflow for FEFF XAS calculations within the open source computational materials science workflow package Atomate.¹⁵ Atomate provides a high-level interface to compose workflows using the widely used open source materials science software such as Pymatgen,¹³ FireWorks,¹⁴ and Custodian. The proposed default FEFF9 parameters have been implemented as “input sets” in Pymatgen,¹³ which ensures reproducible and automated generation of standardized input files for any material. The compounds used in the high-throughput spectra generation were obtained from the Materials Project database.¹² For each compound, the K-edge XANES spectrum was computed with each symmetrically unique site in the structure as the absorbing atom.

All computed spectra, as well as accompanying meta-data (e.g., input structure, absorbing atom, materials project id, etc.), are stored in a MongoDB database for on-demand querying and retrieval of data. So far, K-edge XANES spectra have been computed for more than 40,000 unique materials in the Materials Project database, which amounts to over 800,000 K-edge spectra. This is by far the largest repository of XANES spectra in the world, and is growing rapidly. Future plans include the calculation of XANES for L, M, and N shells as well as EXAFS spectra.

Spectra matching using ensemble learning

To extract the most utility and power from the XASdb, we have developed a novel Ensemble-Learned Spectra Identification (ELSIE) algorithm that allows for rapid identification of matching spectra for any experimental XAS spectra. The main goal of spectral matching is to obtain a list of compounds (the “hit list”) whose spectra are most similar to that of the target spectrum. The success and failure of matching is defined by the characteristics of the spectrum. In the case of XANES spectra, the relevant information to be extracted is the coordination environment and oxidation state of the absorbing atom. As multiple materials can have atoms in the same oxidation state and coordination environment, we define the matching to be successful if the correct coordination environment and oxidation state are within the top entries.

The ELSIE algorithm uses the ensemble method to improve the robustness of XAS identification. In ensemble learning, the core

concept is the combination of multiple weak learners to achieve superior performance. It relies on the assumption that each weak learner is better than a random guess, and each weak learner captures different aspects of the problem. At the core of the algorithm is the process of building individual weak learners. Taking inspiration from the spectra matching algorithms for Raman spectroscopy³³ and other spectra,^{34,35} we broke down the problem of matching XAS spectra into two main steps, namely preprocessing and similarity computations. We define each weak learner to be a combination of a preprocessor (a specific series of preprocessing steps) with a similarity metric. Figure 3 provides an overview of the ELSIE algorithm (see Methods section for the details on the construction of the ELSIE algorithm).

We evaluated the ELSIE algorithm using 13 XANES spectra from EELSdb (Supplementary Table 1), supplemented by six high-quality experimental XANES spectra of V₂O₅, V₂O₃, VO₂, LiNiO₂, LiCoO₂, and NiO from previous studies.^{36,37} The inclusion of this latter dataset is motivated by our desire to improve the diversity of the test data, especially with regards to transition metal species.

The first step is to narrow down the candidate computed reference spectra by the absorbing element (A). Though this information is usually known a priori, the characteristic XAS absorption edge energy follows a power law with the atomic number,^{5,6} which leads to clearly separated energy ranges. Hence, we can identify the absorbing element with 100% accuracy by comparing the energy range of the target spectrum to tabulated X-ray absorption edge data.³⁸

Once the absorbing element A is identified, the computed spectra of all materials within the same chemical system are queried from the XASdb. For example, for the Al K-edge of Al₂O₃, we include the Al K-edge spectra of all Al and Al_xO_y materials as reference spectra. We excluded compounds with energy above hull (E_{hull}) larger than 100 meV/atom since they are not likely to be stable.³⁹ For C K-edge XANES of the diamond structure ($Fd\bar{3}m$), we relaxed the constraint to 200 meV/atom as the corresponding entry (mp-66, diamond) has an E_{hull} of 136 meV/atom. It should be noted that though the individual absorption spectrum for each symmetrically distinct site was computed for all crystal structures in the Materials Project database, the reference spectra used for comparison with the target spectra are constructed by summing these individual spectra taking into account the site multiplicities.

To evaluate the overall performance of ELSIE, we looked at three key metrics: (i) whether the correct structure is within the top 5 ranked computed spectra, (ii) whether the top ranked entry has the absorbing species in the correct oxidation state, and (iii) whether the top ranked entry has the absorbing species in the correct coordination environment, i.e., coordination number and geometry. Where the exact structural information is not available (e.g., in the experimental spectra from EELSdb), it is assumed that those spectra correspond to the ground state structures in the

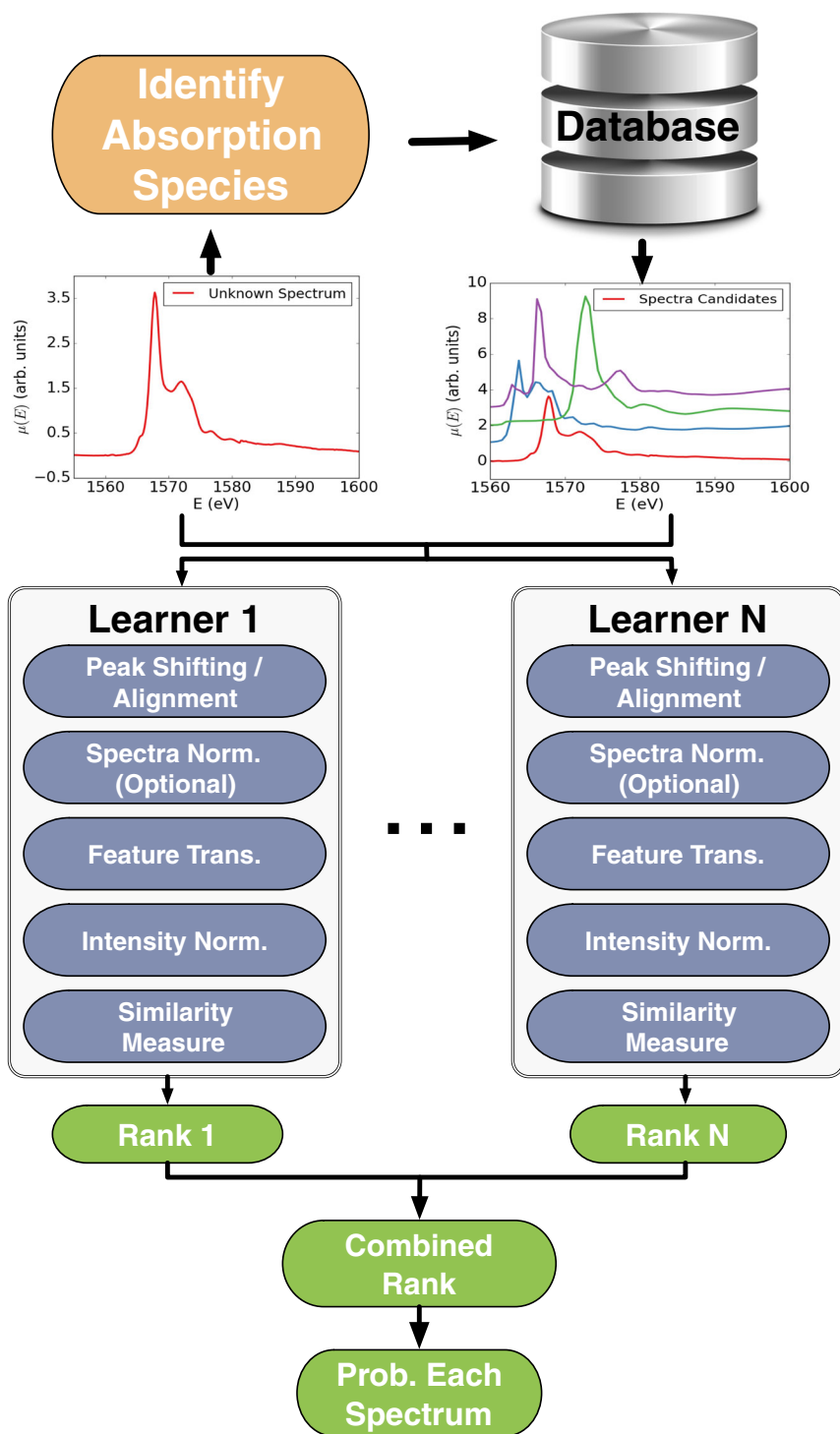


Fig. 3 Workflow schema of the Ensemble-Learned Spectra Identification (ELSIE) algorithm. The algorithm consists of two steps. In the first step, the absorption species is identified and used to narrow down the candidate computed reference spectra. In the second step, the spectral matching ensemble yields a rank-ordered list of computational spectra according to similarity with respect to the target spectrum

Materials Project database with the same chemical composition. It should also be noted that some reference materials may have the same element in multiple oxidation states and coordination environments. Therefore, the application of metrics (ii) and (iii) merely indicates whether at least one of the distinct sites in the top entry have the correct oxidation state and coordination environment. The results are summarized in Table 1.

Of the 19 test spectra, we find that the correct structure is within the top 5 ranked structures for 11 systems, i.e., only 57.9% accuracy. However, the correct oxidation state and coordination environment are in the top entry for 16 and 15 systems, i.e., accuracies of 84.2% and 78.9%, respectively. The best coefficient α is found to be 0.01. Given that XANES is a technique primarily used to extract oxidation state and coordination environment

Table 1. Performance of ELSIE algorithm on 19 test spectra

Formula	Space group	Absorbing species	Correct structure within top 5 rank?	Correct oxidation state in top entries?	Correct coordination environment in top entries?
SiO ₂	<i>P3₂21</i>	Si	No	Yes	Yes
Si	<i>Fd3m</i>	Si	Yes	Yes	Yes
AlPO ₄	<i>I4</i>	Al	No	Yes	Yes
SiC	<i>F43m</i>	Si	No	Yes	Yes
Al ₂ O ₃	<i>R3c</i>	Al	Yes	Yes	Yes
Al	<i>Fm3m</i>	Al	Yes	Yes	Yes
Na ₂ O	<i>Fm3m</i>	Na	Yes	No	No
C	<i>Fd3m</i>	C	No	Yes	No
B ₂ O ₃	<i>P3₂21</i>	B	Yes	No	No
Si ₃ N ₄	<i>P31c</i>	Si	Yes	Yes	Yes
Si ₃ N ₄	<i>P6₃/m</i>	Si	Yes	Yes	Yes
AlN	<i>P6₃mc</i>	Al	Yes	Yes	Yes
NaCl	<i>Fm3m</i>	Na	Yes	Yes	Yes
V ₂ O ₅	<i>Pmmn</i>	V	No	Yes	No
VO ₂	<i>P2₁/c</i>	V	No	Yes	Yes
V ₂ O ₃	<i>R3c</i>	V	No	Yes	Yes
LiNiO ₂	<i>R3m</i>	Ni	No	No	Yes
NiO	<i>Fm3m</i>	Ni	Yes	Yes	Yes
LiCoO ₂	<i>R3m</i>	Co	Yes	Yes	Yes

information, these results are a major validation of the effectiveness of the ELSIE matching algorithm.

To emphasize the effectiveness of the ensemble approach, we also performed the same benchmark using a single learner utilizing just the sigmoid squashing function and cosine similarity measure on spectra that have been pre-normalized with respect to summed intensity. The ELSIE algorithm outperforms the single learner approach by 15.8% in identifying both the correct oxidation state and coordination environment.

We will now illustrate the performance of our spectral matching algorithm with a few case studies on diverse chemistries. For all spectra, we have confined our comparison to the energy range from -10 to 45 eV from the absorption edge, which is the region typically referred to as XANES.

Case study 1: main group metals

Figure 4a, b shows the ELSIE spectral matching results of the Al K-edge XANES of α -Al₂O₃ and Na K-edge XANES of NaCl, respectively. For both target spectra, the correct oxidation states and coordination environments are found in the top candidates. Furthermore, we may observe that our proposed peak shifting approach is effective in aligning the target and reference spectra.

Figure 4c shows a notable case—the Na K-edge of Na₂O—where the ELSIE algorithm fails. Here, the ELSIE algorithm returns elemental Na as the top ranked result, as opposed to Na₂O. The main reason for this failure is that the FEFF-computed spectra is not in good agreement with experimental spectra (see Supplementary Fig. 7 for this and a few other examples). Possible solutions include the use of real-space full-potential multiple-scattering theory or other first principle approaches.⁴⁰ For Na₂O in particular, we find that the experimental Na K-edge XANES of Na₂O is more similar to the computed Na K-edge XANES of Na₂CO₃ (Supplementary Fig. 7(c)), which may indicate possible contamination by the atmosphere in experiments.

Case study 2: transition metal oxides

Figure 5 shows the ELSIE spectra matching results of the Ni K-edge XANES in NiO, Co K-edge XANES in LiCoO₂. From Fig. 5a, we note that although the computed peak positions and amplitude are not in great quantitative agreement with the experimental measured spectra, the ground state NiO entry is nevertheless returned as the top ranked candidate. In particular, the small Ni 1s-3d peak at 8332 eV in the experimental Ni K-edge XANES of NiO is not present in the FEFF calculated spectra. There is, however, a small peak at 8337 eV in the FEFF calculated spectra, which we believe is the Ni 1s-3d peak. The inaccuracy in the position of the peak may be due to the muffin-tin approximation used in FEFF.

For LiCoO₂ (Fig. 5b), the ground state structure of LiCoO₂ (*R3m*) is among the top five entries. All Co³⁺ ions in the top entry (Li(CoO₂)₂) are in octahedral coordination, i.e., the same coordination environment of Co³⁺ ions in LiCoO₂ (*R3m*). We may, therefore, conclude that the ELSIE algorithm performs satisfactorily in both instances.

Figure 5c shows the ELSIE spectra matching results for the V K-edge of V₂O₅ (*Pmmn*). The ELSIE algorithm fails to retrieve the correct square-pyramidal coordination environment of V⁵⁺ in V₂O₅ (*Pmmn*). Indeed, vanadium ions in the top five matches returned by the ELSIE algorithm are in octahedral coordination. Here, the relative similarity of the V K-edge spectra for the different V oxidation states and coordination environments seems to be the key issue. Further structural refinement based on EXAFS simulations, therefore, becomes critical, which will be available in the XASdb in the near future.

In conclusion, we have demonstrated the development of a large database for XAS using high-throughput FEFF calculations. Parameter benchmark results indicate that the overall quality of the FEFF9 calculations with default input parameters is in quantitative agreement with experiments, which is adequate for comparison purposes. We developed a novel spectra-matching algorithm—the ELSIE algorithm—that enables the rapid matching of computed reference spectra to any target spectra. The ensemble learning approach far outperforms any single approach based on a pre-defined set of preprocessing and similarity metric; outstanding ~ 84 and $\sim 79\%$ accuracies in identifying the correct oxidation state and coordination environment are demonstrated based on a diverse test set comprising 19 experimental XANES spectra. The XASdb with the ELSIE algorithm has been integrated into a web application in the Materials Project, providing an important new public resource for the analysis of XAS to all materials researchers, and the ELSIE algorithm itself has been made available as part of *veidt*, an open source machine-learning library for materials science.

METHODS

Benchmarking details

Robust, well-defined datasets are necessary for any benchmarking exercise. We have used the existing high-quality K-edge XAS spectra available in the open EELS Data Base (EELSDb)⁸ as reference data, and matched them with the corresponding materials in the Materials Project¹² using the Materials API⁴¹ and pymatgen.¹³ For materials in the EELSDb without structural information, ground state structures with identical chemical compositions in the Materials Project were used. For spectra in EELSDb taken using the same materials, we selected one and adopted it in our study. Supplementary Table 1 summarizes the 13 unique materials used in this work.

FEFF

The FEFF software calculates X-ray absorption spectra using the RSGF formulation of the multiple-scattering theory.¹¹ The X-ray absorption μ is written in terms of the imaginary part of the one-particle Green's function $G(r, r'; E)$, which incorporate both the inelastic losses and other quasiparticle

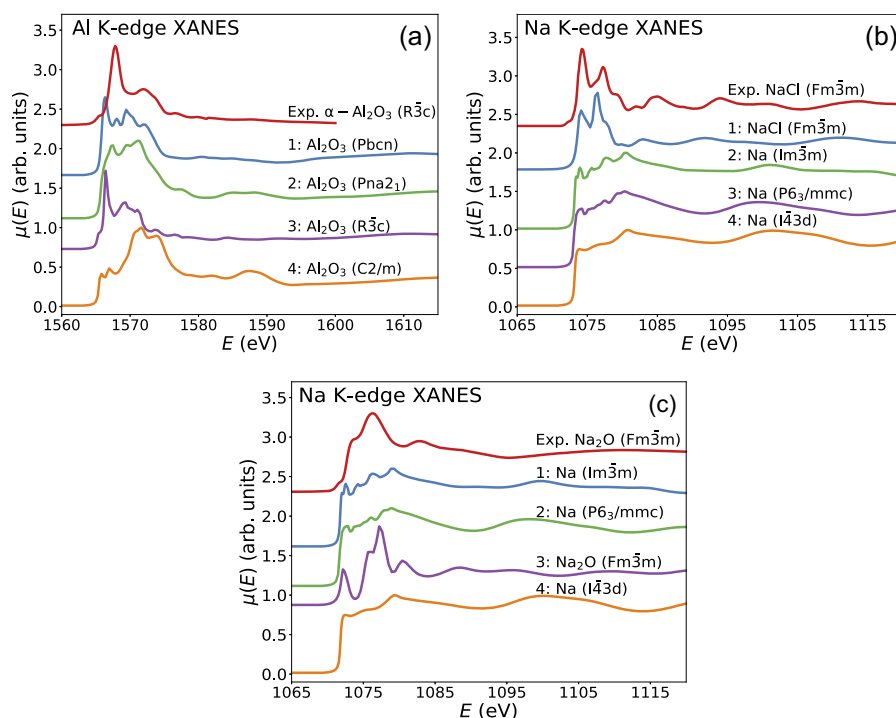


Fig. 4 Results of the similarity ranking returned by the ELSIE matching algorithm on **a** Al K-edge XANES of α - Al_2O_3 entry; **b** Na K-edge XANES of NaCl; and **c** Na K-edge of Na_2O . Detailed information about the retrieved compounds can be found in the Materials Project website, **a** Al_2O_3 (*Pbcn*, mp-1938), Al_2O_3 (*Pna2₁*, mp-2254), Al_2O_3 (*R3c*, mp-1143), and Al_2O_3 (*C2/m*, mp-7048), **b** NaCl (*Fm3m*, mp-22862), Na (*Im3m*, mp-127), Na (*P6₃/mmc*, mp-10172) and Na (*I43d*, mp-567772), and **c** Na (*Im3m*, mp-127), Na (*P6₃/mmc*, mp-10172), Na_2O (*Fm3m*, mp-2352), and Na (*I43d*, mp-567772), in decreasing similarity order

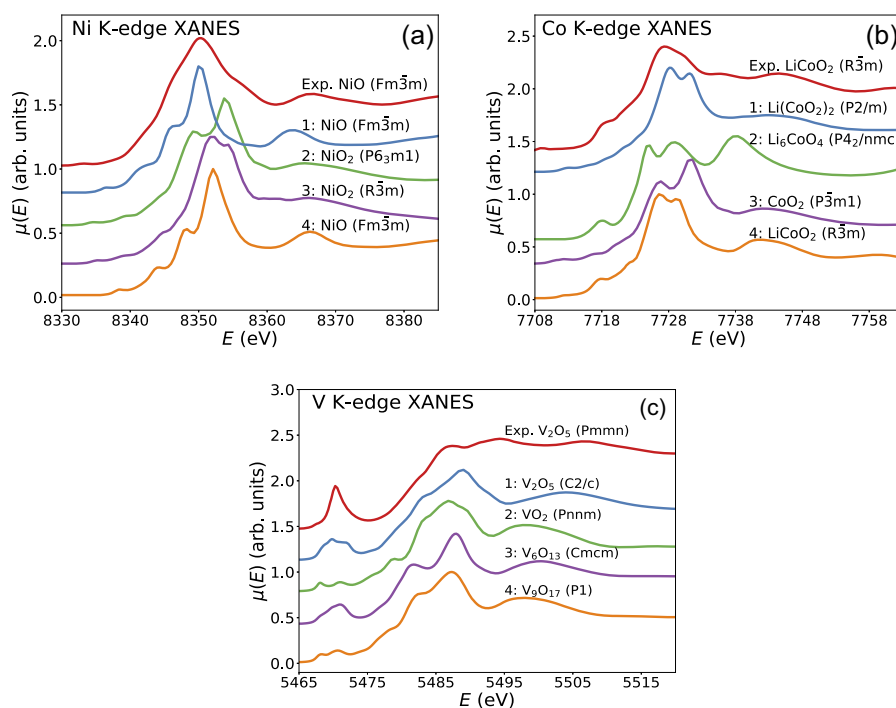


Fig. 5 Results of the similarity ranking returned by the ELSIE matching algorithm on **a** Ni K-edge XANES of NiO; **b** Co K-edge XANES of LiCoO_2 ; and **c** V K-edge of V_2O_5 . Detailed information about the retrieved compounds can be found in the Materials Project website, **a** NiO (*Fm3m*, mp-19009), NiO_2 (*P6₃/m1*, mp-543096), NiO_2 (*R3m*, mp-25593) and NiO (*Fm3m*, mp-715434), **b** $\text{Li}(\text{CoO}_2)_2$ (*P2/m*, mp-553952), Li_6CoO_4 (*P4₂/nmc*, mp-18925), CoO_2 (*P3m1*, mp-714976) and LiCoO_2 (*R3m*, mp-24850), and **c** V_2O_5 (*C2/c*, mp-542844), VO_2 (*Pnnm*, mp-714880), V_6O_{13} (*Cmcm*, mp-715617) and V_9O_{17} (*P1*, mp-716723), in decreasing similarity order

effects. In terms of $G(r, r'; E)$, μ is given by:

$$\mu = -\frac{1}{\pi} \text{Im} \langle c | \hat{\epsilon} \cdot r G(r, r'; E) \hat{\epsilon} \cdot r' | c \rangle \theta_{\Gamma} (E - E_{\Gamma}), \quad (2)$$

where θ_{Γ} is a broadened step function at the Fermi energy E_{Γ} . This yields a unified treatment of EXAFS and XANES. The treatment of X-ray absorption can then be separated into atomic and scattering parts, i.e., $G(r, r'; E) = G^c(r, r'; E) + G^{sc}(r, r'; E)$. The exact result of $G^{sc}(r, r'; E)$ is given by the full matrix inverse, or equivalently, a sum over all multiple-scattering paths.⁴² For the XANES calculation, FEFF implements the FMS technique, which includes the contributions from all orders of scattering within a cluster containing the absorber and scatterers. The FEFF code also incorporates a GW-based self-energy based on the Hedin-Lundqvist plasmon-pole model, which includes effects of electron-electron interactions such as mean-free paths and self-energy shifts. This method has been well tested and is usually a good approximation for EXAFS and reasonable for XANES. FEFF includes a screened corehole and gives results for excitonic enhancements comparable to GW/Bethe-Salpeter equation (BSE) calculations in many materials. FEFF can also incorporate Debye-Waller factors using correlated-Debye or more advanced models. Further details on the FEFF code and its theoretical foundations can be found in ref.¹¹ for interested readers.

In the FEFF input file, parameters are specified in control "cards". The following parameters in FEFF were tested for convergence.

- i. Self-consistent field (SCF): The `rfms1` field in the SCF card specifies the radius of the cluster considered in the FMS calculation. The higher the `rfms1` is, the greater the number of atoms is included in calculation.
- ii. Full multiple scattering (FMS): The `rfms` field in the FMS card determines the total number of multiple-scattering paths considered in the XANES calculation. Default values are used for the other five optional fields in the FMS card.
- iii. EXCHANGE: The EXCHANGE card specifies the exchange-correlation potential model used for XANES calculation. No shift was applied to the Fermi energy level in this work, i.e., the second and third fields of the EXCHANGE card were kept being 0.
- iv. COREHOLE: The COREHOLE card is used to specify the treatment of the core during XAS calculations. "Core hole" is the hole in the orbital formed by the excitation of a single electron from that orbital.⁵ In FEFF9 code, a combination of BSE and time-dependent density functional theory (TDDFT) is used to improve the approximation of the core hole interactions.^{10,20}

ELSIE algorithm construction

We adopted the concept of ensemble method to index the most similar spectra from the database with respect to a target spectrum. Each weak learner has a unique combination of a few spectral preprocessing techniques and one similarity metric, we will describe the preprocessing approaches and similarity metrics in turn.

Each preprocessor comprises a series of steps, designed to emphasize or weaken certain characteristics of the experimental and computed spectra. A preprocessor is generated as follows:

1. *Peak shifting and quantization*: This step is necessary to all preprocessors. Because of the differences in energy sampling intervals and energy ranges, linear interpolation was used to convert each spectrum to a vector of 200 intensity values with identical energy grid. The reference spectra are shifted such that the onset of absorption, which is well-defined by the photoelectric effect, is aligned with that of the target spectra. This onset is determined by ascertaining the lowest incident energy at which the computed absorption intensity reaches 6% of the peak intensity.
2. *Pre-normalization*: We included an optional pre-normalization step to rescale the intensity to a similar range. Given the spectrum X with X_i represents the i th intensity, four normalization approaches are adopted:⁴³

$$X_i^{\text{norm}} = \frac{X_i}{\sum X_i}. \quad (3)$$

$$X_i^{\text{norm}} = \frac{X_i}{\sqrt{\sum X_i^2}}. \quad (4)$$

$$X_i^{\text{norm}} = \frac{X_i - X_{\min}}{X_{\max} - X_{\min}}. \quad (5)$$

$$X_i^{\text{norm}} = (X_i - \mu) / \sigma, \text{ where } \mu = \sum X_i / n \text{ and}$$

$$\sigma = \sqrt{\sum (X_i - \mu)^2 / n}. \quad (6)$$

3. *Feature transformation*: Several feature transformation functions were implemented in the third step, which include the square root and sigmoid squashing functions. The sigmoid squashed spectrum is calculated using $X' = \frac{1 - \cos(\pi X)}{2}$. The squared root squashing uses $X' = \sqrt{X}$, where X is the squashed new spectrum. This technique has shown to improve the response sensitivity with respect to different spectral features.⁴⁴ The feature transformation functions also include taking the first or second order derivative of spectrum, or weighted the spectra with the first and second order derivatives. This step is necessary to make distinct weak learners.
4. *Normalization*: This last step is for all preprocessors. The spectra are all normalized such that the sum of intensities is equal to 1, i.e. $\sum_{i=1}^D X_i = 1$.

Both the computed and target spectra are processed using the same series of steps for each preprocessor.

The preprocessed target and computed spectra are then compared in a pairwise manner using a similarity metric. Only bin-to-bin similarity metrics are used in the ELSIE algorithm development as they are less computationally demanding for high-throughput datasets.⁴⁵ Four commonly used similarity metrics in the literatures are used in the ELSIE algorithm:

1. *Pearson correlation*: as defined in the Benchmarking section.
2. *Euclidean similarity*: In the D-dimensional spectral feature space, the Euclidean distance between two spectra X and Y is given by the following equation:

$$d_{\text{Euc}} = \sqrt{\sum_{i=1}^D |X_i - Y_i|^2}. \quad (7)$$

The spectral similarity measure can be derived from the distance calculated using the following expression:

$$S_{\text{Euc}}(X, Y) = 1 - \frac{d_{\text{Euc}}(X, Y)}{d_{\text{Euc}}^{\max}}, \quad (8)$$

where d_{Euc}^{\max} is the absolute maximum expected Euclidean distance between two probability mass functions.⁴⁵

3. *Cosine similarity*: The cosine similarity measure is the normalized inner product and measures the angle between two spectral vectors.⁴⁶ The cosine similarity between two spectra can be calculated as:

$$S_{\text{Cos}} = \frac{\sum_{i=1}^D X_i Y_i}{\sqrt{\sum_{i=1}^D X_i^2} \sqrt{\sum_{i=1}^D Y_i^2}}. \quad (9)$$

4. *Ruzicka similarity*: The Ruzicka⁴⁵ similarity between two spectra is given by the following equation:

$$S_{\text{Ruz}} = \frac{\sum_{i=1}^D \min(X_i, Y_i)}{\sum_{i=1}^D \max(X_i, Y_i)}. \quad (10)$$

The combination of preprocessors and similarity metrics results in a total of 168 learners that can potentially be used to construct the ELSIE algorithm. To make an ensemble that outperforms individual learners, one prerequisite is that each learner should have an error rate lower than random guessing. We, therefore, filtered the 168 learners to 33 and adopted them in the ELSIE algorithm. The detailed filtering procedure can be found in the Supplementary Information.

For each target spectrum, each learner (one preprocessor + one similarity metric) outputs similarity scores for the reference spectra. However, the quantitative scores for different similarity metrics cannot be compared even for the same target spectrum. In the ELSIE algorithm, we instead combine the reference spectra ranking from each learner to derive an ensemble result. For a mixture of classifiers of various types, ranking-based combination methods have been shown to be more reliable.⁴⁷ Based on the rankings, we compute the Borda count, defined as the number of candidates that are ranked equal and below the specific candidate. For example, the top spectrum among ten computed candidates would receive a Borda count of 10, while the second ranked spectrum has a Borda count of 9. For each target spectrum, the Borda

counts of the reference spectra under all learners are then summed to arrive at a consensus ranking.⁴⁸

Finally, the Borda ranks of all reference spectra are then combined with a penalty term for the peak shift and converted to a probabilistic estimate using the modified softmax function. The probability of a reference spectrum X^k is indicated by $P(X^k)$ where the superscript k indicates the k -th spectrum, and is calculated as follows:

1. The Borda count of each reference (R^k) is normalized with respect to the count sum: $R_{\text{norm}}^k = \frac{R^k}{\sum R^k}$. This step is required to avoid the exponential overflow.
2. $P(X^k)$ is then calculated by the following equation:

$$P(X^k) = \frac{\exp(R_{\text{norm}}^k) \exp\left(-\frac{a|\Delta S^k|}{\delta_s}\right)}{\sum \exp(R_{\text{norm}}^k) \exp\left(-\frac{a|\Delta S^k|}{\delta_s}\right)}, \quad (11)$$

where ΔS^k could be calculated as $\Delta S^k = S^k - \bar{S}$. S^k is the peak shift amount between the reference spectrum X^k and the target spectrum. \bar{S} is the mean peak shift of the reference spectra. δ_s is the standard deviation of S^k .

Coefficient a is fitted to the test dataset. $\exp\left(-\frac{a|\Delta S^k|}{\delta_s}\right)$ is therefore a term that imposes a larger penalty on large peak shifts relative to smaller peak shifts.

The algorithm itself has been highly optimized by leveraging on well-established numerical packages such as numpy and scipy.^{49,50} On a laptop computer with Intel i5 2.6GHz single CPU and 2 GB of RAM, the ELSIE algorithm can perform a comparison between a target and candidate spectrum in about 0.03 s. Typically, 20–30 spectra are selected for comparison according to the rules that the computational reference spectra should have identical absorption species, limited number of elements and $E_{\text{hull}} < 100$ meV/atom. The overall time to perform a complete ranking is, therefore, around 1 s, which allows for on-the-fly matching of uploaded spectra.

Data availability

The computed spectra in the XASdb have been made available in the Materials Project website. A new web application—the XASApp (<https://materialsproject.org/#apps/xas/>)—has been developed which allows any user to compare multiple X-ray absorption spectra and find matches within the XASdb for an uploaded spectrum using the ELSIE algorithm.

The ELSIE algorithm has also been made publicly available as a part of *veidt*, an open-source Python machine-learning library for materials science developed by the Materials Virtual Lab that is available on the Python Package Index and Github (<https://github.com/materialsvirtuallab/veidt>).

ACKNOWLEDGEMENTS

This work is supported by the National Science Foundation's Cyberinfrastructure Framework for 21st Century Science and Engineering (CIF21) program under Award No. 1640899. The Materials Project, supported by the Department of Energy (DOE) Basic Energy Sciences (BES) program, under Grant No. EDCBEE is gratefully acknowledged for web dissemination and data infrastructure. The FEFF project is supported primarily by DOE BES Grant DE-FG02-97ER45623. We also acknowledge computational resources provided by Triton Shared Computing Cluster (TSCC) at the University of California, San Diego, the National Energy Research Scientific Computing Center (NERSC), and the Extreme Science and Engineering Discovery Environment (XSEDE) supported by National Science Foundation under grant number ACI-1053575. L.F.J.P. acknowledges support from the National Science Foundation (DMREF-1627583).

AUTHOR CONTRIBUTIONS

C.Z., K.M., and C.C. performed the workflow design, code implementation, and calculation analysis. Y.C., H.T., and A.D., J.J.K., F.D.V., and J.J.R. helped to the simulations of XAS spectra. L.F.J.P. helped experimental XANES spectra analysis. K.P. and S.P.O. is the primary investigators and supervised the workflow and code development. All authors contributed to the writing and editing of the manuscript.

ADDITIONAL INFORMATION

Supplementary information accompanies the paper on the *npj Computational Materials* website (<https://doi.org/10.1038/s41524-018-0067-x>).

Competing interests: The authors declare no competing interests

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

REFERENCES

1. Lin, Y.-C. et al. Thermodynamics, kinetics and structural evolution of ϵ -LiVOPO₄ over multiple lithium intercalation. *Chem. Mater.* **28**, 1794–1805 (2016).
2. Yu, X. et al. High rate delithiation behaviour of LiFePO₄ studied by quick X-ray absorption spectroscopy. *Chem. Commun.* **48**, 11537–11539 (2012).
3. Cheng, J.-H. et al. Simultaneous Reduction of Co 3+ and Mn 4+ in P2-Na 2/3 Co 2/3 Mn 1/3 O 2 as evidenced by x-ray absorption spectroscopy during electrochemical sodium intercalation. *Chem. Mater.* **26**, 1219–1225 (2014).
4. Koningsberger, D. C. & Prins, R. *X-ray Absorption: Principles, Applications, Techniques of EXAFS, SEXAFS, and XANES* (Wiley, New York, 1988).
5. Bunker, G. *Introduction to XAFS: A Practical Guide to X-ray Absorption Fine Structure Spectroscopy* (Cambridge University Press, New York, 2010).
6. Newville, M. Fundamentals of XAFS. *Rev. Mineral. Geochem.* **78**, 33–74 (2014).
7. Ravel, B. & Newville, M. ATHENA, ARTEMIS, HEPHAESTUS: Data analysis for X-ray absorption spectroscopy using IFEFFIT. *J. Synchrotron Radiat.* **12**, 537–541 (2005).
8. Ewels, P., Sikora, T., Serin, V., Ewels, C. P. & Lajaunie, L. A complete overhaul of the electron energy-loss spectroscopy and X-ray absorption spectroscopy database: eelsdb.eu. *Microsc. Microanal.* **22**, 717–724 (2016).
9. Egerton, R. F. *Electron Energy-Loss Spectroscopy in the Electron Microscope*. (Springer, Boston, MA, 2011).
10. Rehr, J. J., Kas, J. J., Vila, F. D., Prange, M. P. & Jorissen, K. Parameter-free calculations of X-ray spectra with FEFF9. *Phys. Chem. Chem. Phys.* **12**, 5503 (2010).
11. Rehr, J. J. Theoretical approaches to x-ray absorption fine structure. *Rev. Mod. Phys.* **72**, 621–654 (2000).
12. Jain, A. et al. Commentary: The Materials Project: A materials genome approach to accelerating materials innovation. *APL Mater.* **1**, 11002 (2013).
13. 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).
14. Jain, A. et al. FireWorks: a dynamic workflow system designed for high-throughput applications. *Concurr. Comput. Pract. Exp.* **27**, 5037–5059 (2015).
15. 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).
16. Wang, Z. et al. Effects of cathode electrolyte interfacial (CEI) layer on long term cycling of all-solid-state thin-film batteries. *J. Power Sources* **324**, 342–348 (2016).
17. Jia, Q. et al. Experimental observation of redox-induced Fe–N switching behavior as a determinant role for oxygen reduction activity. *ACS Nano* **9**, 12496–12505 (2015).
18. Beharid, F. et al. Structural and electronic properties of micellar Au nanoparticles: Size and ligand effects. *ACS Nano* **8**, 6671–6681 (2014).
19. Jorissen, K. & Rehr, J. J. Calculations of electron energy loss and x-ray absorption spectra in periodic systems without a supercell. *Phys. Rev. B* **81**, 245124 (2010).
20. Vinson, J. & Rehr, J. J. Ab initio Bethe-Salpeter calculations of the x-ray absorption spectra of transition metals at the L-shell edges. *Phys. Rev. B-Condens. Matter Mater. Phys.* **86**, 1–6 (2012).
21. Perdew, J. P., Burke, K. & Ernzerhof, M. Generalized gradient approximation made simple. *Phys. Rev. Lett.* **77**, 3865–3868 (1996).
22. Wu, Z. & Cohen, R. E. More accurate generalized gradient approximation for solids. *Phys. Rev. B-Condens. Matter Mater. Phys.* **73**, 2–7 (2006).
23. Kresse, G. & Harl, J. Accurate bulk properties from approximate many-body techniques. *Phys. Rev. Lett.* **103**, 4–7 (2009).
24. Haas, P., Tran, F. & Blaha, P. Calculation of the lattice constant of solids with semilocal functionals. *Phys. Rev. B-Condens. Matter Mater. Phys.* **79**, 1–10 (2009).
25. Klimeš, J., Bowler, D. R. & Michaelides, A. Van der Waals density functionals applied to solids. *Phys. Rev. B* **83**, 195131 (2011).
26. Alkauskas, A. & Pasquarello, A. Band-edge problem in the theoretical determination of defect energy levels: The O vacancy in ZnO as a benchmark case. *Phys. Rev. B* **84**, 125206 (2011).
27. Perdew, J. et al. Restoring the density-gradient expansion for exchange in solids and surfaces. *Phys. Rev. Lett.* **100**, 136406 (2008).
28. Heyd, J., Scuseria, G. E. & Ernzerhof, M. Hybrid functionals based on a screened Coulomb potential. *J. Chem. Phys.* **118**, 8207–8215 (2003).
29. Paier, J., Asahi, R., Nagoya, A. & Kresse, G. Cu₂ZnSnS₄ as a potential photovoltaic material: A hybrid Hartree-Fock density functional theory study. *Phys. Rev. B-Condens. Matter Mater. Phys.* **79**, 1–8 (2009).
30. Da Silva, J. L. F., Ganduglia-Pirovano, M. V., Sauer, J., Bayer, V. & Kresse, G. Hybrid functionals applied to rare-earth oxides: The example of ceria. *Phys. Rev. B-Condens. Matter Mater. Phys.* **75**, 19–24 (2007).
31. Wróbel, J., Kurzydowski, K. J., Hummer, K., Kresse, G. & Piechota, J. Calculations of ZnO properties using the Heyd-Scuseria-Ernzerhof screened hybrid density functional. *Phys. Rev. B-Condens. Matter Mater. Phys.* **80**, 1–8 (2009).

32. Ong, S. P., Mo, Y. & Ceder, G. Low hole polaron migration barrier in lithium peroxide. *Phys. Rev. B-Condens. Matter Mater. Phys.* **85**, 2–5 (2012).
33. Carey, C., Boucher, T., Mahadevan, S., Bartholomew, P. & Dyar, M. D. Machine learning tools for mineral recognition and classification from Raman spectroscopy. *J. Raman Spectrosc.* **46**, 894–903 (2015).
34. Liu, J. et al. Methods for peptide identification by spectral comparison. *Proteome Sci.* **5**, 3 (2007).
35. Stein, S. E. & Scott, D. R. Optimization and testing of mass spectral library search algorithms for compound identification. *J. Am. Soc. Mass Spectrom.* **5**, 859–866 (1994).
36. Rana, J. et al. Local structural changes in LiMn_{1.5}Ni_{0.5}O₄ spinel cathode material for lithium-ion batteries. *J. Power Sources* **255**, 439–449 (2014).
37. Rana, J. et al. On the structural integrity and electrochemical activity of a 0.5Li₂MnO₃-0.5LiCoO₂ cathode material for lithium-ion batteries. *J. Mater. Chem. A* **2**, 9099 (2014).
38. Bearden, J. A. & Burr, A. F. Reevaluation of X-ray atomic energy levels. *Rev. Mod. Phys.* **39**, 125–142 (1967).
39. Sun, W. et al. The thermodynamic scale of inorganic crystalline metastability. *Sci. Adv.* **2**, e1600225–e1600225 (2016).
40. Xu, J. et al. X-ray absorption spectra of graphene and graphene oxide by full-potential multiple scattering calculations with self-consistent charge density. *Phys. Rev. B* **92**, 125408 (2015).
41. Ong, S. P. et al. The materials application programming interface (API): A simple, flexible and efficient API for materials data based on REpresentational State Transfer (REST) principles. *Comput. Mater. Sci.* **97**, 209–215 (2015).
42. Ravel, B. A practical introduction to multiple scattering theory. *J. Alloy. Compd.* **401**, 118–126 (2005).
43. Zoubir, A. *Raman Imaging*, Vol. 168 (Springer, Berlin Heidelberg, 2012).
44. Hansen, M. E. & Smedsgaard, J. A new matching algorithm for high resolution mass spectra. *J. Am. Soc. Mass Spectrom.* **15**, 1173–1180 (2004).
45. Hernández-Rivera, E., Coleman, S. P. & Tschopp, M. A. Using Similarity Metrics to Quantify Differences in High-Throughput Data Sets: Application to X-ray Diffraction Patterns. *ACS Comb. Sci.* **19**, 25–36 (2017).
46. Deza, M. M. & Deza, E. *Encyclopedia of Distances*. (Springer, Berlin Heidelberg, 2013).
47. Ho, T. K., Hull, J. J. & Srihari, S. N. Decision combination in multiple classifier systems. *IEEE Trans. Pattern Anal. Mach. Intell.* **16**, 66–75 (1994).
48. Black, D. *The Theory of Committees and Elections*. (Springer, Netherlands, 1986).
49. Jones, E., Oliphant, T. & Peterson, P. *Scipy: Open Source Scientific Tools For Python*. <http://www.scipy.org> (2001).
50. van der Walt, S., Colbert, S. C. & Varoquaux, G. The NumPy Array: A structure for efficient numerical computation. *Comput. Sci. Eng.* **13**, 22–30 (2011).



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2018