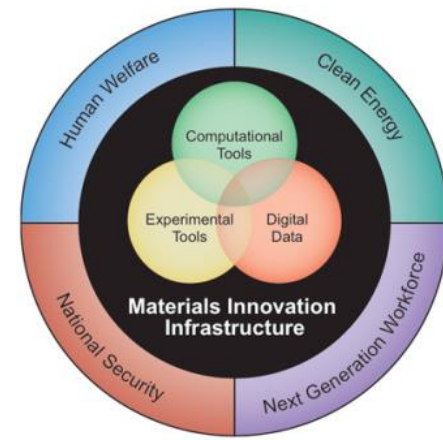


# JARVIS-ML: Physics inspired AI for fast and accurate screening of materials

## Crystals, Surfaces, Grain-boundaries, Molecules, Proteins

Kamal Choudhary, James Hickman, Brian DeCost, Francesca Tavazza  
NIST, Gaithersburg  
Aalto University, May 11



# Motivation

## Materials Genome Initiative



## National Quantum Initiative

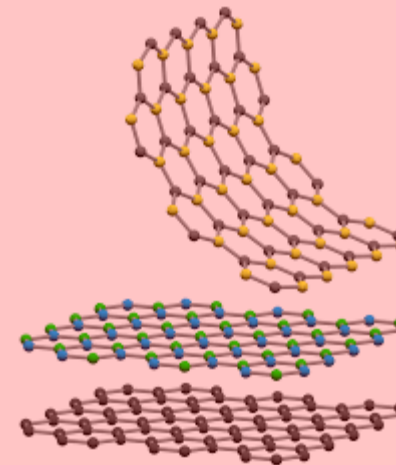
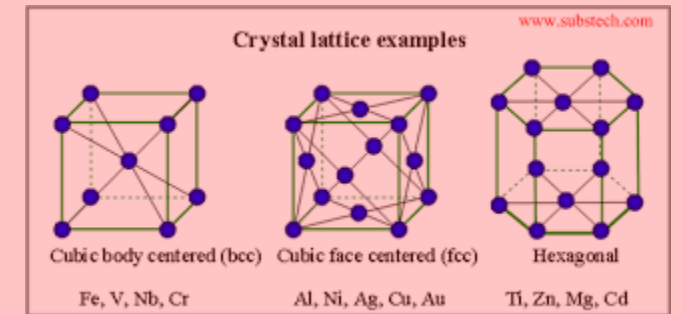
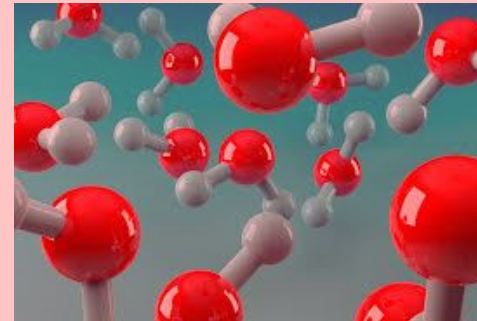
CONGRESS.GOV [Advanced Searches](#) | [Browse](#)

All Legislation

[Home](#) > [Legislation](#) > [115th Congress](#) > H.R.6227

**H.R.6227 - National Quantum Initiative Act**  
115th Congress (2017-2018)

## Unification with ML models



Experiments

Computation

ML?

# JARVIS-DFT, FF and ML datasets and tools

>30000 bulk, 900 monolayer materials



Article | OPEN | Published: 12 July 2017

## High-throughput Identification and Characterization of Two-dimensional Materials using Density functional theory

Kamal Choudhary, Irina Kalish, Ryan Beams & Francesca Tavazza

Scientific Reports 7, Article number: 5179 (2017) | Download Citation

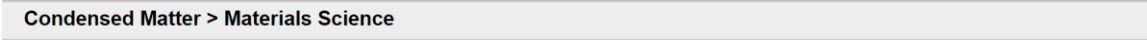


Data Descriptor | OPEN | Published: 08 May 2018

## Computational screening of high-performance optoelectronic materials using OptB88vdW and TB-mBJ formalisms

Kamal Choudhary, Qin Zhang, Andrew C.E. Reid, Sugata Chowdhury, Nhan Van Nguyen, Zachary Trautt, Marcus W. Newrock, Faical Yannick Congo & Francesca Tavazza

Scientific Data 5, Article number: 180082 (2018) | Download Citation



## High-throughput discovery of topological materials using spin-orbit spillage

Kamal Choudhary, Kevin F. Garrity, Francesca Tavazza

(Submitted on 24 Oct 2018)



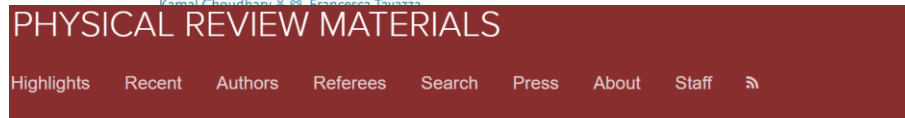
Computational Materials Science

Volume 161, 15 April 2019, Pages 300-308



## Convergence and machine learning predictions of Monkhorst-Pack k-points and plane-wave cut-off in high-throughput DFT calculations

Kamal Choudhary, Irina Kalish, Ryan Beams & Francesca Tavazza



## Machine learning with force-field-inspired descriptors for materials: Fast screening and mapping energy landscape

Kamal Choudhary, Brian DeCost, and Francesca Tavazza  
Phys. Rev. Materials 2, 083801 – Published 3 August 2018

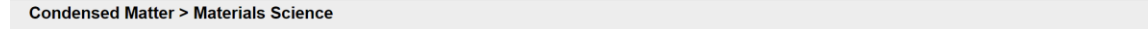
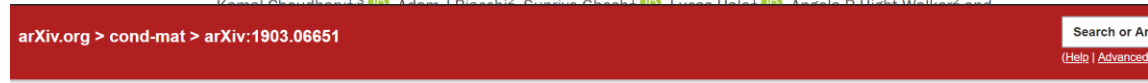


Journal of Physics: Condensed Matter

PAPER

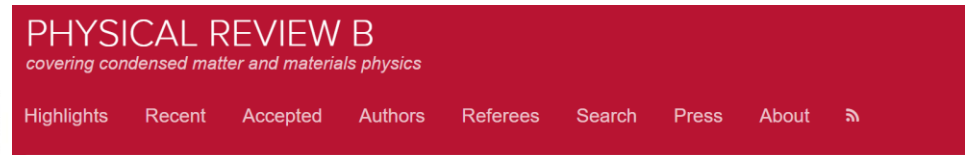
## High-throughput assessment of vacancy formation and surface energies of materials using classical force-fields

Kamal Choudhary, Brian DeCost, Sugata Chowdhury, Nhan Van Nguyen, Zachary Trautt, Marcus W. Newrock, Faical Yannick Congo & Francesca Tavazza



## Accelerated Discovery of Efficient Solar-cell Materials using Quantum and Machine-learning Methods

Kamal Choudhary, Marnik Berx, Jie Jiang, Ruth Pachter, Dirk Lamoen, Francesca Tavazza



## Elastic properties of bulk and low-dimensional materials using van der Waals density functional

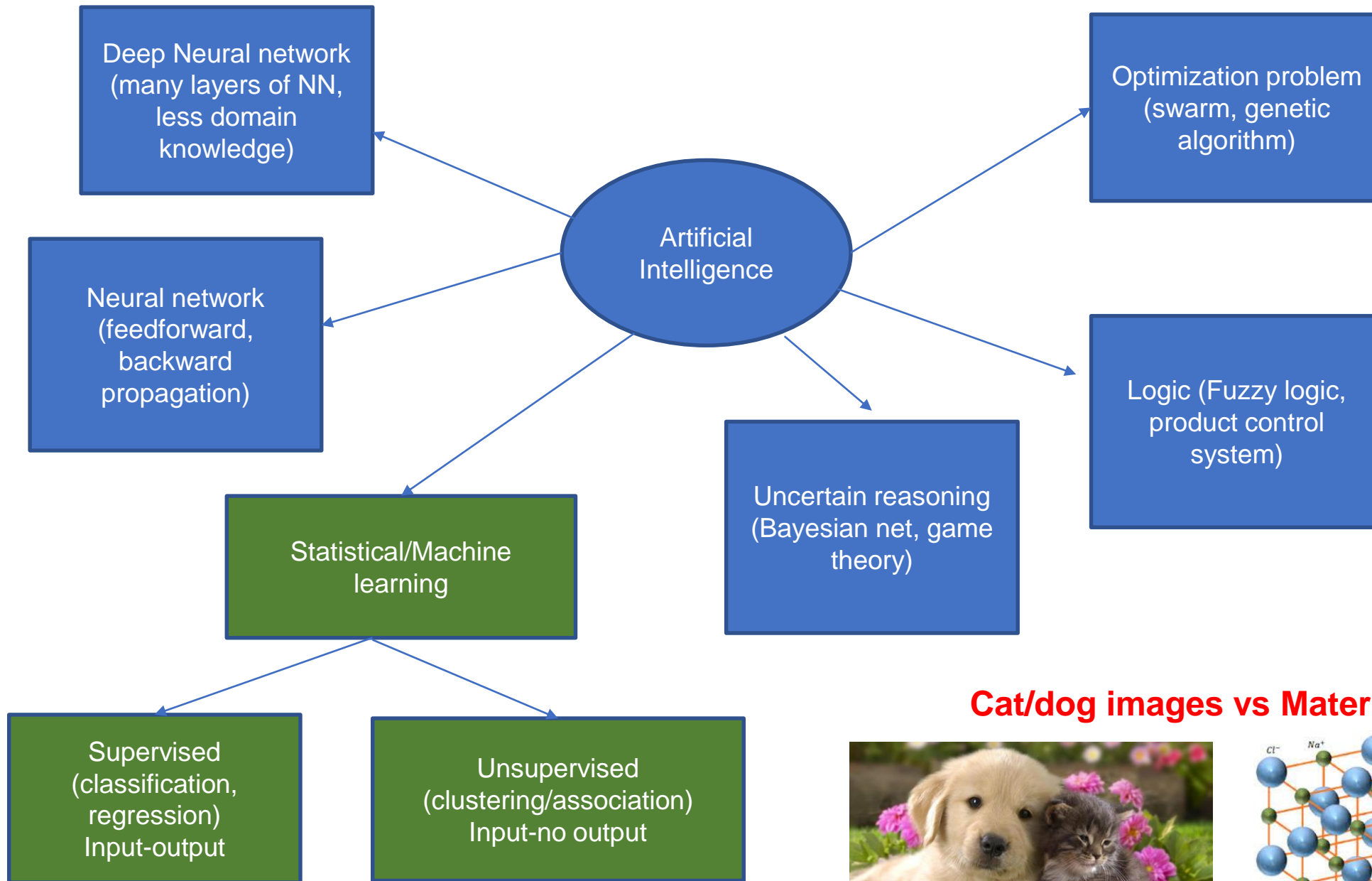
Kamal Choudhary, Irina Kalish, Ryan Beams & Francesca Tavazza



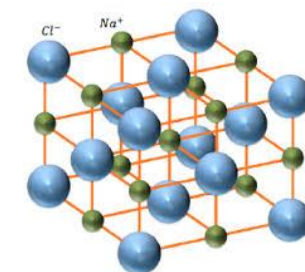
Data Descriptor | OPEN | Published: 31 January 2017

## Evaluation and comparison of classical interatomic potentials through a user-friendly interactive web-interface

Kamal Choudhary, Faical Yannick P. Congo, Tao Liang, Chandler Becker, Richard G. Hennig & Francesca Tavazza

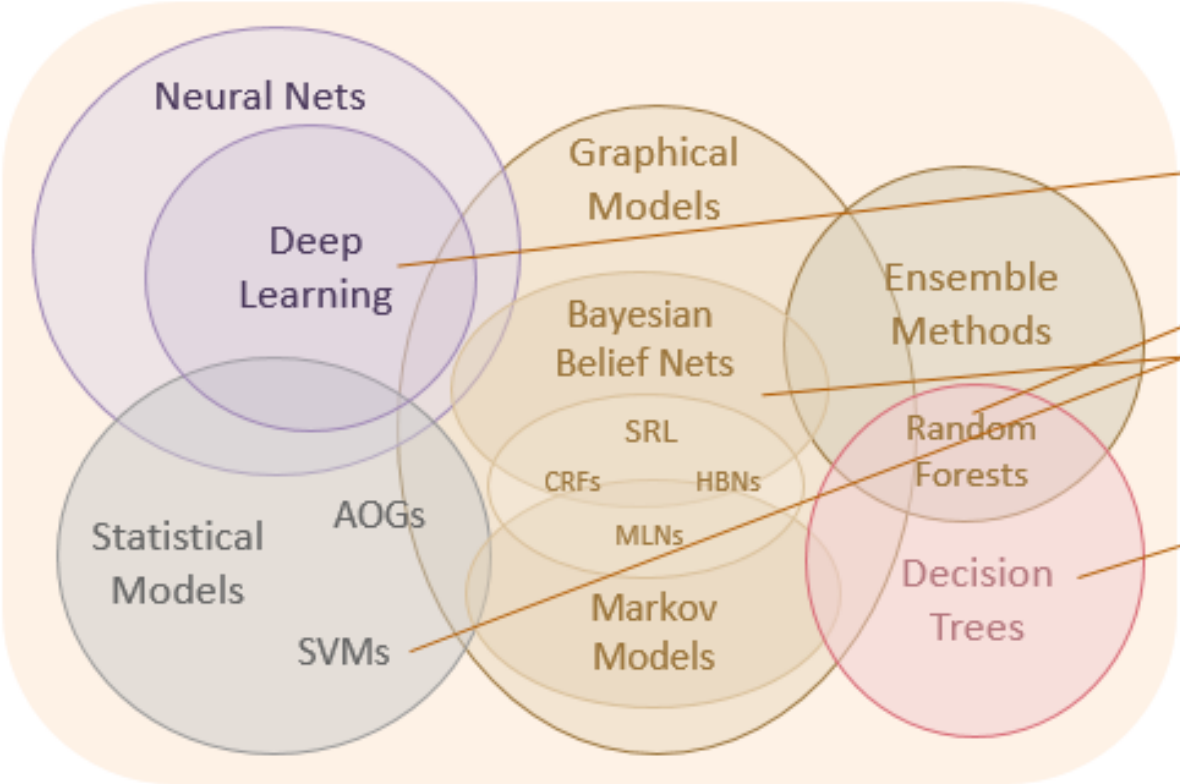


**Cat/dog images vs Materials data**

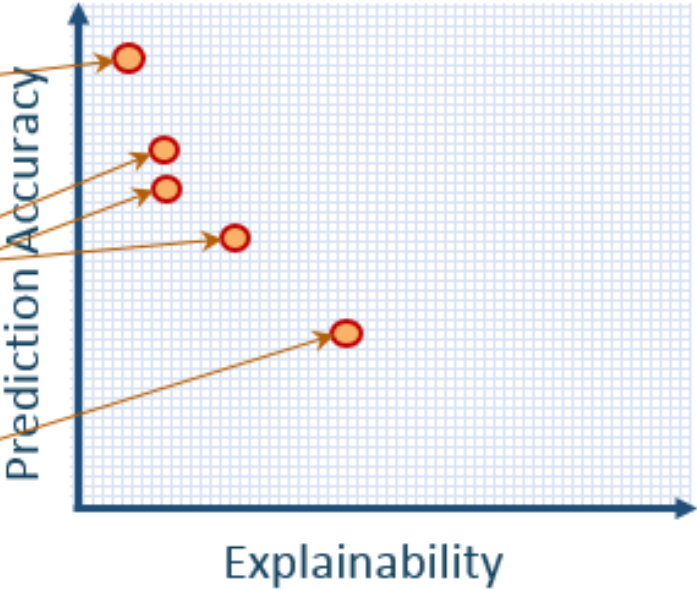


# Explainable AI

Learning Techniques (today)



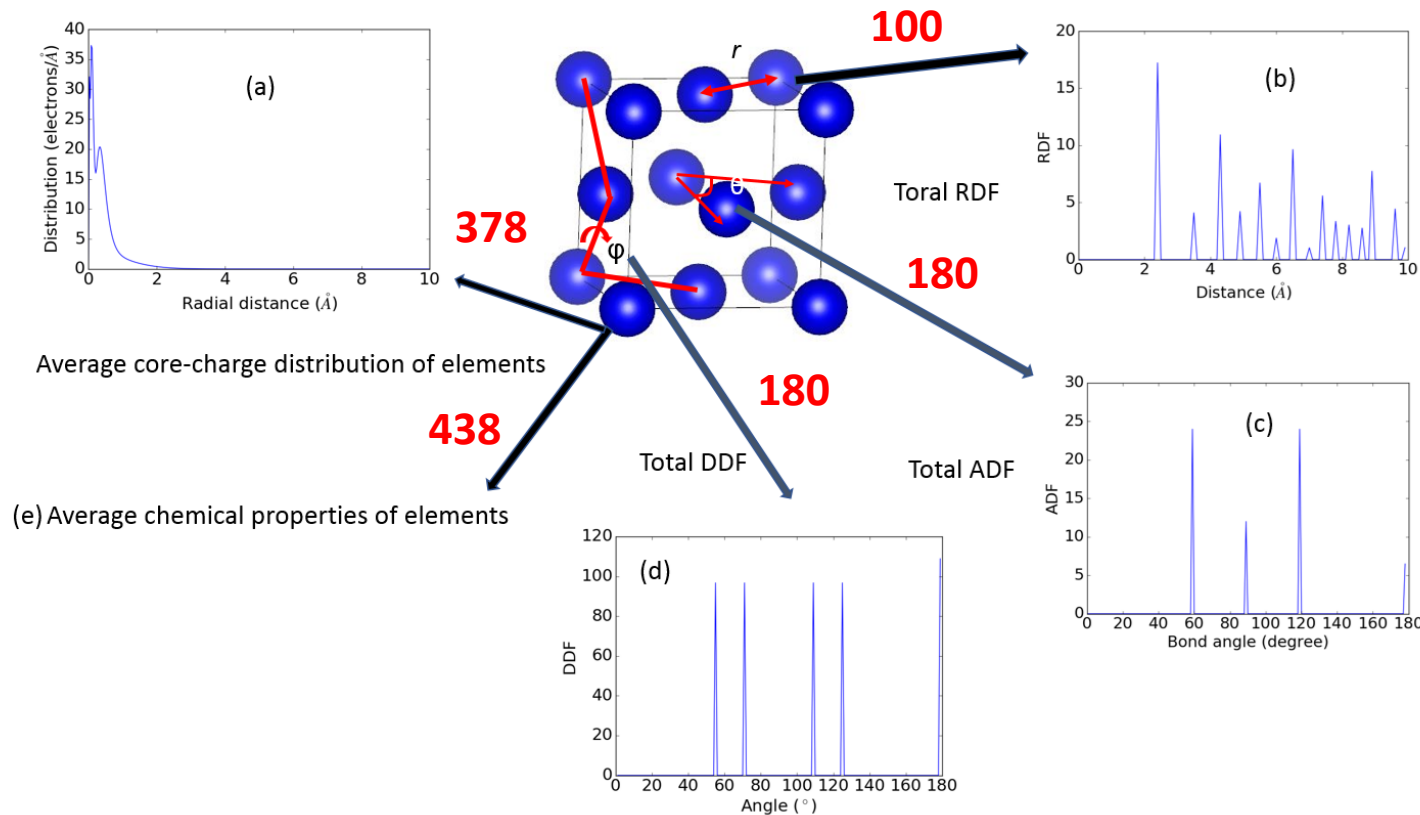
Explainability (notional)



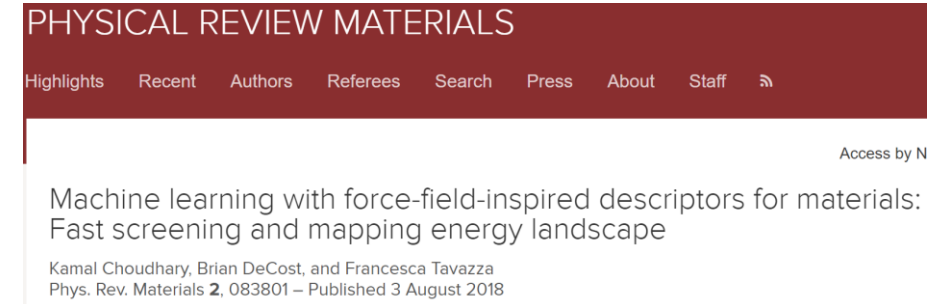


# CFID descriptors

1557 descriptors/features for one material



- Classical force-field inspired descriptors
- Arithmetic operations (mean, sum, std. deviation...) of **electronegativity, atomic radii, heat of fusion,....** of atoms at each site  
(example:  $\text{Electronegativity of Mo+Mo+S+S+S+S}/6 = 0.15$ )
- Atomic bond distance based descriptors
- Angle based descriptors

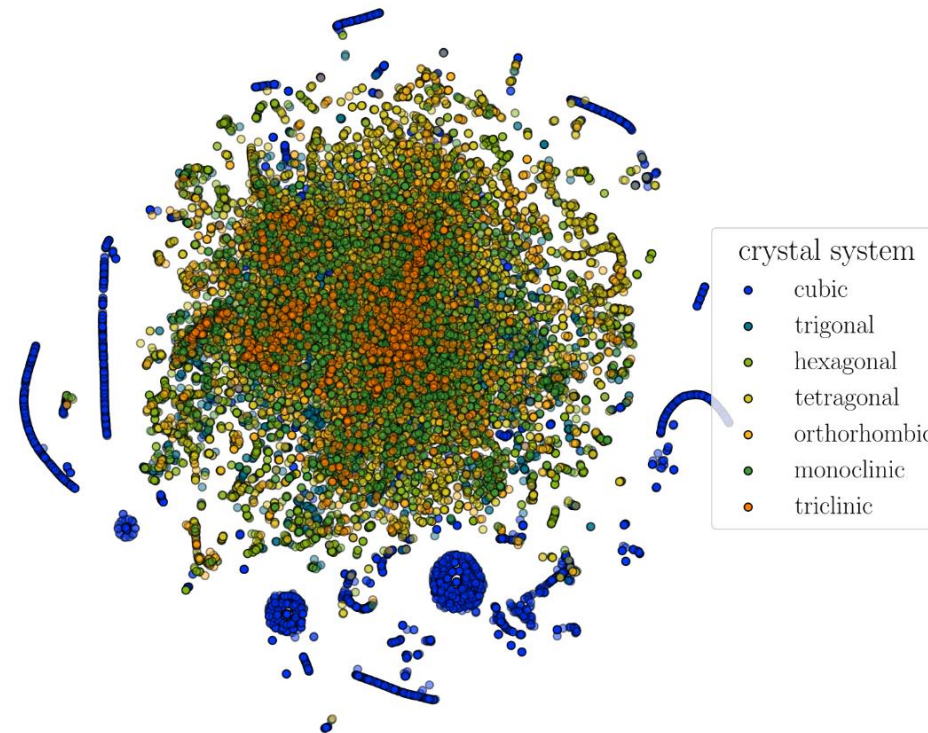
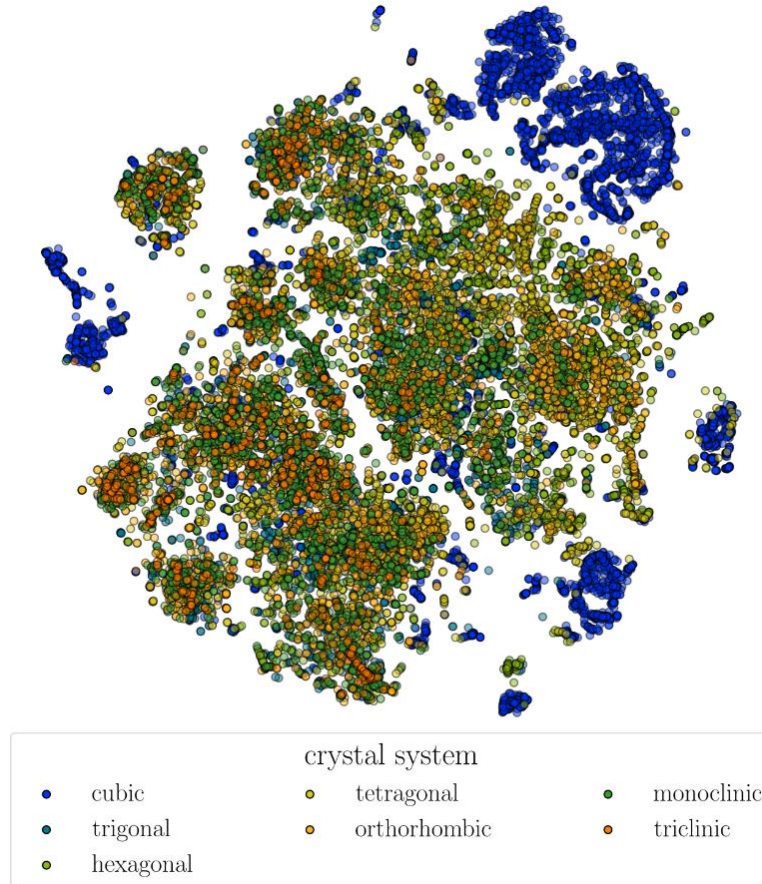


1.5 % unary, 26% binary, 56 % ternary, 13 % quaternary, 2 % quinary and 1% senary compounds, 1-96 atoms

<https://github.com/usnistgov/jarvis>

<https://hackingmaterials.github.io/matminer/index.html>

# Visualizing multi-dimensional data with t-SNE

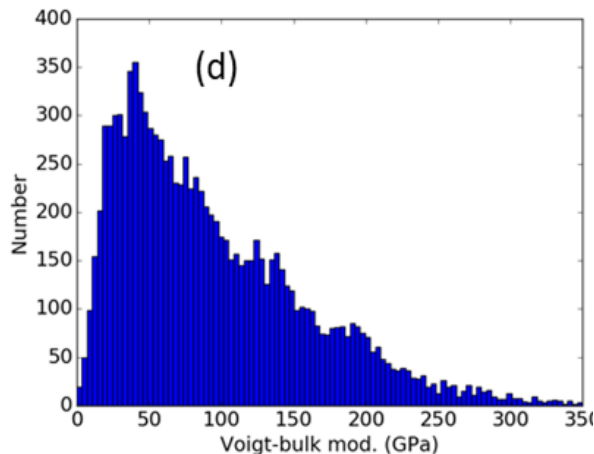
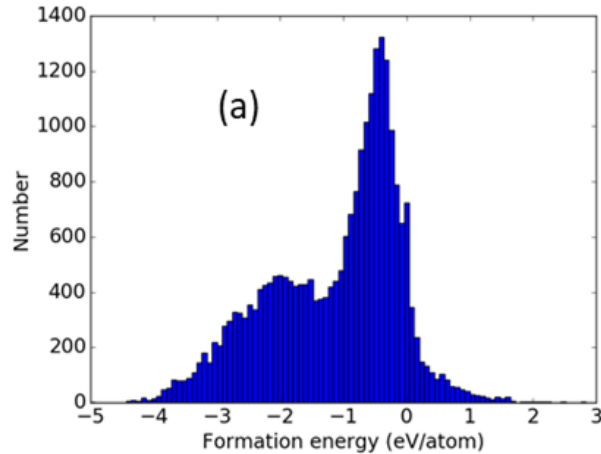


- Converts similarities between data points to joint probabilities
- Visualization with t-SNE for ~25000 materials

<http://holoviews.org/>

<http://scikit-learn.org/stable/modules/generated/sklearn.manifold.TSNE.html>

# Properties of interest & histogram plots

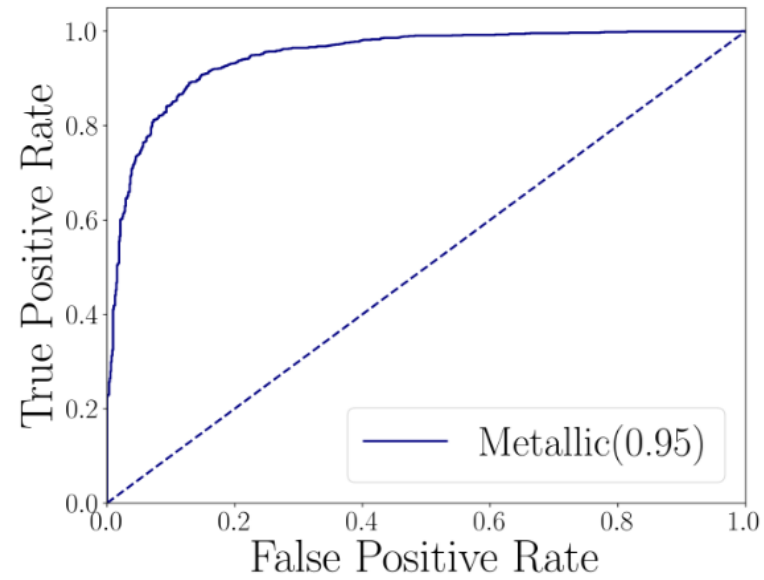
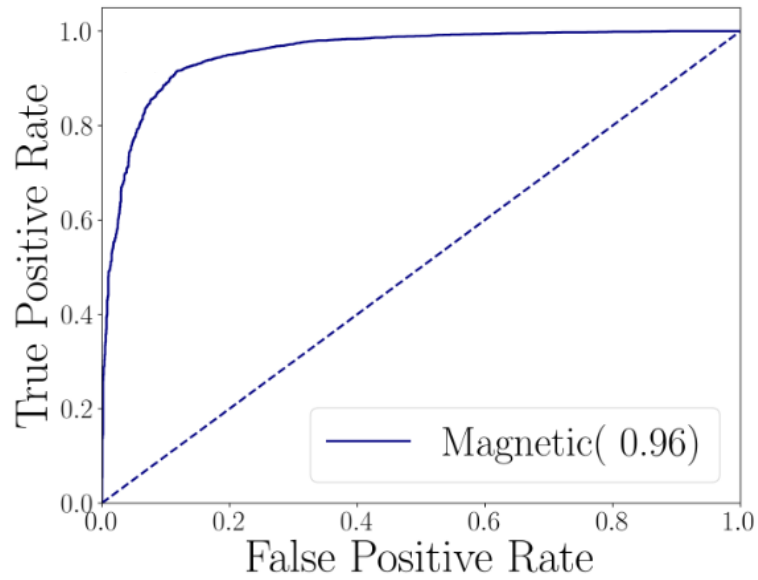


- Formation energy
- Bandgap
- Bulk/shear modulus
- **K-points, cut-off**
- **Thermoelectric metrics**
- **Solar-cell efficiency**
- **Refractive index**
- **2D Exfoliation energy**
- **Surface energy**
- **Grain boundary energy**
- **Topological spillage**

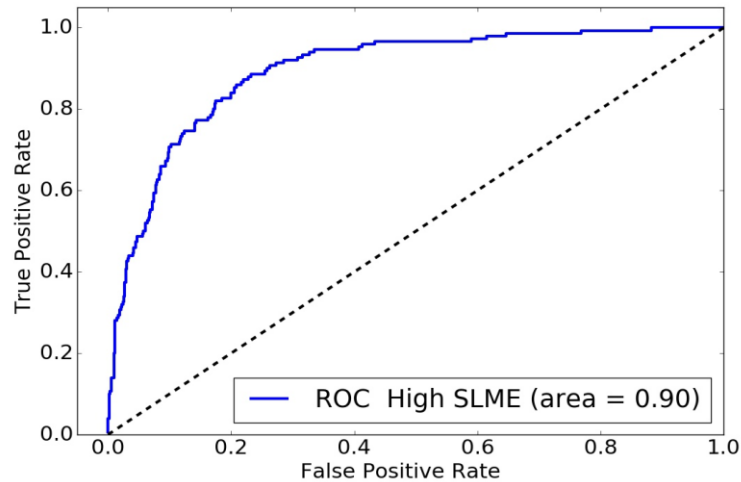
**New**



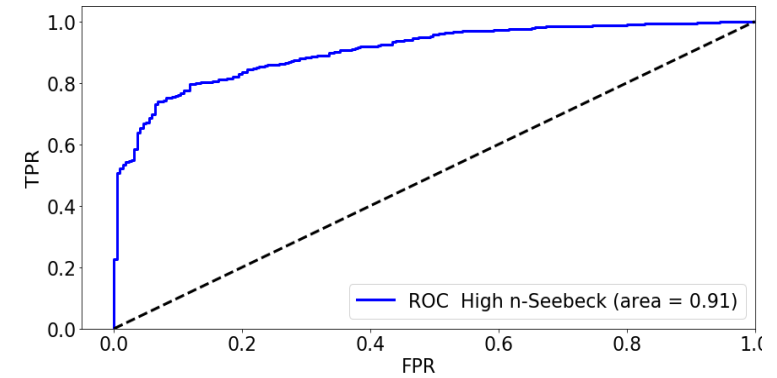
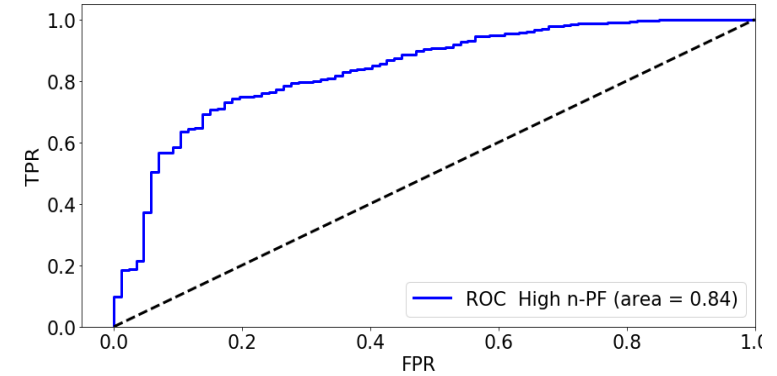
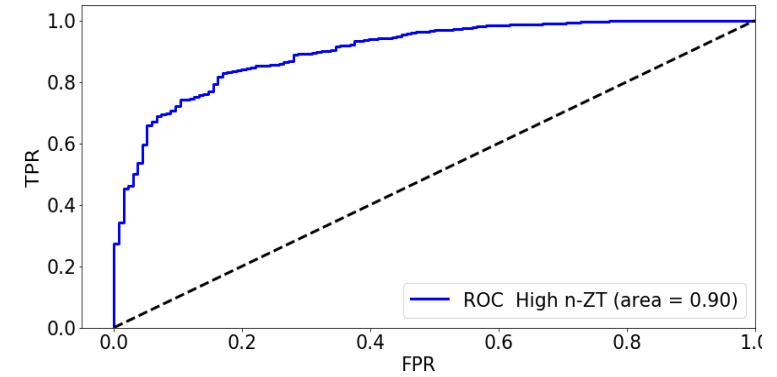
# Classification: ROC curves



Color cell efficiency

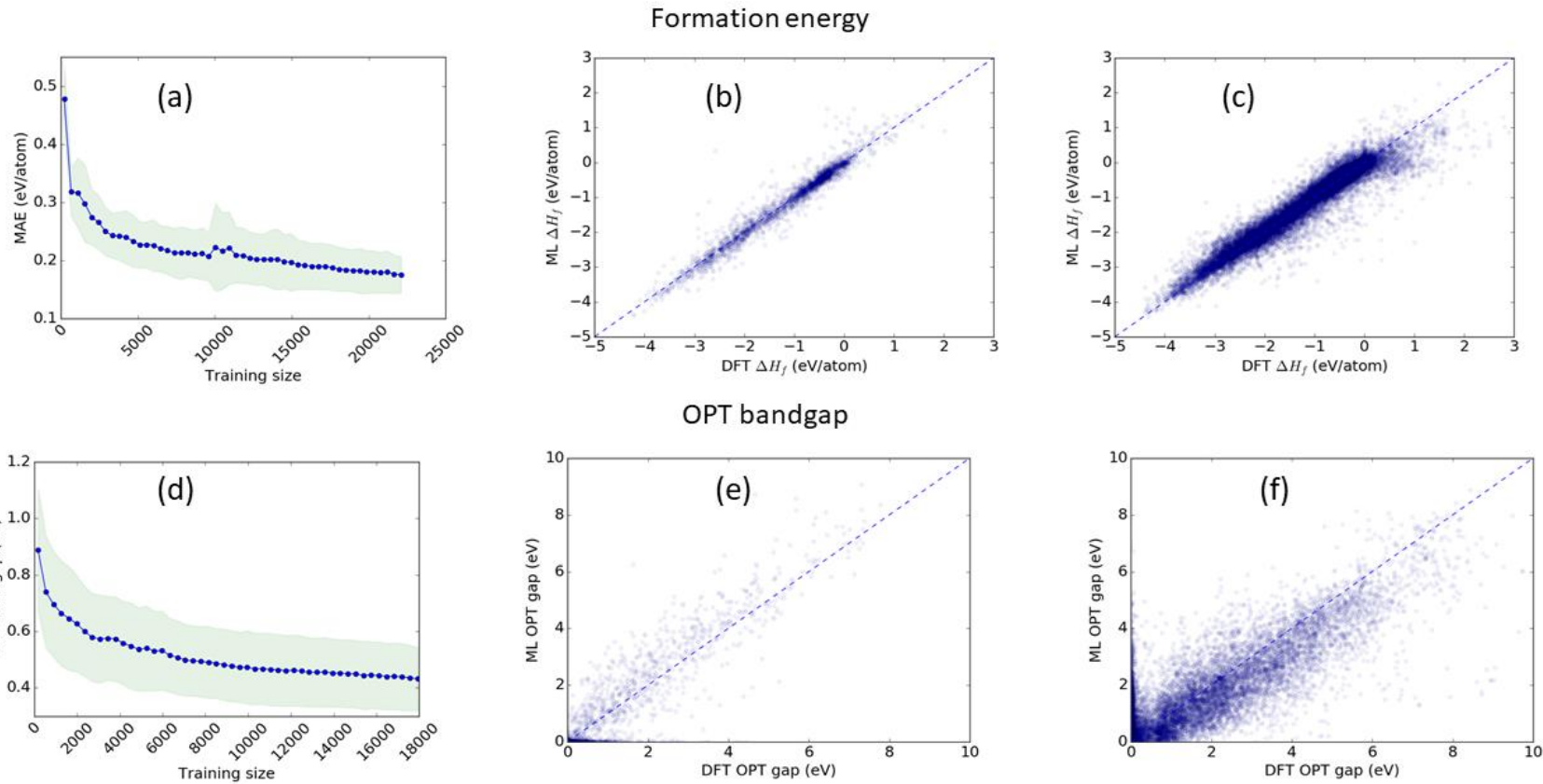


Thermoelectric



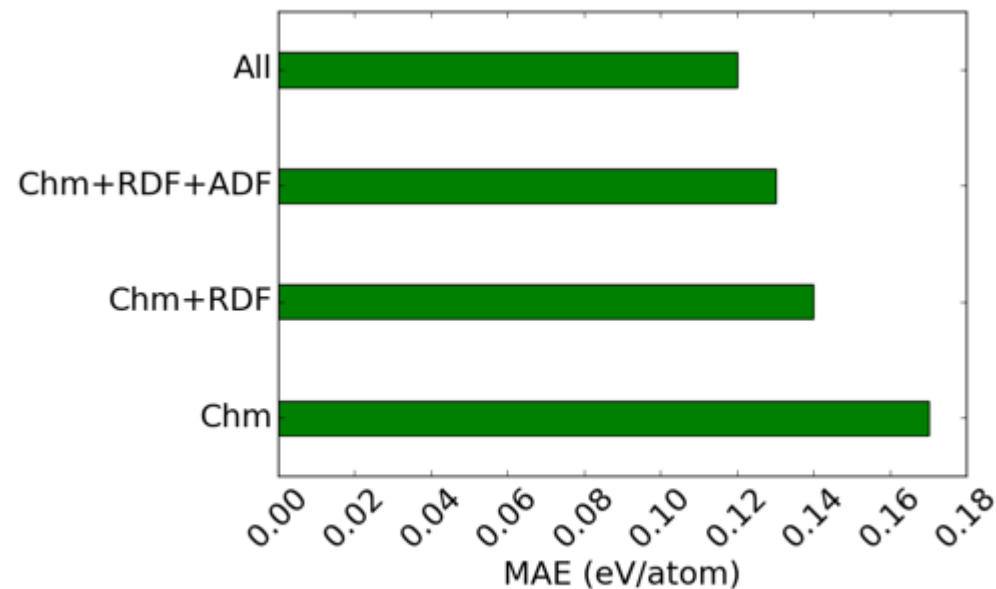
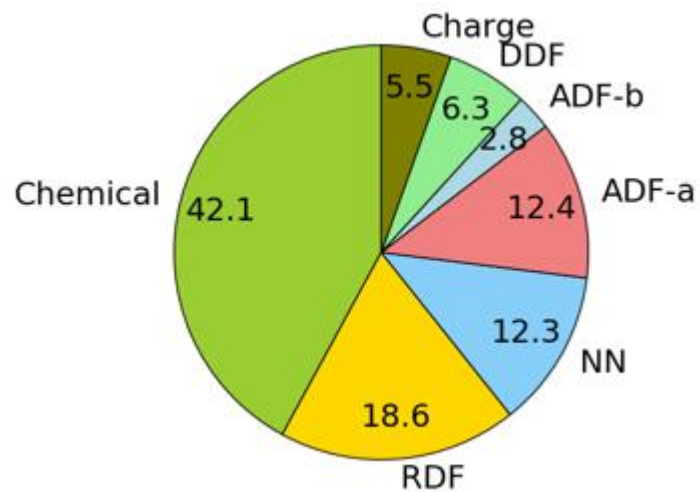
Perfect model area: 1  
Random guess: 0.5

# Regression models: formation energy and bandgap model



Learning curve shows scope of further improvement

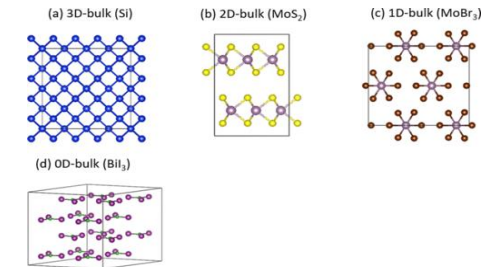
# Explainability: feature importance



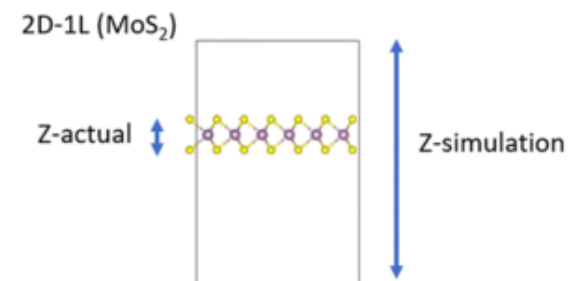
- Chemical features most important followed by RDF and NN
- Incrementally adding structural features decreases MAE

# Regression

Property	#Data-points	MAE <sub>CFID-DFT</sub>	MAE <sub>CFID-DFT</sub> (CV)	MAE <sub>DFT-Exp</sub>
Formation energy (eV/atom)	24549	0.12	0.17±0.05	0.136
OPT-bandgap (eV)	22404	0.32	0.37±0.24	1.33
MBJ-bandgap (eV)	10499	0.44	0.56±0.26	0.51
Bulk modulus (GPa)	10954	10.5	12.63±3.3	10.0
Shear modulus (GPa)	10954	9.5	11.55±3.15	10.0
OPT-n <sub>x</sub> (no unit)	12299	0.54	0.65±0.15	1.78
OPT-n <sub>y</sub> (no unit)	12299	0.55	0.65±0.16	-
OPT-n <sub>z</sub> (no unit)	12299	0.55	0.70±0.18	-
MBJ-n <sub>x</sub> (no unit)	6628	0.45	0.55±0.14	1.6
MBJ-n <sub>y</sub> (no unit)	6628	0.50	0.51±0.15	-
MBJ-n <sub>z</sub> (no unit)	6628	0.46	0.54±0.14	-
Exfoliation energy (meV/atom)	616	37.3	60.13±10.41	-

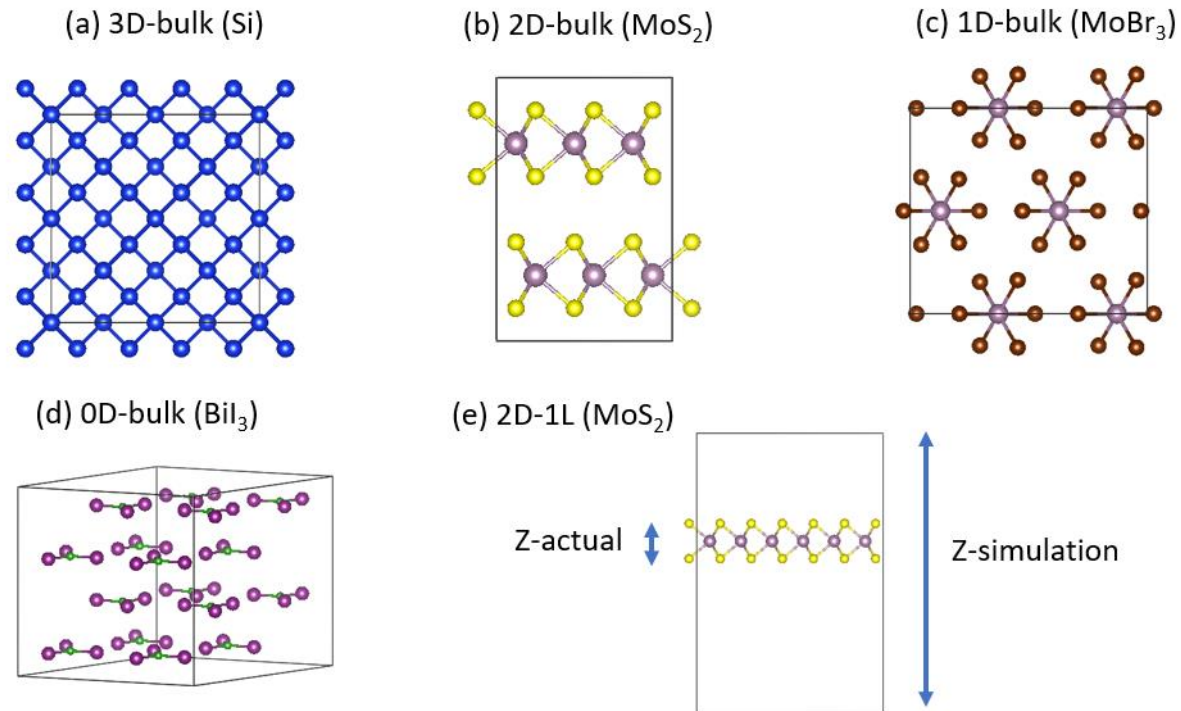


**3D**



**2D**

# 2D materials-screening example



- ~5000 2D materials predicted
- Requires expensive DFT calculations for predicting properties such as bandgap, exfoliation energy etc.
- Use of ML drops down the time to a few seconds
- Using this technique we identified new 2D materials such as CuI, InS etc.
- Validated using DFT

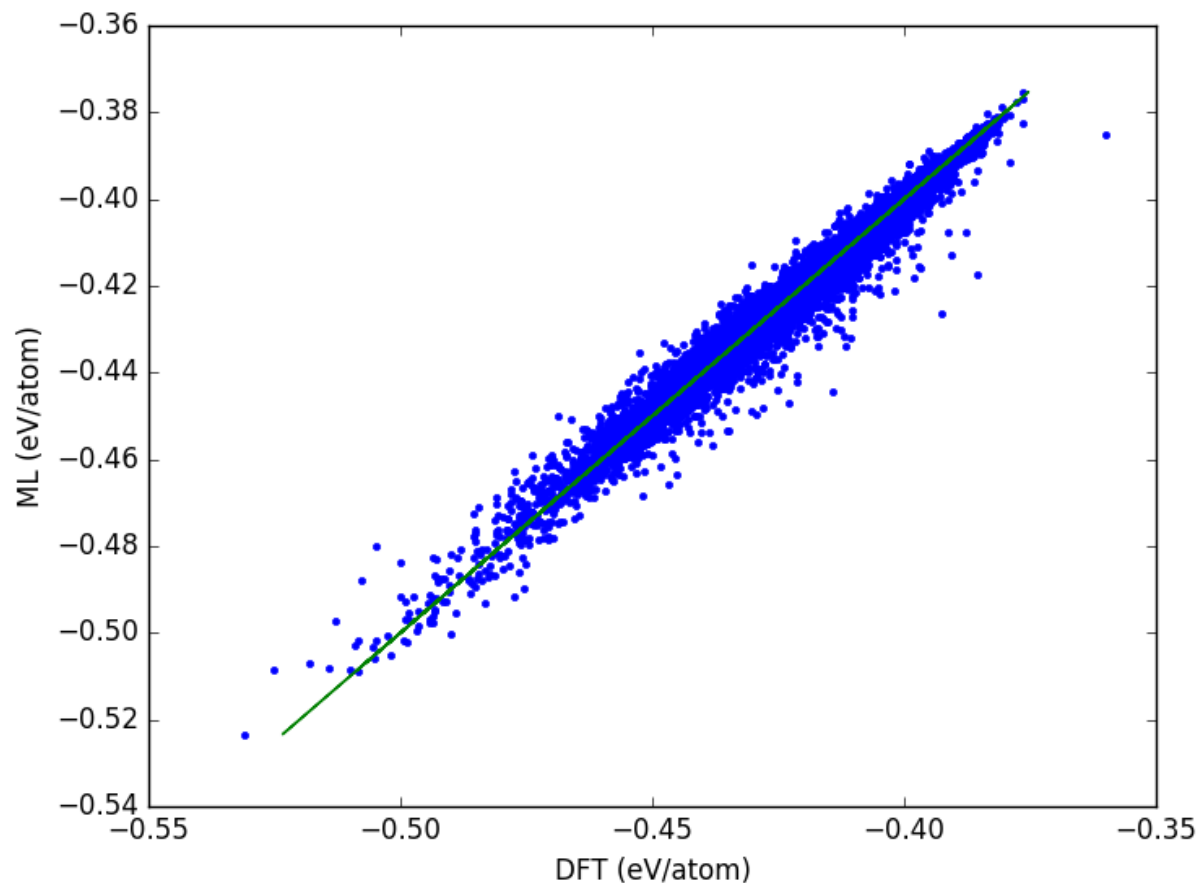


# Molecules

MAE internal energy (eV/atom): 0.002

$r^2$ : 0.97

(On-going work)




Email: kamal.choudhary@nist.gov

MENU ▾

SCIENTIFIC DATA 

Data Descriptor | [OPEN](#) | Published: 05 August 2014

## Quantum chemistry structures and properties of 134 kilo molecules

Raghunathan Ramakrishnan, Pavlo O. Dral, Matthias Rupp & O. Anatole von Lilienfeld 

*Scientific Data* **1**, Article number: 140022 (2014) | [Download Citation](#) 

### Abstract

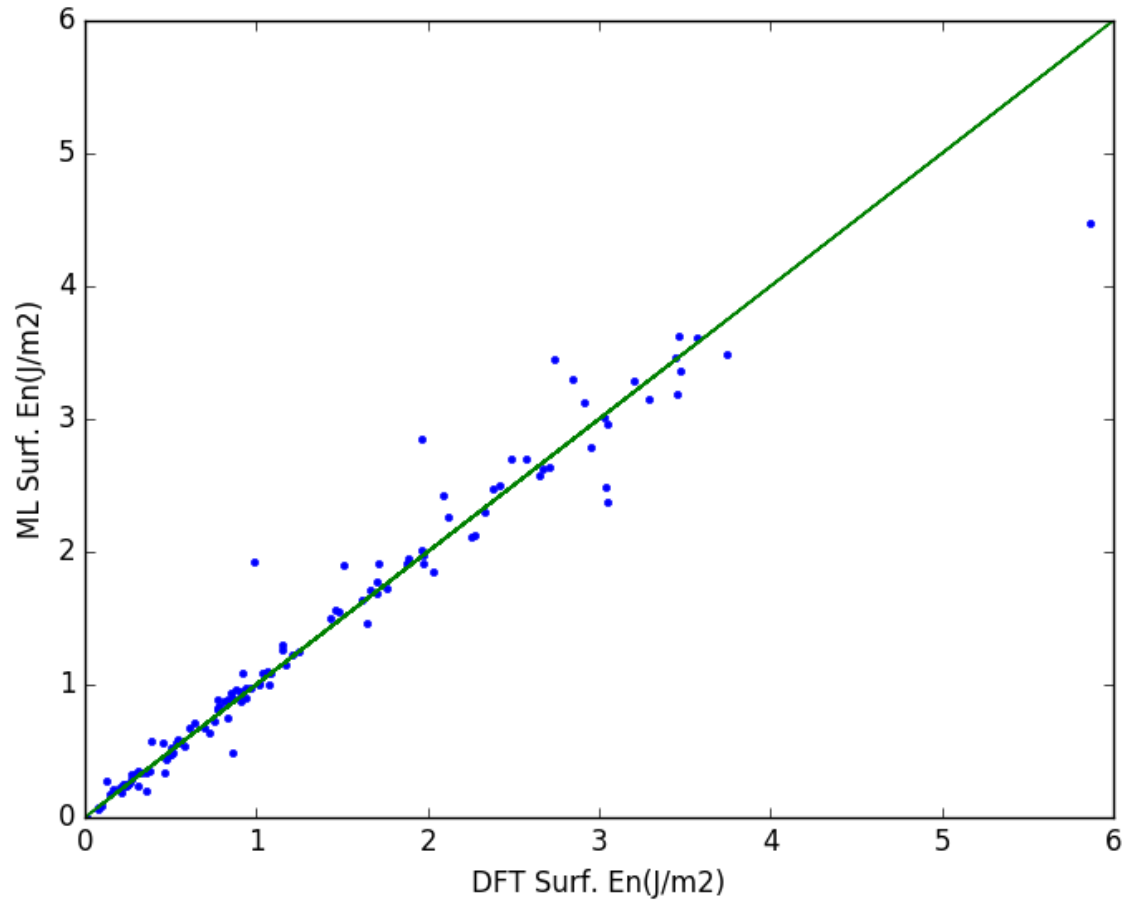
Computational *de novo* design of new drugs and materials requires rigorous and unbiased exploration of chemical compound space. However, large uncharted territories persist due to its size scaling combinatorially with molecular size. We report computed geometric, energetic, electronic, and thermodynamic properties for 134k stable small organic molecules made up of CHONF. These molecules correspond to the subset of all 133,885 species with up to nine heavy

# Surfaces

MAE: 0.13 J/m<sup>2</sup>

r<sup>2</sup>:0.94

(On-going work)



Data Descriptor | [OPEN](#) | Published: 13 September 2016

## Surface energies of elemental crystals

Richard Tran, Zihan Xu, Balachandran Radhakrishnan, Donald Winston, Wenhao Sun, Kristin A. Persson & Shyue Ping Ong

*Scientific Data* **3**, Article number: 160080 (2016) | [Download Citation](#)

### Abstract

The surface energy is a fundamental property of the different facets of a crystal that is crucial to the understanding of various phenomena like surface segregation, roughening, catalytic activity, and the crystal's equilibrium shape. Such surface phenomena are especially important at the nanoscale, where the large surface area to volume ratios lead to properties that are significantly different from the bulk. In this work, we present the largest database of calculated surface energies for elemental crystals to date. This database contains the surface energies of more than 100 polymorphs of about 70 elements, up to a maximum Miller index of two and three for non-cubic and cubic crystals, respectively. Well-known reconstruction schemes are also accounted for. The database is systematically improvable and has been rigorously

# Grain boundaries

MAE: 0.04 J/m<sup>2</sup>

r<sup>2</sup>:0.98

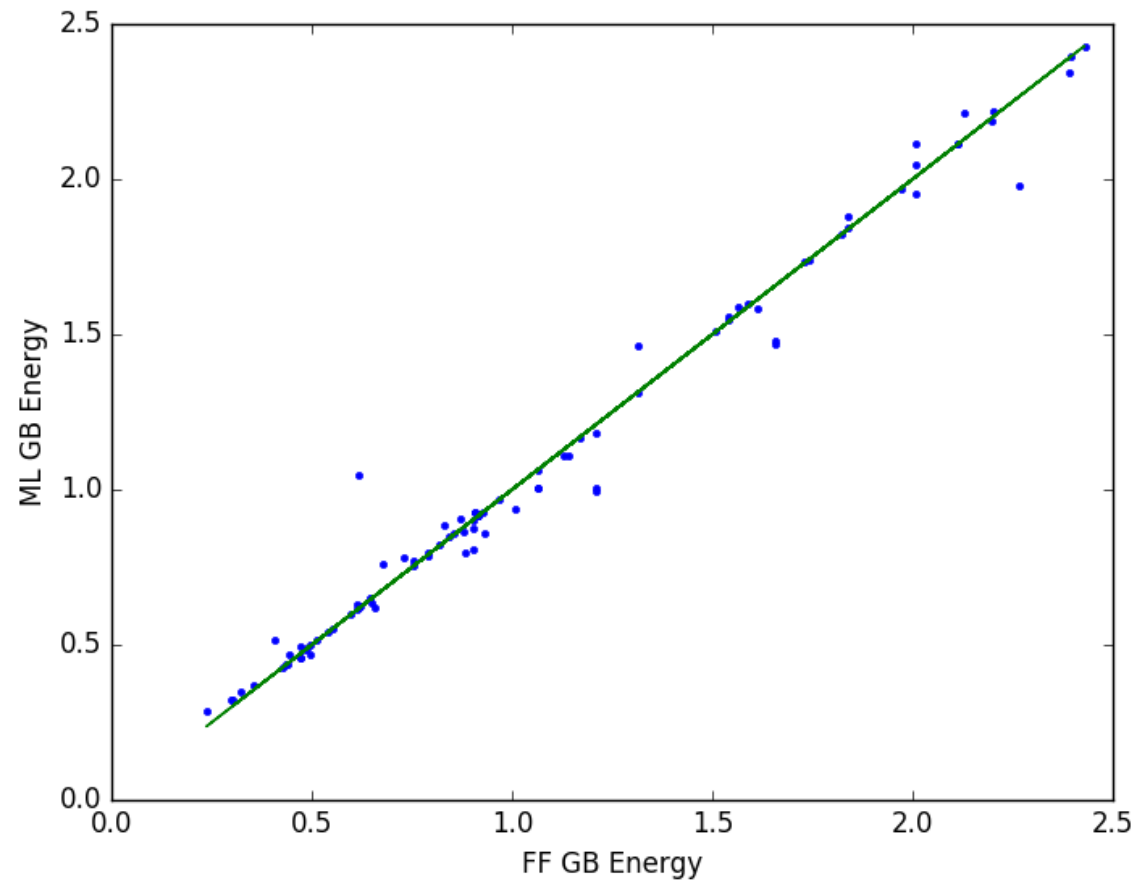
(On-going work)

**Symmetric tilt GBs**

FCC: Al, Ni, Cu, Ag, Au, Pd, Pt

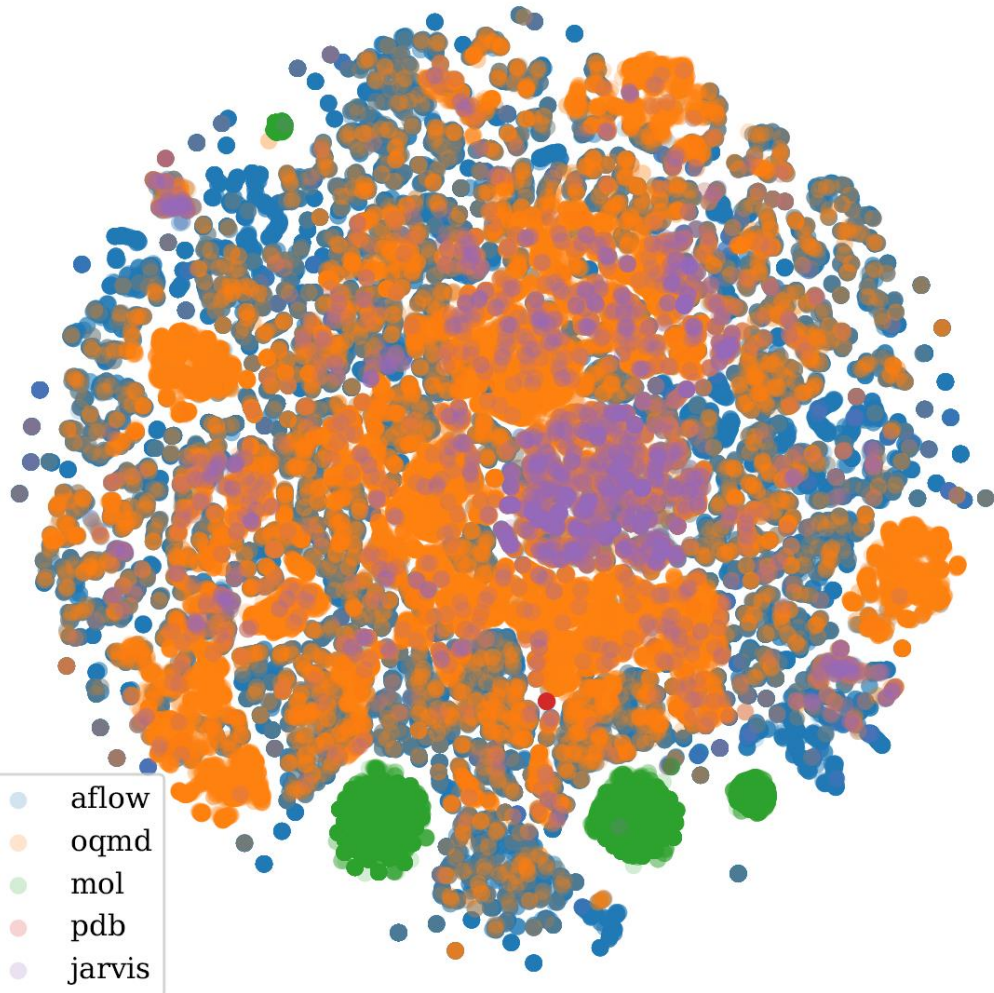
BCC: Fe, W, Ta, Mo

Diamond: Si



# Proteins

t-SNE visualization  
(On-going work)



>10000 proteins  
>630000 AFLOW  
>360000 OQMD  
>111000 COD  
>820000 MP  
> 30000 JV

The screenshot shows the top navigation bar of the RCSB PDB website. It includes a search bar with the text "Search by PDB ID, author, macromolecule, sequence, or ligands" and a "Go" button. Below the search bar are several logos for partner organizations: PDB-101, Worldwide Protein Data Bank, EMDataResource, Nucleic Acid Database, and Worldwide Protein Data Bank Foundation. The main header text reads "RCSB PDB PROTEIN DATA BANK" and "149886 Biological Macromolecular Structures Enabling Breakthroughs in Research and Education".

The screenshot shows the left sidebar of the RCSB PDB website. It contains a vertical list of navigation options: "Welcome", "Deposit", "Search", "Visualize", "Analyze", "Download", and "Learn". Each option is accompanied by a small icon representing its function.

## A Structural View of Biology

This resource is powered by the Protein Data Bank archive-information about the 3D shapes of proteins, nucleic acids, and complex assemblies that helps students and researchers understand all aspects of biomedicine and agriculture, from protein synthesis to health and disease.

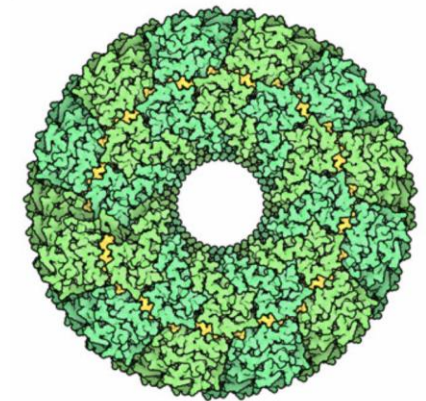
As a member of the wwPDB, the RCSB PDB curates and annotates PDB data.

The RCSB PDB builds upon the data by creating tools and resources for research and education in molecular biology, structural biology, computational biology, and beyond.

### Superbugs! How Bacteria Evolve Resistance to Antibiotics



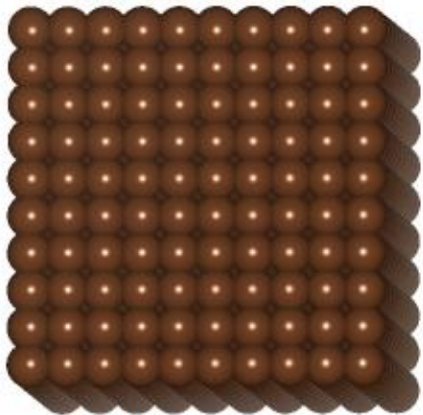
## March Molecule of the Month



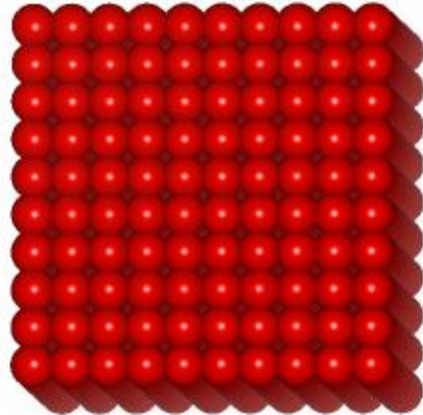
Measles Virus Proteins



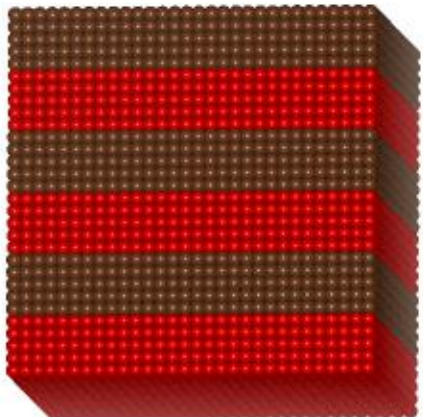
# Stringent validation: Genetic Algorithm with ML



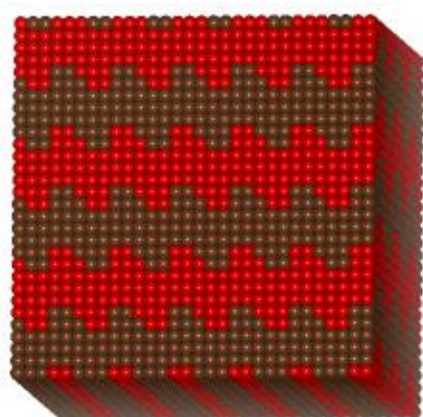
(a) Artificial cubic parent structure.



(b) Artificial cubic parent structure.



(c) Child created by the slicing crossover using a horizontal cut.



(d) Child created by the slicing crossover using a periodic cut.

- Based on **'Survival of the fittest'** theory: fitness of crystal structure based on energy of structure
- Parents to offspring crystal structure
- Generally energy is obtained from DFT, MD...let's try ML ?

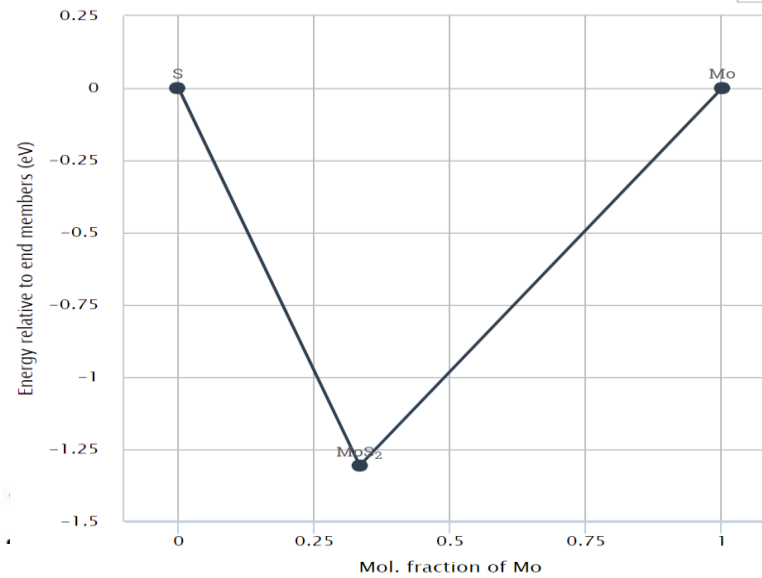
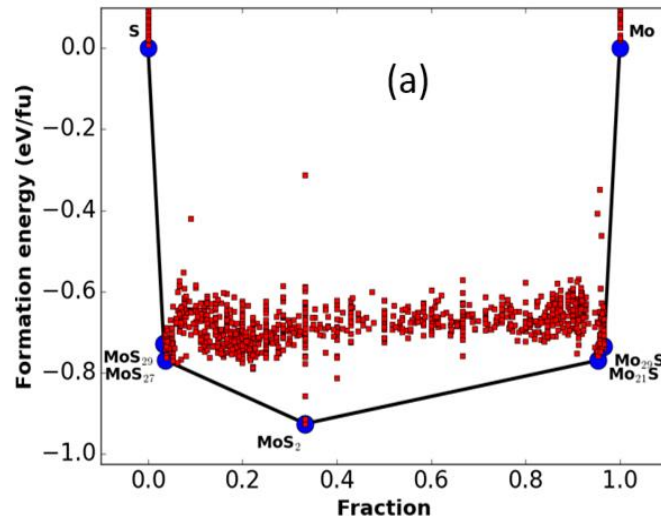
D. M. Deaven, Molecular geometry optimization with a genetic algorithm, Physical Review Letters, 75 (1995)

G. Ceder, Data-mining-driven quantum mechanics for the prediction of structure, MRS Bulletin, 31 (2006)

<https://github.com/henniggroup/GASP-python/>



# Genetic algorithm with ML

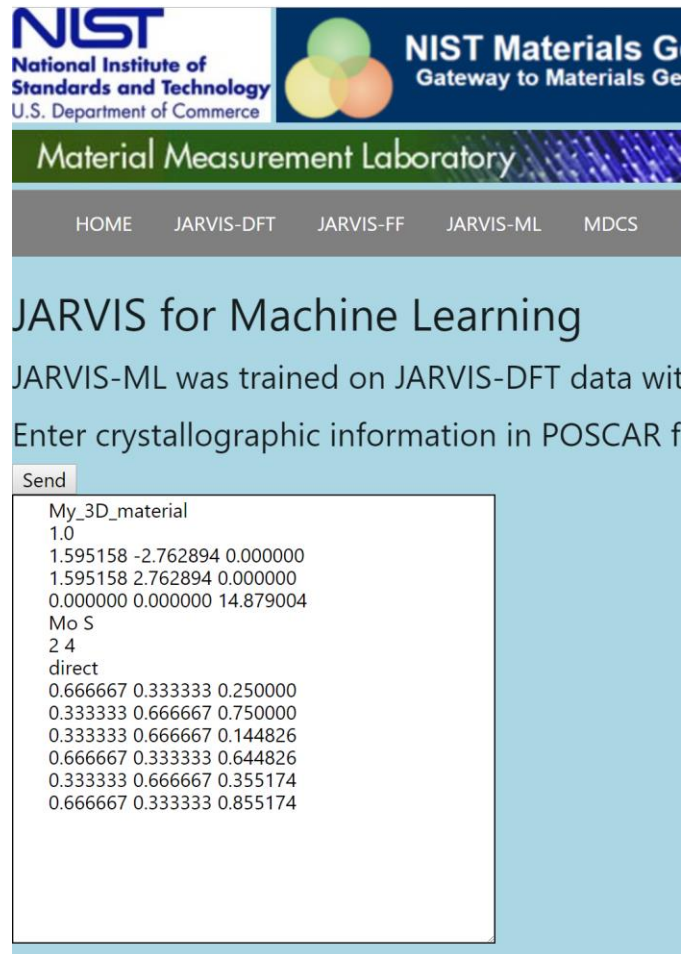


- **New way** of validating ML model for materials
- MoS<sub>2</sub>, WS<sub>2</sub> indeed stable as in DFT and experiments
- Need further verification for low-lying energy structures with DFT

# Web-app: DEMO

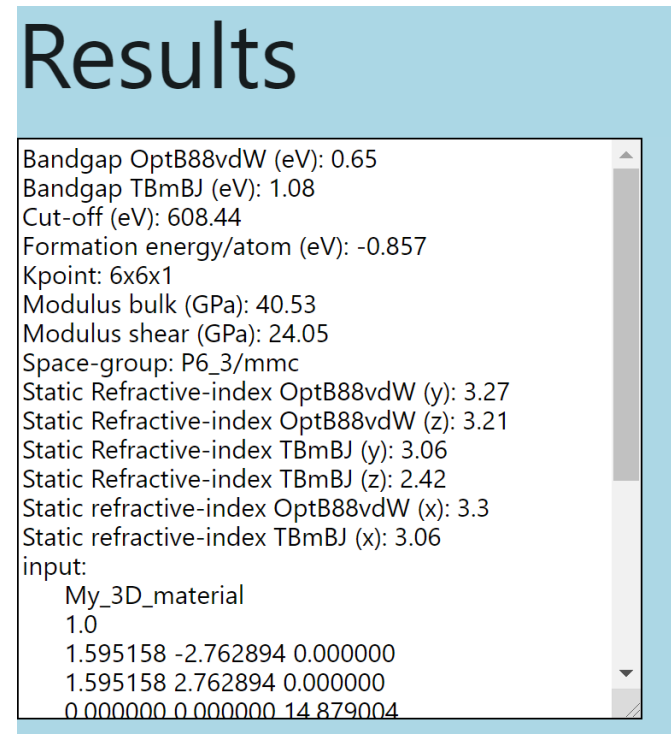
- **How to use the model?**

<https://www.ctcms.nist.gov/jarvisml/>



The screenshot shows the NIST Materials Gateway website. At the top, there is the NIST logo and the text "National Institute of Standards and Technology, U.S. Department of Commerce". To the right is the "NIST Materials Gateway to Materials Ge" logo. Below this is a navigation bar with "HOME", "JARVIS-DFT", "JARVIS-FF", "JARVIS-ML", and "MDCS". The main heading is "JARVIS for Machine Learning". Below the heading, it says "JARVIS-ML was trained on JARVIS-DFT data with" and "Enter crystallographic information in POSCAR format". There is a "Send" button and a text area containing the following POSCAR file content:

```
My_3D_material
1.0
1.595158 -2.762894 0.000000
1.595158 2.762894 0.000000
0.000000 0.000000 14.879004
Mo S
2 4
direct
0.666667 0.333333 0.250000
0.333333 0.666667 0.750000
0.333333 0.666667 0.144826
0.666667 0.333333 0.644826
0.333333 0.666667 0.355174
0.666667 0.333333 0.855174
```



The screenshot shows the "Results" page of the web application. It lists the following properties:

- Bandgap OptB88vdW (eV): 0.65
- Bandgap TBmBJ (eV): 1.08
- Cut-off (eV): 608.44
- Formation energy/atom (eV): -0.857
- Kpoint: 6x6x1
- Modulus bulk (GPa): 40.53
- Modulus shear (GPa): 24.05
- Space-group: P6\_3/mmc
- Static Refractive-index OptB88vdW (y): 3.27
- Static Refractive-index OptB88vdW (z): 3.21
- Static Refractive-index TBmBJ (y): 3.06
- Static Refractive-index TBmBJ (z): 2.42
- Static refractive-index OptB88vdW (x): 3.3
- Static refractive-index TBmBJ (x): 3.06

input:

```
My_3D_material
1.0
1.595158 -2.762894 0.000000
1.595158 2.762894 0.000000
0.000000 0.000000 14.879004
```

# Summary

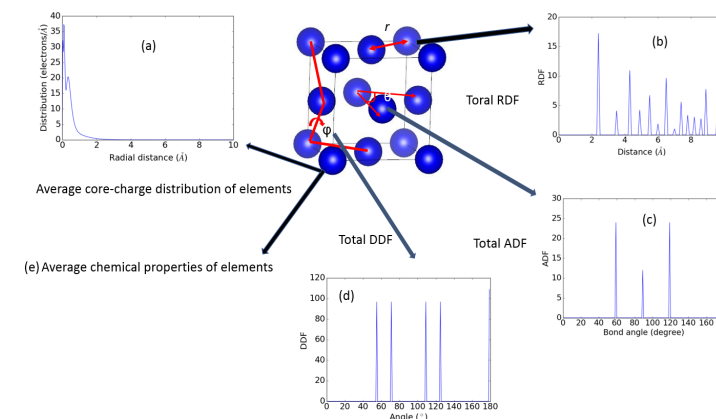


- Unified machine learning descriptors for various classes of materials
- All the code and data publicly available
- Formation energy convex hull test, beyond data-science metric
- Web-app for on-the fly prediction of properties
- **AIMS workshop: August 1-2, 2019, Registration open**
- **More data and tools on the way**
- **Important links:**

✓ <https://jarvis.nist.gov/>

✓ <https://github.com/usnistgov/jarvis>

✓ Slides available at: <https://www.slideshare.net/KAMALCHOUDHARY4/>



## Thank you for your time!