

Data dissemination and material informatics at LBNL Anubhav Jain, Energy Technologies Area, Berkeley Lab

Contact: ajain@lbl.gov

Data Dissemination

The Materials Project

A "science gateway" exposes data to the research community and facilitates knowledge extraction. The Materials Project (www.materialsproject.org) shares simulation data on hundreds of thousands of materials for a community of 22,000 researchers; millions of data points have been downloaded and used in the research articles of its users. A recently introduced feature, MPContribs, accepts data contributions from the user community.



ormula (Fe2O3) (1234) or chemica

Materials Informatics

MIDAS: Materials Informatics & Data Analysis Software

We are building a materials data mining platform called MIDAS (Materials Informatics & Data Analysis Software). MIDAS can retrieve data from several materials databases that have REST APIs or from the JSON format, format that data into a Pandas DataFrame object, automatically generate possible descriptors for the data, and run machine learning algorithms through scikit-learn. Visualization tools are provided through adapters to the plot.ly library. This system will be scaled for "big data analysis" in the future through Apache Spark and similar toolkits such as the FireWorks workflow software developed at LBNL.







CALCULATE EXPLORE MATERIAL calculate the enthal of 10,000+ reactions r lithium batteries. Ge ourbaix diagrams to ith our structure edite and substitution ind compare with or property algorithms experimental values waen evolution data **Database Statistics** 51,169 21,954 67,206 BANDSTRUCTURES MOLECULES INORGANIC COMPOUNDS 530,243 11 Stable Compounds Name ▲ Form ★ Decomp. ★ Id ★ NANOPOROUS MATERIALS



A Representational State Transfer Application Programming Interface (REST API) allows users to download large quantities of data using many popular programming languages. This method uses HTTP requests as a mechanism of accessing data and is employed by many of the top internet companies, including Google, Dropbox, and Twitter.



Data on over 500,000 porous

materials was contributed by



Left: data pipeline for MIDAS Right: example of an interactive plot that can be generated through MIDAS via the plot.ly toolkit.

Materials data mining

Using data sets from the Materials Project, we have derived structure-property relationships that relate fundamental descriptors such as composition, density, and coordination to output properties such as bulk modulus, shear modulus, and the electronic character of the valence and conduction bands. Crucial to this effort has been the development of relevant descriptor combinations as well as new



In this machine learning predictor for the bulk and shear moduli of a material, we use Hölder means to improve the predictive power of descriptors and develop a local linear regression method that fits the tails of distributions (i.e., extreme values) accurately without overfitting.



Several "apps", including the PDApp (which helps assess synthesis and stability) provide domain-specific views of the data. Plots and analyses from such apps are some of the most heavily used and cited aspects of the Materials Project.

MPContribs



the Nanoporous Materials Genome Center. Users can create custom plots based on properties of interest as well as obtain detailed information about specific materials.

Although the Materials Project is largely an open database, a sandboxing system allows for core data to be stored separately from data generated for proprietary or external projects. Access management is controlled via an API.

Unified View

Historically, all data in the Materials Project was generated at LBNL via highthroughput computing and density functional theory (DFT) calculations. MPContribs is a new feature that allows users to contribute their own data sets and

Hierarchical Data	Tables		
- (root): {} withing	data_Ni_XM	CD_Spectra	Sollapso/Expand
- Experiment: () ditema			Search:
+ Preparation: {} 1 Item	2.2.2	Sec.	5.2.4
→ Sample: {} Flums	Energy 11	XAS XM	ICD
Material_Name: Platinum doped Permalloy	821.0	0.0104183	-0.0004518
Form: -20nm film on Si wafer	822.0	0.00931404	-0.0009740
Thickness: ca. 20nm with 2-3 nm Au-capping (nominally)	823.0	0.00821621	-0.0008330
Grower: Ales Hrabec	004.0	0.00007000	0.0000000
Authors: Ales Hrabec, Alpha T. N'Diaye, Elke Arenholz, Christopher Marrows	824.0	0.00827328	-0.00069030
Measurement: {} 5 llems	825.0	0.00735342	-0.0002823
Beamline: {} 7 Items	826.0	0.00655145	-0.0011654
- NI_XMCD: () Telling	827.0	0.00595945	-0.0016242
<pre>> get_xmcd: {} 1 ltem</pre>	000 0	0.00590099	0.0005.407
+ xas_normalization_to_min_and_max:	828.0	0.00080988	-0.0005407
energy_range: 800 1000	Showing 1 to	451 of 451 en	tries
normalization_factor: 0.952002315041	data_Fe_XM	CD_Spectra	Collapse/Expand
offset: 0.358620768783			



machine learning approaches.



This diagram shows the pairwise likelihood for the electronic state on the y-axis to form a greater contribution to the valence or conduction band character versus the state on the x-axis. For example, d orbitals in Cu^{1+} are the most likely to form the VB. For this study, we repurposed algorithms used in ranking sports teams to rank electronic orbitals.

Automated materials design

Materials discovery today is largely performed through targeted experiments based on researcher intuition. However, such intuition is difficult to obtain within the context of highthroughput and combinatorial studies in which tens of thousands of data points may be collected. We are developing automated optimization routines that couple forward models (e.g. density functional theory calculations) with inverse optimizers (e.g., genetic algorithms or Gaussian processes) to build a fully automated materials discovery system.



Left: schematic diagram showing the integration of typical highthroughput screening coupled with an automatic optimization routine.