# JARVIS: An Integrated Infrastructure for Data-driven Materials Design

Kamal Choudhary<sup>1,2</sup>, Kevin F. Garrity<sup>1</sup>, Andrew C. E. Reid<sup>1</sup>, Brian DeCost<sup>1</sup>, Adam J. Biacchi<sup>3</sup>, Angela R. Hight Walker<sup>3</sup>, Zachary Trautt<sup>1</sup>, Jason Hattrick-Simpers<sup>1</sup>, A. Gilad Kusne<sup>1</sup>, Andrea Centrone<sup>3</sup>, Albert Davydov<sup>1</sup>, Jie Jiang<sup>4</sup>, Ruth Pachter<sup>4</sup>, Gowoon Cheon<sup>5</sup>, Evan Reed<sup>5</sup>, Ankit Agrawal<sup>6</sup>, Xiaofeng Qian<sup>7</sup>, Vinit Sharma<sup>8,9</sup>, Houlong Zhuang<sup>10</sup>, Sergei V. Kalinin<sup>11</sup>, Bobby G. Sumpter<sup>11</sup>, Ghanshyam Pilania<sup>12</sup>, Pinar Acar<sup>13</sup>, Subhasish Mandal<sup>14</sup>, Kristjan Haule<sup>14</sup>, David Vanderbilt<sup>14</sup>, Karin Rabe<sup>14</sup>, Francesca Tavazza<sup>1</sup>

- 1. Materials Measurement Laboratory, National Institute of Standards and Technology, Gaithersburg, MD, 20899, U.S.A.
- 2. Theiss Research, La Jolla, CA, 92037, U.S.A.
- 3. Physical Measurement Laboratory, National Institute of Standards and Technology, Gaithersburg, MD, 20899, U.S.A.
- 4. Materials and Manufacturing Directorate, Air Force Research Laboratory, Wright-Patterson Air Force Base, OH 45433, USA.
- 5. Department of Materials Science and Engineering, Stanford University, Stanford, CA, 94305, U.S.A.
- 6. Department of Electrical and Computer Engineering, Northwestern University, Evanston, IL 60208, U.S.A.
- 7. Department of Materials Science and Engineering, Texas A&M University, TX, 77843, U.S.A.
- 8. Joint Institute for Computational Sciences, University of Tennessee, Knoxville, TN, 37996, USA.
- 9. National Institute for Computational Sciences, Oak Ridge National Laboratory, TN 37831, U.S.A.
- 10. School for Engineering of Matter, Transport & Energy, Arizona state university, AZ, 85287, U.S.A.
- 11. Center for Nanophase Materials Sciences, Oak Ridge National Laboratory, TN 37831, U.S.A.
- 12. Materials Science and Technology Division, Los Alamos National Lab, Los Alamos, NM, 87545, U.S.A.
- 13. Department of Mechanical Engineering, Virginia Tech, Blacksburg, VA 24061, U.S.A.
- 14. Department of Physics and Astronomy, Rutgers University, Piscataway, NJ, 08901, USA

**Abstract** 

The Joint Automated Repository for Various Integrated Simulations (JARVIS) is an integrated

infrastructure to accelerate materials discovery and design using density functional theory (DFT),

classical force-fields (FF), and machine learning (ML) techniques. JARVIS is motivated by the

Materials Genome Initiative (MGI) principles of developing open-access databases and tools to

reduce the cost and development time of materials discovery, optimization, and deployment. The

major features of JARVIS are: JARVIS-DFT, JARVIS-FF, JARVIS-ML, and JARVIS-Tools. To

date, JARVIS consists of  $\approx 40,000$  materials and  $\approx 1$  million calculated properties in JARVIS-

DFT,  $\approx$  1,500 materials and  $\approx$  110 force-fields in JARVIS-FF, and  $\approx$  25 ML models for material-

property predictions in JARVIS-ML, all of which are continuously expanding. JARVIS-Tools

provides scripts and workflows for running and analyzing various simulations. We compare our

computational data to experiments or high-fidelity computational methods wherever applicable to

evaluate error/uncertainty in predictions. In addition to the existing workflows, the infrastructure

can support a wide variety of other technologically important applications as part of the data-driven

materials design paradigm. The databases and tools are publicly distributed through the following

major websites http://jarvis.nist.gov/ and https://github.com/usnistgov/jarvis.

**Corresponding author:** Kamal Choudhary ( kamal.choudhary@nist.gov )

2

## Introduction

The Materials Genome Initiative (MGI)<sup>1,2</sup> was introduced in 2011 to accelerate materials discovery using computational<sup>3-9</sup>, experimental<sup>10-13</sup> and data analytics<sup>14-16</sup> approaches. The MGI has revolutionized several fields for materials-applications, such as batteries<sup>17</sup>, thermoelectrics<sup>18</sup>, and alloy-design<sup>19</sup>, thorough open-access public database and tool development<sup>20</sup>. The MGI encourages systematic Process-Structure-Property-Performance (PSPP)<sup>21</sup>-based efficient designapproaches rather than Edisonian trial-error methods<sup>22</sup>.

Especially in the field of computational materials design, quantum mechanics-based density functional theory (DFT)<sup>23</sup> has proven to be an immensely successful technique, and several databases of automated DFT calculations are widely used in materials design applications. Despite their successes, existing DFT databases face limitations due to issues intrinsic to conventional DFT approaches, e.g. the generalized gradient approximation of Perdew-Burke-Ernzerhof (GGA-PBE)<sup>23,24</sup>. Drawbacks of the existing DFT databases include non-inclusion of van der Waals (vdW) interactions<sup>8</sup>, bandgap underestimations<sup>25</sup>, non-inclusion of spin-orbit coupling<sup>7</sup>, overly simplifying magnetic ordering<sup>26</sup>, neglecting defects<sup>27</sup> (point, line, surface and volume), unconverged computational parameters such as k-points<sup>28</sup>, ignoring temperature effects<sup>29</sup> (generally DFT calculations are performed at 0 K), lack of layer/thickness-dependent properties of low dimensional materials<sup>30</sup>, and lacking interfaces/heterostructures of materials<sup>31</sup>, all of which can be critical for realistic material-applications. Additionally, there are several other computational approaches, such as classical force-field (FF)<sup>32</sup>, computational microscopy, phasefield (PF), CALculation of PHAse Diagrams (CALPHAD)<sup>33</sup>, and Orientation Distribution Functions (ODF)<sup>34</sup> which lack the integrated tools and databases that have been developed for DFT-based computational approaches. Finally, the integration of computational approaches with experiments, the application of statistical uncertainty analysis, and the implementation of data analytics and artificial intelligence (AI) techniques require significant developments to meet the goals set forth by the MGI.

Some of the notable materials databases are: Automatic-FLOW for Materials Discovery (AFLOW)<sup>3</sup>, Materials-project<sup>4</sup>, Khazana<sup>17</sup>, Open Quantum Materials Database (OQMD)<sup>5</sup>, Novel Materials Discovery (NOMAD)<sup>9</sup>, Computational Materials Repository (CMR)<sup>35</sup>, NIMS-MatNavi<sup>36</sup>, NREL-MatDB<sup>37</sup>, Inorganic Crystal Structure Database (ICSD)<sup>38</sup>, Materials-Cloud<sup>39</sup>, Citrine<sup>40</sup>, OpenKIM<sup>41</sup>, Predictive Integrated Structural Materials Science (PRISMS)<sup>42</sup>, and Phase-Field hub (PFhub)<sup>43</sup>. Some of the commonly used computational-tools are Python Materials Genomics (PYMATGEN)<sup>44</sup>, Atomic Simulation Environment (ASE)<sup>45</sup>, Automated Interactive Infrastructure and Database (AIIDA)<sup>6</sup> and MPinterfaces<sup>46</sup>. The data most commonly included in these databases consists of crystal structures, formation energies, bandgaps, elastic constants, Poisson ratios, piezoelectric constants, and dielectric constants. These material properties can be used directly to screen for potentially interesting materials for a given application as candidates for experimental synthesis and characterization, as well as part of a PSPP design approach to better understand the factors driving material performance. Beyond the directly calculated material properties mentioned above, several new selection metrics are also being developed to aid materials design, such as scintillation attenuation length<sup>47</sup>, thermoelectric complexity factor<sup>48</sup>, spectroscopy limited maximum efficiency<sup>49,50</sup>, exfoliation energy<sup>8</sup>, and spin-orbit spillage<sup>7,26,51</sup>. Akin to DFT-like standard computational approaches that are used as screening tools for experiments, machine learning (ML) <sup>14-16,52</sup> models for materials design are being developed as pre-screening tools for other conventional computational methods such as DFT. In addition, ML

tools are proposed to accelerate experimental methods directly based on computational data<sup>53</sup>. All of the above developments show immense promise for accelerating materials design.

The principles mentioned above constitute the foundations of the Joint Automated Repository for Various Integrated Simulations (JARVIS) (<a href="https://jarvis.nist.gov">https://jarvis.nist.gov</a>) infrastructure, a set of databases and tools to meet some of the current material-design challenges. The main components of JARVIS are: JARVIS-DFT, JARVIS-FF, JARVIS-ML, and JARVIS-Tools. JARVIS is developed and hosted at the National Institute of Standards and Technology (NIST)<sup>54</sup> as part of the MGI.

Started in 2017, JARVIS-DFT<sup>7,8,25-27,30,31,49,53,55</sup> is a repository based on DFT calculations that mainly uses the vdW-DF-OptB88 van der Waals functional<sup>56</sup>. The database also uses beyond-GGA approaches for a subset of materials, including the Tran-Blaha modified Becke-Johnson (TBmBJ) meta-GGA<sup>57</sup>, the hybrid functional PBE0, the hybrid range-separated functional Heyd-Scuseria-Ernzerhof (HSE06), Dynamical Mean Field Theory (DMFT), and G<sub>0</sub>W<sub>0</sub>. In addition to hosting conventional properties such as formation energies, bandgaps, elastic constants, piezoelectric constants, dielectric constants, and magnetic moments, it also contains unique datasets, such as exfoliation energies for van der Waals bonded materials, the spin-orbit coupling (SOC) spillage, improved meta-GGA bandgaps, frequency-dependent dielectric functions, the spectroscopy limited maximum efficiency (SLME), infrared (IR) intensities, electric field gradients (EFG), heterojunction classifications, and Wannier tight-binding Hamiltonians. These datasets are compared to experimental results wherever possible to evaluate their accuracy as predictive tools. JARVIS-DFT also introduced protocols such as automatic k-point convergence, which can be critical for obtaining precise and accurate results. JARVIS-DFT is distributed through the website: https://www.ctcms.nist.gov/~knc6/JVASP.html.

The JARVIS-FF<sup>27,58</sup> database, also started in 2017, is a repository of classical force-field/potential computational data intended to help a user select the most appropriate force-field for a specific application. Many classical force-fields are developed for a particular set of properties (such as energies), and may not have been tested for properties not included in training (such as elastic constants, or defect formation energies). JARVIS-FF provides an automatic framework to consistently calculate and compare basic properties, such as the bulk modulus, defect formation energies, phonons, etc., that may be critical for specific molecular-dynamics simulations. JARVIS-FF relies on DFT and experimental data to evaluate accuracy. JARVIS-FF is distributed through the website: <a href="https://www.ctcms.nist.gov/~knc6/periodic.html">https://www.ctcms.nist.gov/~knc6/periodic.html</a>.

The JARVIS-ML<sup>49,53,55,59,60</sup> is a repository of machine learning (ML) model parameters, descriptors, and ML-related input and target data. JARVIS-ML introduced Classical Force-field Inspired Descriptors (CFID) in 2018 as a universal framework to represent a material's chemistry-structure-charge related data. With the help of CFID and JARVIS-DFT data, several high-accuracy classification and regression ML models were developed, with applications to fast materials-screening and energy-landscape mapping. Some of the trained property models include formation energies, exfoliation energies, bandgaps, magnetic moments, refractive indexes, dielectric constants, thermoelectric performance, and maximum piezoelectric and infrared modes. Also, several ML interpretability analyses have provided physical-insights beyond intuitive materials-science knowledge<sup>59</sup>. These models, the workflow, the datasets, etc. are disseminated to enhance the transparency of the work. Recently, JARVIS-ML was expanded to include ML models to analyze STM-images in order to directly accelerate the interpretation of experimental images. Graph convolution neural network models are currently being developed for automated handling

of images and crystal-structure analysis in materials science. JARVIS-ML is distributed through the website: https://www.ctcms.nist.gov/jarvisml.

JARVIS-Tools is the underlying computational framework used for automation, data-generation, data-handling, analysis and dissemination of all the above repositories. JARVIS-Tools uses cloud-based continuous integration, low-software dependency, auto-documentation, Jupyter and Google-Colab notebook integration, pip installation and related strategies to make the software robust and easy to use. JARVIS-Tools also hosts several examples to enable a user to reproduce the data in the above repositories or to apply the tools for their own applications. JARVIS-Tools are provided through the GitHub page: <a href="https://github.com/usnistgov/jarvis">https://github.com/usnistgov/jarvis</a>.

This paper is organized as follows: 1) we introduce the main computational techniques, organized by the time and length scales, 2) we illustrate JARVIS-Tools and its functionalities, 3) we discuss the contents of the major JARVIS databases, 4) we demonstrate some of the derived applications, and 5) we discuss outstanding challenges and future work.

#### **Results and discussion**

# Overview of computational techniques

There are many computational tools for simulating realistic materials depending on the time and length scales of interest<sup>61</sup>. Before we discuss the details of JARVIS, we will provide a brief list of these techniques and highlight their range of applicability, as summarized in Fig. 1. Relevant techniques include quantum mechanical computations, classical/molecular mechanics, mesoscale modeling, finite element analysis, and engineering design. Each of these methodologies has its own ontology and semantics for describing themselves and the PSPP relationship. For example, 'structure' may imply electronic configurations in the quantum regime, atomic arrangement in

molecular mechanics, microstructure, segments in phase field-based mesoscale modeling, and mesh-structure in finite element analysis. Material properties are calculated using corresponding physical laws such as the Schrödinger equation in the quantum regime, or Newton's laws of motion for classical regimes. For realistic material design, it is important to integrate these methods. A major challenge for multiscale modeling is propagating the results of one simulation into another while capturing the relevant physics. Artificial Intelligence (AI) techniques have been applied in each of these domains and can be used to integrate the methods to a certain extent <sup>14</sup>. In JARVIS, we primarily focus on atomistic-based classical and quantum simulations and machine-learning, but we also attempt to integrate other simulation methods with our atomistic data for a few specific applications.

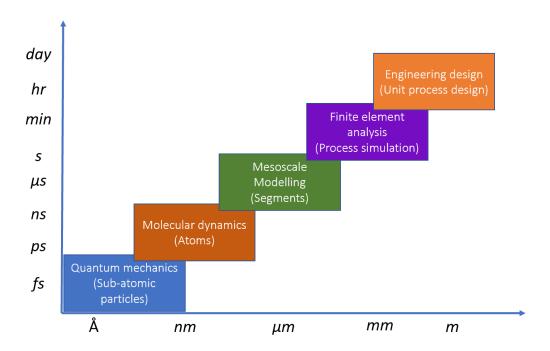


Fig. 1 Length and time-scale based computational materials design techniques.

#### Software and databases

The JARVIS infrastructure (Fig. 2) is a combination of databases and tools for running and integrating some of the computational methods mentioned above. The general procedure for adding a dataset to JARVIS is as follows. We start with the goal of finding or designing a material to display or optimize a given property. Then, we decide on an appropriate computational method, as well as a computationally efficient way to screen for the best candidate materials. The screening process can proceed in several steps, with computationally inexpensive methods applied first, followed by more computationally intensive methods on the remaining materials. Whenever possible, the data is compared with available experiments to evaluate the accuracy and quality of the database. Once a large enough dataset is generated, machine learning techniques can be utilized to accelerate the traditional computational approaches.

As an example, we consider the goal of finding materials to maximize solar-cell efficiency, where the appropriate computational tool is DFT. We develop a screening criterion (Spectroscopic Limited Maximum Efficiency, SLME) and calculate the necessary properties (dielectric function and band gap). We test the method by comparing known materials to experiment, and we perform more accurate meta-GGA and GW calculations as additional screening and validation steps. Finally, we develop a machine learning model to accelerate future materials design.

The database component of JARVIS consists of JARVIS-DFT for DFT calculations and JARVIS-FF for molecular dynamics simulations. JARVIS-ML hosts several machine learning models based on our datasets. JARVIS-Tools contains tools for automating, post-processing and disseminating generated data, as well as several derived applications such as JARVIS-Heterostructure. We also include precision and accuracy analyses of the generated data, which consists of comparing DFT data with experiments, comparing FF data with DFT, comparing ML models with DFT, etc. As a

lower-level technique (see Fig. 1), JARVIS-DFT data can be fed into JARVIS-FF and JARVIS-ML models, but not vice versa. We use JARVIS-ML to accelerate both JARVIS-DFT and JARVIS-FF. In this way, the JARVIS-infrastructure establishes a joint integration for automation and generation of repositories. We provide several social-media platforms to build a community of interest. Some of the key resources for the JARVIS-infrastructure are shown in Table 1.

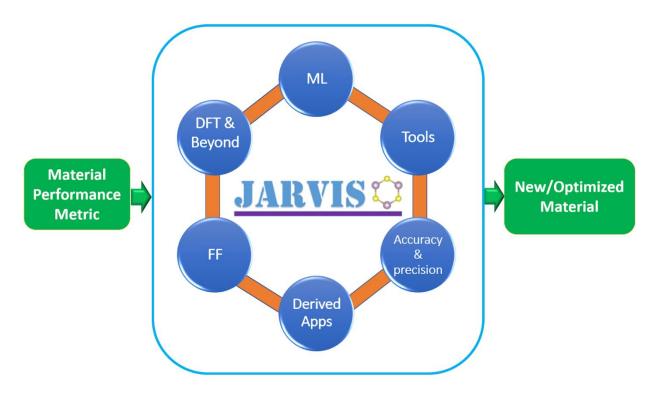


Fig. 2 An overview of the JARVIS infrastructure. For a given materials performance metric, several JARVIS components can work together to design optimized or completely new materials.

Table: 1 An overview of resources available in the JARVIS infrastructure.

Resource	Website	Brief description	
Homepage	https://jarvis.nist.gov/	Description and API	
FF	https://www.ctcms.nist.gov/~knc6/periodic.html	Evaluation of classical force field	
DFT	https://www.ctcms.nist.gov/~knc6/JVASP.html	Density functional theory data	
ML	https://www.ctcms.nist.gov/jarvisml/	Machine learning models	
Tools	https://github.com/usnistgov/jarvis	Scripts for running simulations	
Downloads	https://www.ctcms.nist.gov/~knc6/downloads.html	Downloadable metadata	
Notebooks	https://github.com/JARVIS-Materials- Design/jarvis-tools-notebooks	Jupyter/Google-Colab notebooks	
Heterostruct	https://www.ctcms.nist.gov/jarvish/	2D heterostructure properties	
WannierTB	https://www.ctcms.nist.gov/jarviswtb/	Wannier tight binding models	
BeyondDFT	https://www.ctcms.nist.gov/~knc6/BDFT.html	High-level ab-initio methods	
Publications	https://www.ctcms.nist.gov/~knc6/pubs.html	JARVIS-related publication	
Tools docs	https://jarvis-tools.readthedocs.io/en/latest/	Documentation (docs) of tools	
DB docs	https://www.ctcms.nist.gov/~knc6/documentation.html	Documentation on the database	
Tools pypi	https://pypi.org/project/jarvis-tools/	Pypi repository of tools	
Workshops	https://www.ctcms.nist.gov/~knc6/workshops.html	JARVIS-related workshops	
ResearchG.	https://www.researchgate.net/project/NIST- JARVIS	Social media researchgate page	
Twitter	https://twitter.com/jarvisnist	Social media twitter page	
Facebook	https://www.facebook.com/jarvisnist/	Social media Facebook page	
Linkedin	https://www.linkedin.com/company/jarvisnist	Social media Linkedin page	
YouTube	https://www.youtube.com/channel/UClChK_t7km Vx_QMStQH_T9g	Social media Youtube page	
Google group	https://groups.google.com/forum/#!forum/jarvis- nist	Social media google-group	

#### **JARVIS-Tools**

JARVIS-Tools is a python-based software package with ≈ 20,000 lines of code and consisting of several python-classes and functions. JARVIS-Tools can be used for a) the automation of simulations and data-generation, b) post-processing and analysis of generated data, and c) the dissemination of data and methods, as shown in Fig. 3. It uses cloud-based continuous integration checking including GitHubAction, CircleCI, TravisCI, CodeCov, and PEP8 linter to maintain consistency in the code and its functionalities. The JARVIS-Tools is distributed through an open GitHub repository: https://github.com/usnistgov/jarvis.

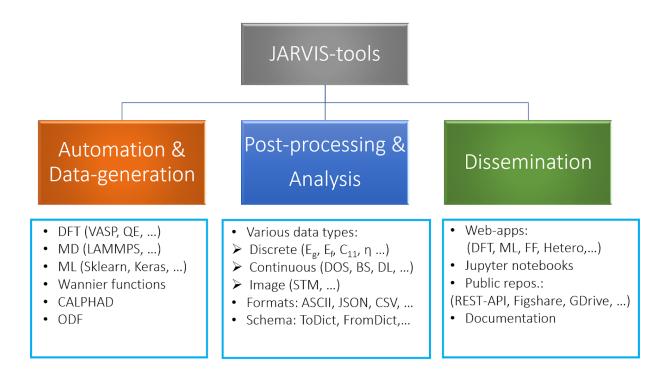


Fig. 3 Three main components of the JARVIS-Tools package and their capabilities.

An example python class in JARVIS-Tools is 'Atoms'. It uses atomic coordinates, element types and lattice vectors to build an 'Atoms' object from which several properties, such as density and chemical formula, can be calculated. This 'Atoms' class, along with several other modules

(discussed later), can be used for setting up calculations with external software packages. An example of the 'Atoms' class is shown in Fig. 4.

```
>>> from jarvis.core.atoms import Atoms
>>> box = [[2.715, 2.715, 0], [0, 2.715, 2.715], [2.715, 0, 2.715]]
>>> coords = [[0, 0, 0], [0.25, 0.25, 0.25]]
>>> elements = ["Si", "Si"]
>>> Si = Atoms(lattice_mat=box, coords=coords, elements=elements)
>>> density = round(Si.density,2)
>>> print (density)
2.33
>>>
>>> from jarvis.db.figshare import data
>>> dft_3d = data(dataset='dft_3d')
>>> print (len(dft_3d))
36099
```

Fig. 4 Examples of using python classes in JARVIS-Tools for constructing 'Atoms' class and downloading data.

The 'Atoms' class along with many other modules in JARVIS-Tools are used to generate input files for automating software codes. Currently, JARVIS-Tools can be used to automate DFT calculations with packages such as Vienna Ab-initio simulation package (VASP)<sup>62,63</sup>, Quantum Espresso (QE)<sup>64</sup>; MD with Large-scale Atomic/Molecular Massively Parallel Simulator (LAMMPS)<sup>65</sup>; ML with Scikit-learn<sup>66</sup>, Keras<sup>67</sup>, and LightGBM<sup>68</sup>; Wannier calculations with Wannier90<sup>69</sup> and Wanniertools<sup>70</sup>. A number of predefined workflows are available in JARVIS-Tools that are continuously being used to calculate properties of new or existing materials in the database. Three workflows are shown in Fig. 5. For DFT calculations, an input Atoms class is used to generate input files for VASP (Fig. 5a) with the 'VaspJob' class in order to calculate the desired properties, such as the energy. We automatically perform calculations to converge

numerical parameters like the k-points and plane-wave cut-off for individual materials. Geometry optimization is then carried out with energy, force, and stress relaxation. We have chosen a particular set of pseudopotentials or PAWs as tested and recommended by the software developers of various codes. Subsequent properties, such as band structure, dielectric function, elastic constants, piezoelectric constants or spin-orbit spillage are computed on the relaxed structure. Later, custom jobs can also be run on the optimized structure using 'VaspJob', such as Wannier90 calculations using the 'Wannier90Win' class, which generates the input files for an Atom class and a chosen set of pseudopotentials, disentanglement window and other controlling parameters. All of these steps produce a JavaScript Object Notation (JSON) file once the calculations are done as a signature of their completion. The workflows can be restarted from intermediate computations, making the calculations robust to interruptions due to computer failure, etc. We also add several error-handlers in the workflows to automatically re-submit a calculation if a typical error is encountered.

A similar workflow is shown for an example of FF based on LAMMPS calculations in Fig. 5b. Here, for a particular force-field such as Ni-Al<sup>58</sup>, for example, all the structures related to Ni, Al and Ni-Al are obtained from the DFT database and converted into a LAMMPS input format using 'Atoms', 'LammpsData' and 'LammpsJob' objects. Then a series of geometry optimization, vacancy formation energy, surface energy, and phonon-related calculations are run, based on the symmetry of the structure. All of these steps use a set of ".mod" module files with input parameters that control respective LAMMPS calculations. The obtained results are compared with corresponding DFT data, to evaluate the quality of an FF for a particular system or simulation.

In machine learning calculations, the input materials-data is transformed into several machinereadable descriptors<sup>71</sup> such as CFID dataset or STM image 'numpy' arrays. As we are not going to generate new data for testing ML models, we split the dataset into training and testing sets in a 90: 10 or similar split. Using k-fold cross-validation, we obtain hyperparameters for the chosen algorithm, for example, the number of trees, learning rate, etc. in the case of Gradient Boosting Decision Tree (GBDT). We choose the optimized parameters and train on 90 % train data and test on the 10 % test data to evaluate the truly predictive performance on unseen data. We also carry out k-fold cross-validation using the finalized model to get model uncertainty. Later, we can analyze interpretability with techniques such as feature importance in tree-based algorithms or filters in neural networks. These models are saved in Pickle, cPickle and Joblib modules for model persistency. We also carry out uncertainty analysis using methods such as prediction interval and Monte-Carlo dropouts<sup>72</sup>. A few examples and Jupyter notebooks are provided on the GitHub page to illustrate the above-mentioned methods. More details about the individual python modules mentioned above can be found in the JARVIS-Tools documentation (https://jarvistools.readthedocs.io/en/latest/).

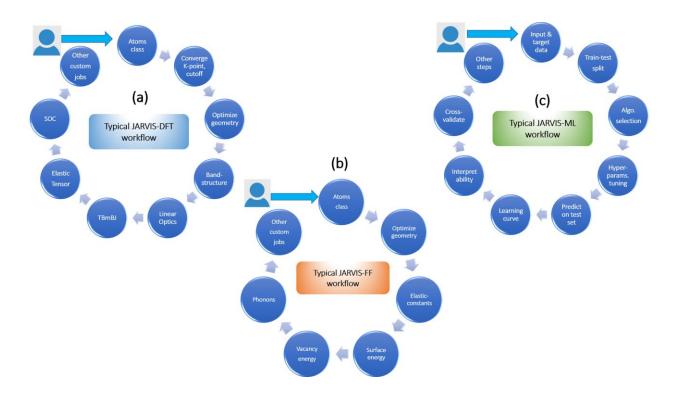


Fig. 5 Flowcharts showing some of the main steps used in most-commons calculations a) JARVIS-DFT, b) JARVIS-FF and c) JARVIS-ML workflows.

After running the automated calculations, the data is post-processed to predict various material properties (such as bandgap, formation energy, spin-orbit spillage, SLME, density of states, phonons, dielectric function, or STM image). Many of the python classes use 'ToDict' and 'FromDict' methods that help store the metadata. These metadata are then used with HTML<sup>73</sup>, Javascript, Flask<sup>74</sup> and other related software to make web-pages and web-apps. The metadata is also shared in public repositories **Figshare** such as (https://figshare.com/authors/Kamal Choudhary/4445539), and JARVIS-Representational state transfer (REST) API, based on the MGI philosophy of creating and using interoperable datasets. Note that through the JARVIS-REST API, a user can download JARVIS data and can also store their own data. The data generated in JARVIS is mainly stored in JavaScript Object Notation (JSON), Comma-Separated Values (CSV) or American Standard Code for Information Interchange (ASCII) format and, again, JARVIS-Tools can be used to analyze the pre-calculated data for materials design. An example of downloading precalculated dataset with JARVIS-Tools is shown in Fig. 4. JARVIS-Tools, along with the various software shown in Fig. 3, has led to several novel databases shown in Fig. 6.

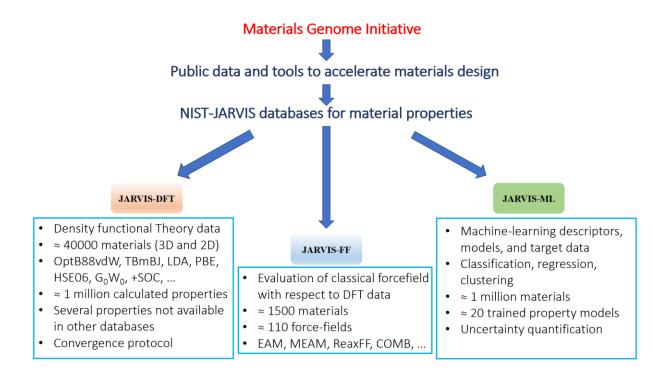


Fig. 6 Three main databases in JARVIS and a summary of their contents.

#### **JARVIS-DFT**

Density functional theory is one of the most commonly used techniques in condensed-matter physics to solve real-world materials problems. In DFT, instead of solving the fully interacting Schrödinger equation, we solve the Kohn-Sham equations, which describe an effective non-interacting problem, greatly improving computational efficiency. Although exact in principle, DFT requires several approximations in practice. In particular, various levels of approximation to the exchange-correlation functional are possible, which require different computational effort. Most existing DFT databases use the common GGA-PBE throughout all the material-classes. JARVIS-DFT can be viewed as an attempt to build a repository beyond existing DFT databases. JARVIS-DFT<sup>7,8,25-27,30,31,49,53,55</sup> was started in 2017 and contains data for  $\approx 40,000$  materials, with  $\approx 1$  million calculated properties, mainly based on the VASP package. Although there are several DFT-functionals adopted in JARVIS-DFT, we use vdW-DF-OptB88 consistently for all the 3D,

2D, 1D and 0D materials. This functional has been shown to provide accurate predictions for lattice-parameters and energetics for both vdW and non-vdW bonded materials<sup>30</sup>. In addition to hosting 3D bulk materials, the database consists of 2D monolayer, 1D-nanowire, and 0Dmolecular materials (as shown in Table 2). However, to date, 3D and 2D materials have primarily been distributed publicly. Moreover, other exchange-correlation functionals are considered (as shown in Table 3), which can help estimate the prediction uncertainty. While vdW-DF-OptB88 can predict accurate lattice parameters and formation energies, bandgaps are still underestimated. Calculations with hybrid functionals (such as range-separated HSE06 and PBE0) and many-body approaches (such as  $G_0W_0$ ) remain too computationally expensive<sup>23</sup> to use in a high-throughput methodology for thousands of materials. Hence, a meta-GGA Tran-Blaha-modified Becke-Johnson (TBmBJ) potential is used to provide a good balance between computational expense and accuracy. The TBmBJ accuracy is shown to be close enough to the high-level methods such as HSE06 at up to ten times lower computational expense<sup>57</sup>. Accurate prediction of optical gaps by calculation of the frequency-dependent dielectric function is important for several applications, for example, solar-cell efficiency calculations. Accurate prediction of bandgaps also helps in obtaining accurate frequency-dependent dielectric functions, which can be critical for solar-cell efficiency calculations; however, TBmBJ cannot describe the excitonic nature of electron-hole pairs in lowdimensional materials. In addition to TBmBJ, we are generating HSE06, PBE0, G<sub>0</sub>W<sub>0</sub> and DMFT datasets, which can be considered as beyond-DFT methods discussed in the next section. Next, SOC is varied to analyze the differences introduced by this coupling. These differences are used to discover 3D and 2D topological materials. In addition, several new DFT databases are developed including properties such as frequency-dependent dielectric function and electric field gradient. A few important protocols such as k-point automatic convergence are also introduced. A

snapshot of the JARVIS-DFT website along with a list of properties that are available is shown in Fig. 7. JARVIS-DFT has several filtering options on the website to screen candidate materials. We provide the input files as downloadable .zip files, especially for the users who do not have much expertise in using python-based codes. Raw input and output files (on the order of 1 terabyte) will soon be made publicly available through the Figshare repository, NIST-Materials data repository, and Materials Data Facility (MDF). A summary table, with the number of data available with vdW-DF-OptB88 and other methods, is shown in Tables 2 through 4.

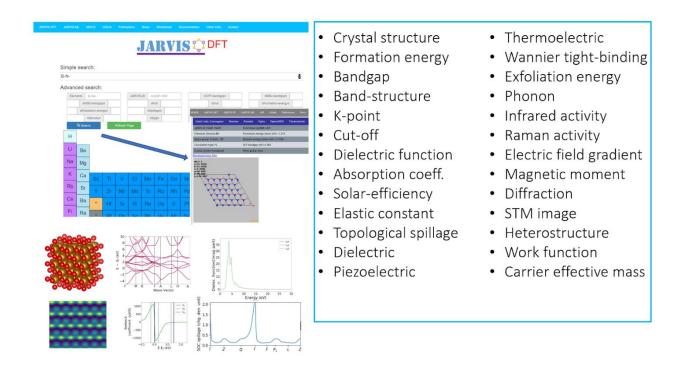


Fig. 7 A snapshot of JARVIS-DFT website and summary of its contents.

Table. 2 A brief summary of datasets available in the JARVIS-DFT.

Material classes	Numbers
3D-bulk	33482
2D-bulk	2293
1D-bulk	235
0D-bulk	413
2D-monolayer	1105
2D-bilayer	102
Molecules	12
Heterostructure	3
<b>Total DFT calculated systems</b>	37646

Table. 3 A brief summary of functionals used in optimizing crystal geometry in the JARVIS-DFT.

Functionals	Numbers	
vdW-DF-OptB88 (OPT)	37646	
vdW-DF-OptB86b (MK)	109	
vdW-DF-OptPBE (OR)	111	
PBE	99	
LDA	92	

*Table. 4 A brief summary of material-properties available in the JARVIS-DFT. The database is continuously expanding*<sup>7,8,25-27,30,31,49,53,55</sup>.

Property	Numbers
Optimized crystal-structure (OPT)	37646
Formation-energy (OPT)	37646
Bandgap (OPT)	37646
Exfoliation energy (OPT)	819
Bandgap (TBmBJ)	15655
Bandgap (HSE06)	40
Bandgap (PBE0)	40
Bandgap (G <sub>0</sub> W <sub>0</sub> )	15
Bandgap (DMFT)	11
Frequency dependent dielectric tensor (OPT)	34045
Frequency dependent dielectric tensor (TBmBJ)	15655
Elastic-constants (OPT)	15500
Finite-difference phonons at Γ-point (OPT)	15500
Work-function, electron-affinity (OPT)	1105
Theoretical solar-cell efficiency (SLME) (TBmBJ)	5097
Topological spin-orbit spillage (PBE+SOC)	11500

Wannier tight-binding Hamiltonians (PBE+SOC)	1771
Seebeck coefficient (OPT, BoltzTrap)	22190
Power factor (OPT, BoltzTrap)	22190
Effective mass (OPT, BoltzTrap)	22190
Magnetic moment (OPT)	37528
Piezoelectric constant (OPT, DFPT)	5015
Dielectric tensor (OPT, DFPT)	5015
Infrared intensity (OPT, DFPT)	5015
DFPT phonons at Γ-point (OPT)	5015
Electric field gradient (OPT)	15187
Non-resonant Raman intensity (OPT, DFPT)	250
Scanning tunneling microscopy images (PBE+SOC)	770

# **JARVIS-Beyond-DFT**

While quantum mechanical methods in single-particle theories such as DFT or DFT+U methods (mainly GGA) are fast and can predict accurate results for most structural parameters, even when relatively strong electron correlations are present, qualitative predictions of excited state properties may require beyond-DFT methods<sup>75</sup>. Beyond-DFT calculations have been applied to many materials systems, including cuprates and Fe-based high-temperature superconductors, Mott insulators, heavy Fermion systems, semiconductors, photovoltaics, and topological Mott insulators<sup>75</sup>. In the last few decades, both perturbative and stochastic approaches have been developed to understand these strongly correlated materials. These approaches, such as the GW approximation, Dynamical Mean Field Theory (DMFT)<sup>76</sup>, or hybrid functionals are often called beyond-DFT methods since they go beyond the limit of semilocal DFT. The materials design community needs to have a way of answering the question of whether, in a particular case, it is necessary to use a beyond-DFT method, and most importantly which method to use. In the JARVIS-Beyond-DFT database we are building a database of spectral functions and related quantities as computed using meta-GGA, GW, hybrid functionals, and LDA+DMFT for head-tohead comparison on 100+ materials.

In the JARVIS-Beyond-DFT<sup>75</sup> database we try to answer a few key questions regarding discoveries through a materials database for quantum materials. First, where is it necessary to use a beyond-DFT method, and which method to be use? Second, how do different "beyond-DFT" methods compare with experiments? Target materials include but are not limited to various transition metal oxides, perovskites and mixed perovskites, nickelates, transition metal dichalcogenides, and a wide range of metals starting from alkali metals to transition metals, and various Iron-based superconductors. JARVIS-Beyond-DFT will be distributed through the website: https://www.ctcms.nist.gov/~knc6/BDFT.html.

#### **JARVIS-FF**

Classical force-field-/interatomic-potential-based simulations are the workhorse technique for large scale atomistic simulations. They are especially suited for temperature-dependent and defect-related phenomena. Several varieties of FFs differ based on the materials system and the underlying phenomena under investigation, e.g., whether they include bond-angle information and fixed or dynamic charges. Also, they are generally designed for particular applications and phases, making it difficult to ascertain whether they will perform well in simulations for which they were not explicitly trained. JARVIS-FF<sup>27,58</sup> is a collection of LAMMPS calculation-based data consisting of crystal structures, formation energies, phonon densities of states, band structures, surface energies and defect formation energies. There are  $\approx 110$  FFs in the database, for which the corresponding crystal structures are obtained from JARVIS-DFT, converted to LAMMPS format inputs, and used in a series of LAMMPS calculations to produce the aforementioned properties. These properties, when compared with corresponding DFT data, can help a user analyze the quality of a force-field for a particular application. Examples include the comparison of DFT convex hull with FF, elastic modulus, surface energy and vacancy formation energy data. Some types of FFs

included are EAM, MEAM, Bond-order and Tersoff, COMB, and ReaxFF as shown in Table. 5<sup>-</sup> Furthermore, we plan to include several recently developed machine learning force-fields into JARVIS-FF. A snapshot of the JARVIS-FF website is also shown in Fig. 8.

*Table. 5 A summary of various types of force-fields available in the JARVIS-FF*<sup>27,58</sup>.

Force-fields	Numbers
EAM	92
Tersoff	9
ReaxFF	5
COMB	6
AIREBO	2
MEAM	1
EIM	1

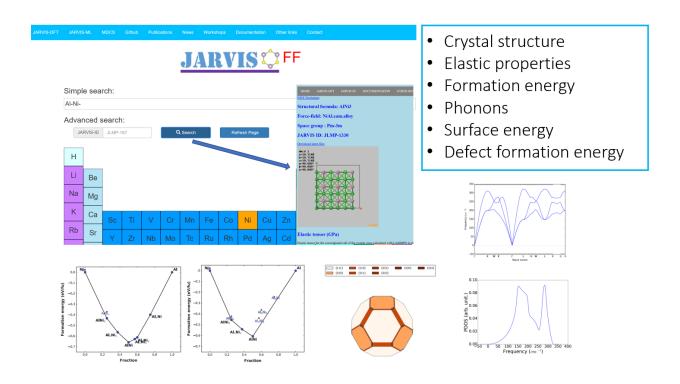


Fig. 8 A snapshot of JARVIS-FF website and summary of its contents.

#### **JARVIS-ML**

Machine learning has several applications in materials science and engineering <sup>14,80,81</sup>, such as automating experimental data analysis, discovering new functional materials, optimizing known ones by accelerating conventional methods such as DFT, automating literature searches, discovering new physical equations, and efficient clustering of materials and their properties. There are several data types that can be used in ML such as scalar data (e.g., formation energies, bandgaps), vector/spectra data (e.g., density of states, dielectric function, charge density, X-ray diffraction patterns, etc.), image-based data (such as scanning tunneling microscopy and transmission electron microscopy images), and natural language processing-based data (such as scientific papers). In addition, ML can be applied on a variety of materials classes such as bulk crystals, molecules, proteins and free-surfaces.

Currently, there are two types of data that are machine-learned in JARVIS-ML<sup>49,53,55,59,60</sup>: discrete and image-based. The discrete target is obtained from the JARVIS-DFT database for 3D and 2D materials. There have been several descriptor developments as attempts to capture the complex chemical-structural information of a material<sup>71</sup>. We compute CFID descriptors for most crystal structures in various databases (as shown in Table. 6). Many of these structures are non-unique but can still be used for pre-screening applications<sup>49</sup>. The CFID can also be applied to other materials classes such as molecules, proteins, point defects, free surfaces, and heterostructures, which are currently ongoing projects. These descriptor datasets, along with JARVIS-DFT and other databases, act as input and outputs for machine learning algorithms. The CFID consists of 1557 descriptors for each material: 438 average chemical, 4 simulation-box-size, 378 radial charge-distribution, 100 radial distribution, 179 angle-distribution up to first neighbor, and another 179 for the second neighbor, 179 dihedral angle up to fist neighbor and 100 nearest neighbor

descriptors. More details can be found in Ref. <sup>59</sup>. Currently, we provide CFID descriptors only, but other descriptors such as Coulomb-matrix, and sine-matrix will be provided soon. With CFID descriptors, we trained several classification and regression tasks. Once these models are trained, parameters are stored that can predict the properties of an arbitrary compound quickly. We developed a web-based application to host the trained models, as shown in Fig. 9, and a list of the trained properties are displayed there as well. We note that classical quantities such as bulk modulus, maximum infrared (IR) active mode, and formation energies can be accurately trained, especially with regression models. For other properties such as bandgaps, magnetic moments, piezoelectric coefficients, thermoelectric coefficients, high accuracy models are obtained for classification tasks only. In addition to the descriptor-based data, we develop Scanning Tunneling Microscopy (STM)<sup>53</sup> image classification models that can be used to accelerate the analysis of STM data. The images are converted into a black/white image to identify spots with/without atoms. The model's accuracy is compared with respect to DFT data or experiments wherever applicable.

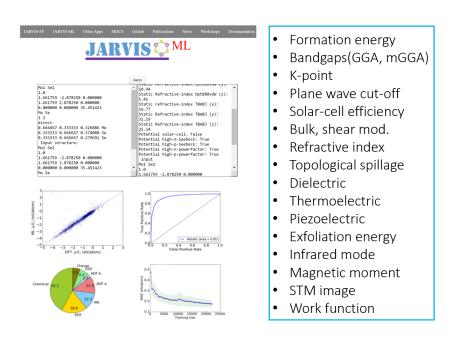


Fig. 9 A snapshot of JARVIS-ML website and summary of its contents.

Table. 6 A summary of Classical Force-field Inspired Descriptors (CFID)-descriptor datasets available in the JARVIS-ML.

CFID-Dataset	Number of materials
JARVIS-DFT 3D	35984
JARVIS-DFT 2D	1105
AFLOW	639262
OQMD	360802
Materials-project	82125
Crystallography Open Database (COD)	11783
QM9	13385
Total	1144446

## **Derived apps**

The knowledge developed through the above-mentioned databases and tools can serve as static content, as well as accessed through dynamic user-defined inputs. Derived applications (apps) are designed to help a user analyze the combinatorics in the data. Based on the databases and tools discussed above, several apps are derived from JARVIS such as JARVIS-Heterostructure<sup>31</sup>, JARVIS-Wannier TB, and JARVIS-ODF. JARVIS-Heterostructure (as shown in Fig. 10a) can be used to characterize heterojunction type and modeling interfaces for exfoliable 2D materials. We classify these heterostructures into type-I, II and III systems according to Anderson's rule, which is based on the band-alignment with respect to the vacuum potential of non-interacting monolayers, obtained from JARVIS-DFT. The app also generates crystallographic positions for the heterostructure that could be used as input for subsequent calculations. JARVIS-WannierTB (as shown in Fig. 10b) can be used to solve Wannier Tight Binding Hamiltonians on arbitrary k-points for 3D and 2D materials. Properties such as the band structure and the density of states can be predicted on the fly from this app. Additionally, many other apps are being developed, which are primarily based on the Flask python package<sup>74</sup>.

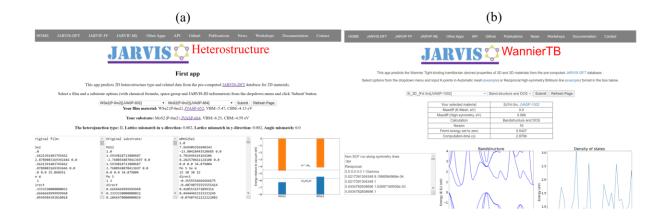


Fig. 10 Snapshots of JARVIS-DFT derived apps: JARVIS-Heterostructure and JARVIS-Wannier Tight Binding.

The JARVIS-ODF (Orientation Distribution Function) library is under development, which aims to calculate volume-averaged (meso-level) material properties, including the elasto-plastic deformation behavior, using the property data available for single crystals in the JARVIS database. Once generated, the JARVIS-ODF library will be capable of obtaining such material properties for all crystalline structures.

# Accuracy and precision analysis

In simulations, accuracy refers to the degree of closeness between a calculated value and a reference value, which can be from an experiment or a high-fidelity theory. Precision refers to the degree of closeness between numerical approaches to solving a certain model, including the effect of convergence and other simulation parameters.

In JARVIS-DFT, the accuracy of the DFT data is obtained by comparing it to available experimental results. The accuracy of JARVIS-FF and JARVIS-ML, instead, is given with respect to DFT results. Note that the numbers of high-quality experimental measurements or high-fidelity

calculations for a given property are often low. Therefore, the accuracy metrics we derive in our works are obtained only for the few cases we can directly compare, not for the entire dataset. Below, we provide accuracy metrics for some material properties in the JARVIS-DFT, with respect to experiments. In addition to the scalar data, vector/continuous data, such as frequency dependent dielectric function and Scanning Tunneling Microscopy (STM) images, are compared to a handful of experimental data points as well. Details of individual properties can be found in Ref. 8,30,49,52,53,55,59,60

Table. 7 Mean Absolute Error (MAE) for JARVIS-DFT data with respect to available experimental data for various material properties.

Property	#Materials	MAE	Typical range
Formation energy (eV/atom)	1317	0.128	-4 to 2
OptB88vdW-bandgaps (eV)	54	1.33	0 to 10
TBmBJ-bandgaps (eV)	54	0.51	0 to 10
Bulk modulus (GPa)	21	5.75	0 to 250
Electronic ( $\varepsilon_{11}$ ) OPT	28	3.2	0 to 60
Electronic (ε <sub>11</sub> ) MBJ	28	2.62	0 to 60
Solar-eff. (SLME) (%) (MBJ)	5	6.55	0 to 33.7
Max. piezoelectric strain coeff (Cm <sup>-2</sup> )	16	0.21	0 to 2
Dielectric constant $(\varepsilon_{11})$ (DFPT)	16	2.46	0 to 60
Seebeck coefficient (µV/K)	14	54.7	-600 to 600
Electric field gradient V <sub>zz</sub> (10 <sup>21</sup> Vm <sup>-2</sup> )	37	1.17	0 to 100
IR mode (cm <sup>-1</sup> )	8	8.36	0 to 4000

JARVIS-FF data accuracy is calculated with respect to the DFT data, for properties such as the convex hull, bulk modulus, phonon frequencies, vacancy formation energies and surface energies. In Refs.27,58, we showed this through several examples, including the comparison of Ni-Al and Cu-O-H systems convex hulls to DFT data. We also showed examples of comparing defect formation energies, surface energies and its effects on Wulff-shape. Although these accuracy

analyses are based on 0K DFT data, they are useful in predicting temperature-dependent and dynamical behavior because we consider several crystal prototypes of a system.

JARVIS-ML model accuracy is evaluated on the test-set (usually 10 %) representing newly computed and previously unseen DFT data for both regression and classifications models. Accuracy of regression and classification models are reported in terms of mean absolute error (MAE) and Receiver Operating Characteristic (ROC) Area Under Curve (AUC) metric respectively. A brief summary of regression and classification model accuracy results is given below in Table. 8 and 9. Details of the accuracy analyses are provided in Refs. <sup>49,53,55,59,60</sup>

Table. 8 Performance of regression machine learning models in JARVIS-ML with JARVIS-DFT data using OptB88vdW (OPT) and TBmBJ (MBJ) with mean absolute error (MAE). The mean absolute deviation (MAD) of properties are also included.

Property	Training data	MAE	MAD
Formation energy (eV/atom)	24549	0.12	0.81
OPT bandgap (eV)	22404	0.32	1.05
MBJ bandgap (eV)	10499	0.44	1.60
Bulk mod., Kv (GPa)	10954	10.5	49.95
Shear mod., Gv (GPa)	10954	9.5	23.26
Refr. Index(x) (OPT)	12299	0.54	1.15
Refr. Index(x) (MBJ)	6628	0.45	1.03
IR mode (OPT) (cm <sup>-1</sup> )	3411	77.84	316.7
Max. Born eff. charge (OPT)(e)	3411	0.60	1.48
Plane-wave cutoff (OPT)(eV)	24549	85.0	370.6
K-point length (OPT)(Å)	24549	9.09	22.23
2D-Exfoliation energy(OPT) (eV/atom)	616	37.3	46.09

Table. 9 Performance of the classification machine learning models in JARVIS-ML with JARVIS-DFT data using OptB88vdW (OPT) and TBmBJ (MBJ) with Receiver Operating Characteristic (ROC) Area Under Curve (AUC) metric. Random guessing and perfect ROC AUC are 0.5 and 1 respectively.

Property	Number of datapoints	ROC AUC
Metal/non-metal (OPT)	24549	0.95
Magnetic/Non-magnetic (OPT)	24549	0.96
High/low solar-cell efficiency (TBmBJ)	5097	0.90
High/low piezoelectric coeff	3411	0.86
High/low Dielectric	3411	0.93
High/low n-Seebeck coeff	21899	0.95
High/low n-power factor	21899	0.80
High/low p-Seebeck coeff	21899	0.96
High/low p-power factor	21899	0.82

Precision analysis can refer to a wide variety of optional selections of simulation set-ups. Examples of precision analysis in JARVIS-DFT are using our convergence protocols for k-points and plane-wave cutoff, and the convergence of Wannier tight-binding Hamiltonians. Using a converged k-point mesh and plane-wave cutoff<sup>28</sup> for each individual material is necessary to obtain high-quality data. Note that these DFT convergences are carried out for energies of the system only, and not for other properties. However, we impose tight convergence parameters for both k-points and energy cutoff (0.001 eV/cell), which typically results in other physical quantities being converged as well. In JARVIS-FF, comparison across structure-minimization methods for calculating surface and vacancy formation energy values are examples of precision analysis<sup>27</sup>. We find that the FF simulation setups ('refine' and 'box' methods) have minimal effect on the FF-based predictions.

For classification ML models, precision is the ratio  $\frac{TP}{TP + FP}$  where TP is the number of true positives and FP the number of false positives, which can be derived from the confusion. Precision analysis for classification ML model for STM Bravais-lattices are available in Ref. <sup>53</sup> . We find high precision (more than 0.87) for all of the 2D-Bravais lattices. Precision analysis for regression tasks are still ongoing and will be available soon.

#### **Future work**

Given that the number of all possible materials <sup>77</sup> could be of the order of 10<sup>100</sup>, and furthermore existing materials properties can be computed at increasing levels of accuracy/cost, the JARVIS databases will always be incomplete. This represents an opportunity for JARVIS to be drastically expanded in the future. Future work will be aimed at addressing some of the limitations of the existing databases, and may include additions like defect/disorder properties, magnetic ordering, non-linear optoelectronics, more beyond-DFT calculations, temperature-dependent properties, integration with experiments, and more detailed uncertainty analysis. Moreover, new ML models and methods for data-prediction and uncertainty quantification will be developed for 'explainable AI' (XAI) and transfer-learning (TL)-based research. Other derived apps such as JARVIS-ODF, JARVIS-Beyond-DFT, JARVIS-GraphConv, and JARVIS-STM are also being developed. In addition to the technical aspects, the broader impact of the infrastructure will be to provide a research platform that will allow maximum participation of worldwide researchers. NIST-JARVIS currently hosts pre-computed data and would host on-the-fly calculation resources also.

In summary, we described the Joint Automated Repository for Various Integrated Simulations (JARVIS) platform, which consists of several databases and computational tools to help accelerate materials design and enhance industrial growth. JARVIS includes three major databases: JARVIS-

DFT for density functional theory calculations, JARVIS-FF for classical force-field calculations, and JARVIS-ML for ML predictions. In addition, we provide JARVIS-Tools, which is used to generate the databases. The generated data is provided publicly with several example notebooks, documentation and calculation examples to illustrate different components of the infrastructure. We believe the publicly available data and resources provided here will significantly accelerate futuristic materials-design in various areas of science and technology.

## Methods

The entire study was managed, monitored, and analyzed using the modular workflow, which we have made available<sup>54</sup> on our JARVIS-Tools GitHub page (https://github.com/usnistgov/jarvis). The DFT calculations are mainly carried out using the Vienna Ab-initio simulation package (VASP)<sup>62,63</sup>. We use the projected augmented wave method and OptB88vdW functional<sup>56</sup>, which gives accurate lattice parameters for both van der Waals (vdW) and non-vdW solids<sup>30</sup>. Both the internal atomic positions and the lattice constants are allowed to relax in spin-unrestricted calculations until the maximal residual Hellmann–Feynman forces on atoms are smaller than 0.001 eV Å-1 and energy-tolerance of 10-7 eV. We do not consider magnetic orderings besides ferromagnetic yet, because of a high computational cost. We note that nuclear spins are not explicitly considered during the DFT calculations. The list of pseudopotentials used in this work is given on the GitHub page. The k-point mesh and plane-wave cut-off were converged for each material using the automated procedure described in Ref<sup>28</sup>. The elastic constants are calculated using the finite difference method with six finite symmetrically distinct distortions. The thermoelectric coefficients such as power factor and Seebeck coefficients are obtained with the BoltzTrap code with Constant Relaxation Time approximation (CRTA)<sup>78</sup>. Optoelectronic properties such as dielectric function and solar-cell efficiency are calculated using linear-optics

methods mainly using OptB88vdW and TBmBJ. We also compared such data with HSE06 and  $G_0W_0$ . The piezoelectric, dielectric and phonon modes at  $\Gamma$ -point are calculated using Density Functional Perturbation Theory (DFPT). Topological spillage for identifying topologically non-trivial materials is calculated by comparing DFT wave functions with/without SOC<sup>7,26</sup>. 2D exfoliation energies are calculated by comparing bulk and 2D monolayer energy per atom. The 2D heterostructure<sup>31</sup> behavior is predicted using Zur and Anderson methods. Wannier tight binding Hamiltonians are generated using the Wannier90 code<sup>69</sup>. 2D STM images are predicted using the Tersoff-Hamman method<sup>53</sup>.

Classical force-field calculations are carried out with the LAMMPS software package<sup>65</sup>. In our structure minimization calculations, we used  $10^{-10}$  eVÅ<sup>-1</sup> for force convergence and 10000 maximum iterations. The geometric structure is minimized by expanding and contracting the simulation box with 'fix box/relax' command and adjusting atoms until they reach the force convergence criterion. These are commonly used computational set-up parameters. After structure optimization point vacancy defects are created using Wycoff-position data. Free surfaces for maximum miller indices up to 3 are generated. The defect structures were required to be at least 1.5 nm long in the x, y and z directions to avoid spurious self-interactions with the periodic images of the simulation cell. We enforce the surfaces to be at least 2.5 nm thick and with 2.5 nm vacuum in the simulation box. The 2.5 nm vacuum is used to ensure no self-interaction between slabs, and the slab-thickness is used to mimic an experimental surface of a bulk crystal. Using the energies of perfect bulk and surface structures, surface energies for a specific plane are calculated. We should point out that only unreconstructed surfaces without any surface-segregation effects are computed, as our high-throughput approach does not allow for taking into account specific,

element dependent reconstructions yet. Phonon structures are generated mainly using the Phonopy package interface<sup>79</sup>.

Machine learning models are mainly trained using Scikit-learn<sup>66</sup>, Keras<sup>67</sup>, and LightGBM<sup>68</sup> (TensorFlow backend) software. For DFT generated scalar data such as formation energies, bandgaps, exfoliation energies etc. the crystal structures are converted into a Classical Force-field Inspired Descriptors (CFID) input array and the DFT data is used as target data, which is then train-test split in a ratio of 90: 10. Preprocessing such as 'VarianceThreshold', 'StandardScalar' are used before ML training. Regression models' performance are generally reported in terms of Mean Absolute Error (MAE) or r<sup>2</sup>, while that for classification models using the Receiver Operating Characteristic (ROC) Area Under Curve (AUC) value which lie between 0.5 and 1.0. Several other analyses such as feature importance, k-fold cross validation and learning curve are carried out after the model training. The trained model is saved in pickle and joblib formats for model persistence. All the web-apps are developed using JavaScript, Flask and Django packages<sup>74</sup>.

## Data availability

JARVIS-related data is available at the JARVIS-API (<a href="https://jarvis.nist.gov">https://jarvis.nist.gov</a>), JARVIS-DFT (<a href="https://www.ctcms.nist.gov/~knc6/JVASP.html">https://www.ctcms.nist.gov/~knc6/JVASP.html</a>), JARVIS-FF

(https://www.ctcms.nist.gov/~knc6/periodic.html), JARVIS-ML

(<a href="https://www.ctcms.nist.gov/jarvisml/">https://www.ctcms.nist.gov/jarvisml/</a>) websites. The metadata is also available at the Figshare repository, see <a href="https://figshare.com/authors/Kamal Choudhary/4445539">https://figshare.com/authors/Kamal Choudhary/4445539</a>.

# **Code Availability**

Python-language based codes with examples are available at JARVIS-Tools page: <a href="https://github.com/usnistgov/jarvis">https://github.com/usnistgov/jarvis</a>.

# Acknowledgements

K.C., K.F.G., and F.T. thank the National Institute of Standards and Technology for funding, computational, and data-management resources. K.C. also thanks the computational support from XSEDE computational resources under allocation number TG-DMR 190095. Contributions by S.M., K.H., K.R., and D.V. were supported by NSF DMREF Grant No. DMR-1629059 and No. DMR-1629346. X.Q. was supported by NSF Grant No. OAC-1835690. B.G.S and S.V.K acknowledge work performed at the Center for Nanophase Materials Sciences, a US Department of Energy Office of Science User Facility. K.C. thanks for helpful discussion with several researchers including Faical Y. Congo, Daniel Wheeler, James Warren, Carelyn Campbell, Chandler Becker, Marcus Newrock, Ursula Kattner, Kevin Brady, Lucas Hale, Eric Cockayne, Philippe Dessauw from National Institute of Standards and Technology; Karen Sauer, Igor Mazin, Nirmal Ghimire, Patrick Vora from George Mason University; Rama Vasudevan, Maxim Ziatdinov from Oak Ridge National Lab, Deyu Lu and Matthew Carbone from Brookhaven National Lab; Marnik Bercx, Dirk Lamoen from University of Antwerp; Yifei Mo from University of Maryland; Anubhav Jain and Sinead Griffin from Lawrence Berkeley National Laboratory; Surya Kalidindi from Georgia Tech.; Tyrel McQueen and David Elbert from Johns Hopkins University; Richard Hennig from University of Florida; Giulia Galli and Ben Blaiszik from University of Chicago; Qiang Zhu from University of Nevada-Las Vegas; Dilpuneet Aidhy from University of Wyoming; Susan B. Sinnott, Tao Liang from Pennsylvania State University.

#### **Contributions**

KC designed the JARVIS workflows, carried out high-throughput calculations, analysis and developed the websites. FT contributed to the development of k-point and other convergence protocol, Beyond-DFT development and several other analyses. KG contributed to the development of topological materials discovery and Wannier-tight binding Hamiltonian projects. ACER assisted in the deployment of the web-apps. BDC, AA and AGK contributed to the machine-learning tasks. AJB, AHR, AC, VS, AD contributed to the phonon data analysis. ZT contributed to the development of the JARVIS-API website. JHS contributed to the experimental validation of some of the screened materials. JJ and RP contributed in the solar-cell and topological materials discovery tasks. GC, ER, XQ, HZ, SVK, BS, GP contributed to the discovery and characterization of low-dimensional materials. PA contributed to the elastic constant analysis task. SM, KR, DV and KH contributed to the Beyond-DFT project. All authors contributed to writing the manuscript.

# **Competing interests**

The authors declare no competing interests.

#### References

- 1 https://mgi.gov/.
- 2 https://www.nist.gov/mgi.
- 3 Curtarolo, S. *et al.* AFLOWLIB. ORG: A distributed materials properties repository from high-throughput ab initio calculations. *Computational Materials Science* **58**, 227-235 (2012).
- Jain, A. *et al.* Commentary: The Materials Project: A materials genome approach to accelerating materials innovation. *Apl Materials* **1**, 011002 (2013).
- Kirklin, S. *et al.* The Open Quantum Materials Database (OQMD): assessing the accuracy of DFT formation energies. npj Comp.Mat. **1**, 15010 (2015).
- Pizzi, G., Cepellotti, A., Sabatini, R., Marzari, N. & Kozinsky, B. AiiDA: automated interactive infrastructure and database for computational science. Comp. Mat. Sci. **111**, 218-230 (2016).

- 7 Choudhary, K., Garrity, K. F. & Tavazza, F. High-throughput Discovery of topologically Non-trivial Materials using spin-orbit spillage. Scientific Reports, **9**, 1-8 (2019).
- 8 Choudhary, K., Kalish, I., Beams, R. & Tavazza, High-throughput Identification and Characterization of Two-dimensional Materials using Density functional theory. Scientific Reports, **7**, 5179 (2017).
- 9 Draxl, C. & Scheffler, M. The NOMAD laboratory: from data sharing to artificial intelligence. J. Phys. Mats. **2**, 036001 (2019).
- 10 Chung, Y. G. *et al.* Computation-ready, experimental metal—organic frameworks: A tool to enable high-throughput screening of nanoporous crystals. Chem. Mater. **26**, 6185-6192 (2014).
- Green, M. L. *et al.* Fulfilling the promise of the materials genome initiative with high-throughput experimental methodologies. *J Applied Physics Reviews* **4**, 011105 (2017).
- Hattrick-Simpers, J. R., Gregoire, J. M. & Kusne, A. G. Perspective: Composition–structure–property mapping in high-throughput experiments: Turning data into knowledge. *APL Materials* **4**, 053211 (2016).
- Zakutayev, A. *et al.* An open experimental database for exploring inorganic materials.Sci. Data. **5**, 180053 (2018).
- Vasudevan, R. K. *et al.* Materials science in the artificial intelligence age: high-throughput library generation, machine learning, and a pathway from correlations to the underpinning physics. MRS Communications **9**, 821-838 (2019).
- Agrawal, A. & Choudhary, A. Perspective: Materials informatics and big data: Realization of the "fourth paradigm" of science in materials science. APL Maters. **4**, 053208 (2016).
- Schleder, G. R., Padilha, A. C., Acosta, C. M., Costa, M. & Fazzio, A. J. From DFT to machine learning: recent approaches to materials science—a review. J. Phys. Mater. **2**, 032001 (2019).
- 17 Ceder, G. J. Opportunities and challenges for first-principles materials design and applications to Li battery materials. MRS Bulletin **35**, 693-701 (2010).
- 18 Xi, L. *et al.* Discovery of high-performance thermoelectric chalcogenides through reliable high-throughput material screening. J. Am. Chem. Soc. **140**, 10785-10793 (2018).
- Olson, G. B. & Kuehmann, C. Materials genomics: from CALPHAD to flight. Scripta Materialia, **70**, 25-30 (2014).
- Aykol, M. *et al.* The Materials Research Platform: Defining the Requirements from User Stories. Matter, **1**, 1433-1438 (2019).
- Callister, W. D. & Rethwisch, D. G. *Materials science and engineering*. Vol. 5 (John wiley & sons NY, 2011).
- de Pablo, J. J. *et al.* The materials genome initiative, the interplay of experiment, theory and computation. Current Opinion in Solid State and Materials Science **18**, 99-117 (2014).
- 23 Sholl, D. & Steckel, J. A. *Density functional theory: a practical introduction*. (John Wiley & Sons, 2011).
- Perdew, J. P., Burke, K. & Ernzerhof, M. J. Generalized gradient approximation made simple. Phys. Rev. Lett. **77**, 3865 (1996).
- 25 Choudhary, K. *et al.* Computational screening of high-performance optoelectronic materials using OptB88vdW and TB-mBJ formalisms. Scientific Data, **5**, 180082 (2018).
- 26 Choudhary, K., Garrity, K. F., Jiang, J., Pachter, R. & Tavazza, F. Computational search for magnetic and non-magnetic 2D topological materials using unified spin—orbit spillage screening. npj Computational Materials **6**, 1-8 (2020).
- 27 Choudhary, K. *et al.* High-throughput assessment of vacancy formation and surface energies of materials using classical force-fields. J. Physics: Condensed Matter **30**, 395901 (2018).
- 28 Choudhary, K. & Tavazza, F. Convergence and machine learning predictions of Monkhorst-Pack k-points and plane-wave cut-off in high-throughput DFT calculations. Computational Materials Science **161**, 300-308 (2019).

- Cooper, M. *et al.* Development of Xe and Kr empirical potentials for CeO2, ThO2, UO2 and PuO2, combining DFT with high temperature MD. J. Phys.:Cond. Matt. **28**, 405401 (2016).
- 30 Choudhary, K., Cheon, G., Reed, E. & Tavazza, F. Elastic properties of bulk and low-dimensional materials using van der Waals density functional. Phys. Rev. B **98**, 014107 (2018).
- Choudhary, K., Garrity, K. F., Pilania, G. & Tavazza, F. Efficient Computational Design of 2D van der Waals Heterostructures: Band-Alignment, Lattice-Mismatch, Web-app Generation and Machine-learning. arXiv 2004.03025 (2020).
- 32 Allen, M. P. & Tildesley, D. J. Computer simulation of liquids. (Oxford university press, 2017).
- 33 Kattner, U. R. Phase diagrams for lead-free solder alloys. JOM **54**, 45-51 (2002).
- Acar, P., Ramazani, A. & Sundararaghavan, V. Crystal plasticity modeling and experimental validation with an orientation distribution function for ti-7al alloy. Metals, **7**, 459 (2017).
- Castelli, I. E. *et al.* New light-harvesting materials using accurate and efficient bandgap calculations. Adv. En. Mater. **5**, 1400915 (2015).
- NIMS-MatNavi database: https://mits.nims.go.jp/index\_en.html.
- 37 Stevanović, V., Lany, S., Zhang, X. & Zunger, A. Correcting density functional theory for accurate predictions of compound enthalpies of formation: Fitted elemental-phase reference energies. Phys. Rev. B **85**, 115104 (2012).
- Belsky, A., Hellenbrandt, M., Karen, V. L. & Luksch, P. New developments in the Inorganic Crystal Structure Database (ICSD): accessibility in support of materials research and design. Acta Crystallographica Section B, **58**, 364-369 (2002).
- Talirz, L. *et al.* Materials Cloud, a platform for open computational science. arXiv 2003.12510 (2020).
- 40 Ctrine informatics, https://citrine.io.
- Tadmor, E. B., Elliott, R. S., Sethna, J. P., Miller, R. E. & Becker, C. A. J. The potential of atomistic simulations and the knowledgebase of interatomic models. JOM **63**, 17 (2011).
- 42 Aagesen, L. *et al.* Prisms: An integrated, open-source framework for accelerating predictive structural materials science. JOM, **70**, 2298-2314 (2018).
- Wheeler, D. et al. PFHub: The Phase-Field Community Hub. 7 (2019).
- Ong, S. P. *et al.* Python Materials Genomics (pymatgen): A robust, open-source python library for materials analysis. Comp. Mat. Sci., **68**, 314-319 (2013).
- Larsen, A. H. *et al.* The atomic simulation environment—a Python library for working with atoms. **29**, 273002 (2017).
- Mathew, K. *et al.* MPInterfaces: A Materials Project based Python tool for high-throughput computational screening of interfacial systems. Comp. Mat. Sci., **122**, 183-190 (2016).
- 47 Setyawan, W., Gaume, R. M., Lam, S., Feigelson, R. S. & Curtarolo, S. High-throughput combinatorial database of electronic band structures for inorganic scintillator materials. ACS Comb. Sci., **13**, 382-390 (2011).
- 48 Gibbs, Z. M. *et al.* Effective mass and Fermi surface complexity factor from ab initio band structure calculations. npj Comp.Mat. **3**, 1-7 (2017).
- Choudhary, K. *et al.* Accelerated Discovery of Efficient Solar-cell Materials using Quantum and Machine-learning Methods. Chemistry of Materials, 31 (15), 5900, (2019).
- Yu, L. & Zunger, A. Identification of potential photovoltaic absorbers based on first-principles spectroscopic screening of materials. Phys. Rev. Lett., **108**, 068701 (2012).
- Liu, J. & Vanderbilt, D. Spin-orbit spillage as a measure of band inversion in insulators. Phys. Rev.B **90**, 125133 (2014).
- Jha, D. *et al.* Enhancing materials property prediction by leveraging computational and experimental data using deep transfer learning. Nature Communications, **10**, 1-12 (2019).

- Choudhary, K. *et al.* Density Functional Theory and Deep-learning to Accelerate Data Analytics in Scanning Tunneling Microscopy. arXiv:1912.09027 (2019).
- Please note that commercial software is identified to specify procedures. Such identification does not imply recommendation by the National Institute of Standards and Technology.
- Choudhary, K., Garrity, K. & Tavazza, F. Data-driven Discovery of 3D and 2D Thermoelectric Materials. arXiv:1906.06024 (2019).
- Klimeš, J., Bowler, D. R. & Michaelides, A. Chemical accuracy for the van der Waals density functional. J. Phys. Cond Matt. **22**, 022201 (2009).
- 57 Tran, F. & Blaha, P. Accurate band gaps of semiconductors and insulators with a semilocal exchange-correlation potential. Phys. Rev. Lett., **102**, 226401 (2009).
- Choudhary, K. *et al.* Evaluation and comparison of classical interatomic potentials through a user-friendly interactive web-interface. Scientific Data, **4**, 1-12 (2017).
- Choudhary, K., DeCost, B. & Tavazza, F. Machine learning with force-field-inspired descriptors for materials: Fast screening and mapping energy landscape. *Physical Review Materials* **2**, 083801, (2018).
- 60 Choudhary, K. *et al.* High-throughput Density Functional Perturbation Theory and Machine Learning Predictions of Infrared, Piezoelectric and Dielectric Responses. npj Computational Materials, 6, 64 (2020).
- 61 Saito, T. Computational materials design. Vol. 34 (Springer Science & Business Media, 2013).
- Kresse, G. & Furthmüller, J. Efficient iterative schemes for ab initio total-energy calculations using a plane-wave basis set. Phys. Rev. B **54**, 11169 (1996).
- Kresse, G. & Furthmüller, J. Efficiency of ab-initio total energy calculations for metals and semiconductors using a plane-wave basis set. Comp. Mat. Sci., 6, 15-50 (1996).
- Giannozzi, P. *et al.* QUANTUM ESPRESSO: a modular and open-source software project for quantum simulations of materials. J. Phys: Cond. Matt., **21**, 395502 (2009).
- 65 Plimpton, S. Fast parallel algorithms for short-range molecular dynamics. (Sandia National Labs., Albuquerque, NM (United States), 1993).
- Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. Journal of Machine Learning Research **12**, 2825-2830 (2011).
- 67 Gulli, A. & Pal, S. Deep learning with Keras. (Packt Publishing Ltd, 2017).
- 68 Ke, G. et al. in Advances in neural information processing systems. 3146-3154.
- 69 Mostofi, A. A. *et al.* wannier90: A tool for obtaining maximally-localised Wannier functions. Comp. Phys. Comm., **178**, 685-699 (2008).
- Wu, Q., Zhang, S., Song, H.-F., Troyer, M. & Soluyanov, A. A. WannierTools: An open-source software package for novel topological materials. Comp. Phys. Comm., **224**, 405-416 (2018).
- Ward, L. *et al.* Matminer: An open source toolkit for materials data mining. Comp. Mat. Sci., **152**, 60-69 (2018).
- Hammer, B. & Villmann, T. in *ESANN*. 79-90 (Citeseer).
- 73 Musciano, C. & Kennedy, B. *HTML, the definitive Guide*. (O'Reilly & Associates, 1996).
- Grinberg, M. *Flask web development: developing web applications with python*. ("O'Reilly Media, Inc.", 2018).
- Mandal, S., Haule, K., Rabe, K. M. & Vanderbilt, D. Systematic beyond-DFT study of binary transition metal oxides. npj Comp. Mat. **5**, 1-8 (2019).
- Kotliar, G. *et al.* Electronic structure calculations with dynamical mean-field theory. Rev. Mod. Phys., **78**, 865 (2006).
- Walsh, A. The guest for new functionality. Nat. Chem. **7**, 274-275 (2015).
- 78 Madsen, G. K. & Singh, D. BoltzTraP. A code for calculating band-structure dependent quantities. Comp. Phys. Comm., **175**, 67-71 (2006).

- Togo, A. & Tanaka, I. First principles phonon calculations in materials science. Scripta Mater., **108**, 1-5 (2015).
- Schmidt, Jonathan, et al. "Recent advances and applications of machine learning in solid-state materials science." npj Computational Materials 5.1 (2019): 1-36.
- Ramprasad, R., Batra, R., Pilania, G., Mannodi-Kanakkithodi, A. and Kim, C., 2017. Machine learning in materials informatics: recent applications and prospects. npj Computational Materials, 3(1), pp.1-13.