

# Class Overview

---

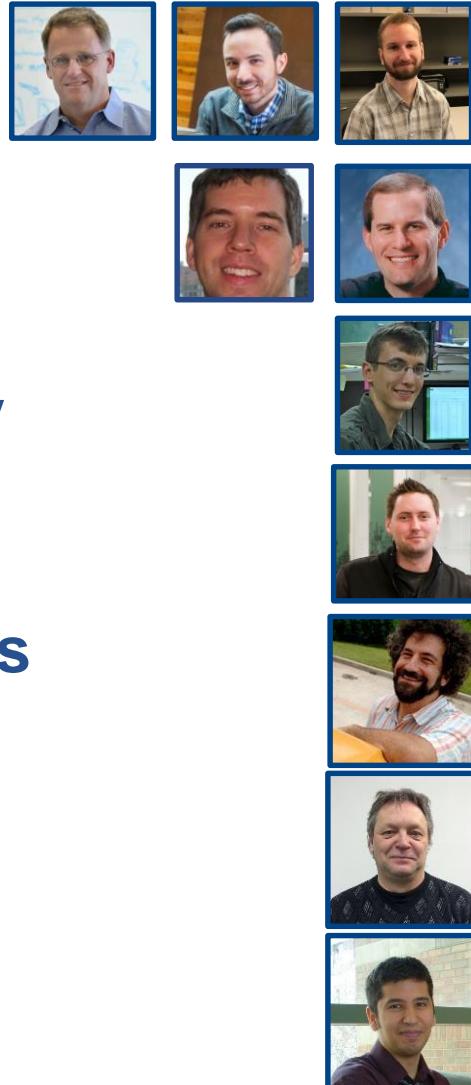
- **Materials Data Facility**
  - Data publication
  - Data discovery
  - Data service integration
  - Automation
- **Overview of ML/DL in Materials Science**
- **Model Serving with DLHub**
- **Lots of time for questions, so please ask!**

# MDF Overview

## Build data services to

- Empower researchers to publish data, regardless of size, type, and location
- Automate data and metadata ingest, to enable capture of many valuable materials datasets
- Enable unified search and discovery across disparate materials data sources

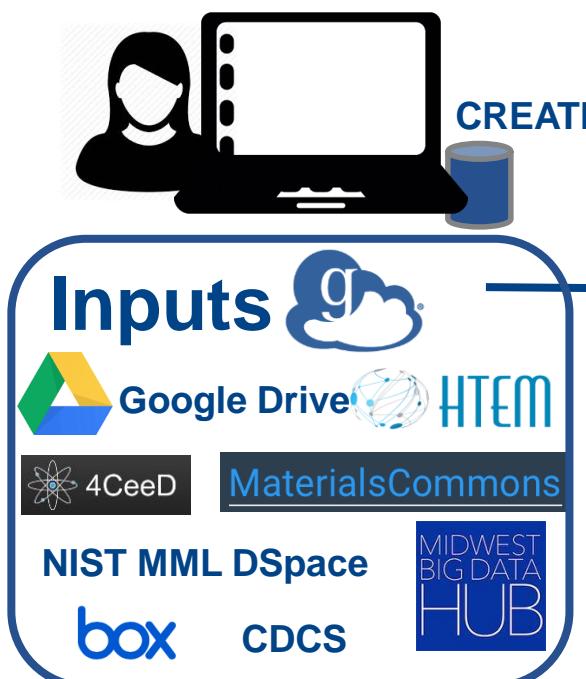
Deploy with APIs to simplify connection to other data efforts and to enable automation



# MDF Connect - Connecting Community Services

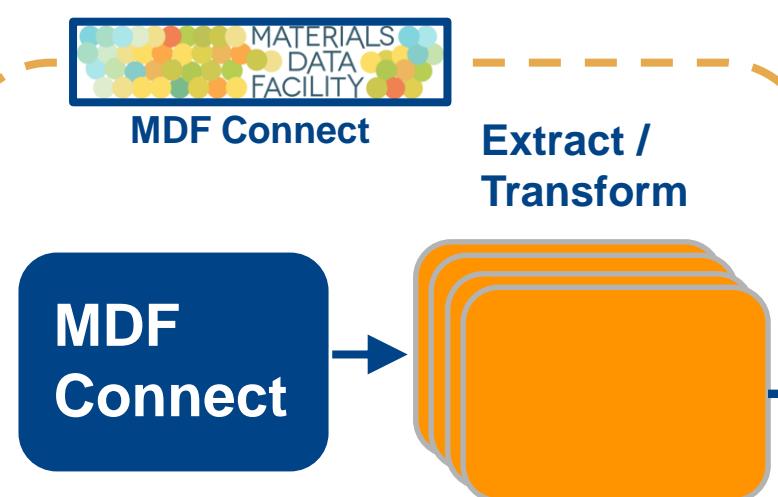
- Make it easy to deposit into many services from one location
- Strictly opt-in for cross-posting datasets

## Submit Data [UI or API]



## Enrich Data

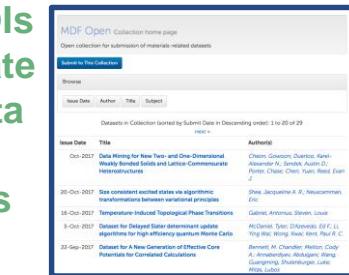
Extract / Transform



## Send to Community

- Mint DOIs
- Associate metadata
- Persist datasets

In progress → NanoMine



- Query
- Browse
- Aggregate



NIST  
MRR



CHIMaD NIST

...

# MDF Connect Prototype

The screenshot shows a web browser window for 'MDF Connect - Home'. The top navigation bar includes 'Transfer Files | Globus' and a 'LOGIN' button. The main content area features the 'MATERIALS DATA FACILITY' logo, the title 'MDF CONNECT', a subtext about sharing data, a call-to-action button 'Become a Contributor', and a 'HOW TO GET STARTED' section with three circular icons.

MDF Connect - Home

Transfer Files | Globus

LOGIN

MATERIALS DATA FACILITY

## MDF CONNECT

It has never been easier to share your data with the community. Deposit data once, send to partner services.

Tell your research story.

Become a Contributor

### HOW TO GET STARTED



# MDF Connect Alpha

## MDF Connect Alpha invites going out next week

- Functionality includes dataset indexing (MDF), dataset publication (MDF), MRR registration, Citrine submission
- Send data from Globus, Google Drive, or HTTP using Python or Web UI

The screenshot shows the MDF Connect homepage. At the top right, there is a 'LOGOUT' link, an email address 'BLAISZIK@GLOBUSID.ORG', and a profile icon. Below this, the 'MATERIALS DATA FACILITY' logo is displayed. The main heading 'MDF CONNECT' is centered. Below it, a sub-headline reads 'It has never been easier to share your data with the community. Deposit data once, send to partner services.' A call-to-action button labeled 'Become a Contributor' is visible at the bottom left.

The screenshot shows the 'ADD A NEW DATASET' form. At the top right, there is a 'LOGOUT' link, an email address 'BLAISZIK@GLOBUSID.ORG', and a profile icon. The form fields include:

- Title \***: My New Dataset
- Authors \***: Ben Blaiszik (with an 'add author' link)
- Institutions**: University of Chicago, Argonne National Laboratory (with an 'add institution' link)
- Data Locations \***: Two URLs: [https://www.globus.org/app/transfer?origin\\_id=e38ee745-6d04-11e5-ba46-22000b92c6ec&origin\\_path=%2F](https://www.globus.org/app/transfer?origin_id=e38ee745-6d04-11e5-ba46-22000b92c6ec&origin_path=%2F) and [https://www.globus.org/app/transfer?origin\\_id=e38ee745-6d04-11e5-ba46-22000b92c6ec&origin\\_path=%2FBlaiszik%2F](https://www.globus.org/app/transfer?origin_id=e38ee745-6d04-11e5-ba46-22000b92c6ec&origin_path=%2FBlaiszik%2F) (with an 'add location' link)
- Tags**: DFT, defect, diffusion (with an 'add tag' link)
- Description**: A large text area with placeholder text 'Add some tags to help users find the data'.
- Sync Data**: Checkboxes for 'MDF Publish', 'Materials Resource Registry', and 'Citrination'.

A 'Submit' button is located at the bottom left of the form area.

# MDF Connect Alpha

## MDF Connect Alpha invites going out next week

- Functionality includes dataset indexing (MDF), dataset publication (MDF), MRR registration, Citrine submission
- Send data from Globus, Google Drive, or HTTP using Python or Web UI

```
pip install mdf_connect_client
```

```
from mdf_connect_client import MDFConnectClient

mdf = MDFConnectClient()

authors = ["Blaiszik, Ben", "Jonathon Gaff"]
affs = ["University of Chicago", "Argonne National Laboratory"]

mdf = MDFConnectClient(title="My Dataset Title",
                      authors=authors,
                      affiliations=affs)

mdf.add_data(["googledrive:///mydata.zip", "globus://1a2b3c/data/"])
mdf.add_services("mdf_publish", "mrr", "citrine")
mdf.submit()
```

```
mdf.create_mrr_block({"dataOrigin": "experiment"})
mdf.mrr
```

# MDF Connect Input Integration Highlights

## Google Drive (operational)

Drive

Share with others

Get shareable link

Link sharing on [Learn more](#)

Anyone with the link can edit

Copy link

<https://drive.google.com/drive/folders/0Bz9P8tnaPRx-VmtUT2psQ2J3UVU?usp=sharing>

People

materialsdatalab@gmail.com Add more people...

Add a note

Shared with Rachana Ananthakrishnan, Ian Foster and 12 others

Advanced

## Box (prototype)

box

Search Files and Folders

All Files Synced Trash Notifications Notes Admin Console Dev Console Favorites Drag items here for quick access

9xLwvaT.jpg is selected

Name Updated Size

Name	Updated	Size
DESJ223233.0-514630.9.tif.png	Aug 14, 2018 by B...	32.6 KB
DESJ232352.3-581526.1.tif.png	Aug 14, 2018 by B...	31.5 KB
DESJ033735.8-402739.1.tif.png	Aug 14	
9xLwvaT.jpg	Aug 10	

Share Upload New Version Download Favorite Move or Copy Lock Properties More Actions Integrations

Materials Data Facility Send to Chatter Send with DocuSign

## DropBox (operational)



## Figshare (operational)

w/ Ben Galewksky @ UIUC

# MDF Connect Input Integration Highlights

## 4CeeD

4CeeD You Shared Create Trash Help Uploaders Search Jupyter Session Logout

< Dashboard

### GaN Etch Pressure 3.5 mTorr

Created by Patrick Su  
Created on May 22, 2018  
Access:  Space Default (Private)  Private  
 Public

Add a description

+ Add Files Download All Files Delete Collaborators

Space containing the Dataset  
4CeeD Demo: GaN Etch Recipe Optimization Figure 2  
4 datasets | Remove

Copy Dataset to Spaces  
Select a Space + COPY

Collections containing the Dataset

Publish  
SELECT REPOSITORY Materials Data Facility Zenodo TAG

Files Metadata Comments (0)



4CeeD You Shared Create Trash Help Uploaders Search Jupyter Session Logout

### Create Data Citation

Title: GaN Etch Pressure 3.5 mTorr Required

Author: Patrick Su Required

Publication Year: 2018 Numeric

Description:

Related DOIs:

PUBLISH CANCEL

# Materials Data Discovery

Start your search here



- Data are often locked:
  - In archives (zip, tar, etc.) [e.g., Zenodo, NIST DSpace ]
  - Behind web forms [many research group resources]
  - Within clunky, partially functional APIs
- Data are located in many distributed locations
- Data file formats are heterogeneous
  - How would you generally search by material composition??

The screenshot shows a Zenodo dataset page. At the top, there's a navigation bar with 'zenodo' logo, 'Search' input, 'Upload', and 'Communities' buttons. Below the navigation is a timestamp 'January 13, 2014'. A 'Dataset' button and 'Open Access' badge are visible. The main content area describes the dataset as 'Elemental vacancy diffusion database from high-throughput first-principles calculations for fcc and hcp structures' by Angsten, Thomas; Mayeshiba, Tam; Wu, Henry; Morgan, Dane. It shows publication details ('Publication date: January 13, 2014'), a DOI link ('DOI 10.5281/zenodo.17888'), and a license ('Creative Commons Zero - CC0 1.0'). The file section lists four files: 'bulk\_modulus.tar.gz' (117.1 MB), 'fcc\_hvf\_hvm.tar.gz' (163.4 MB), 'figure\_excel\_files.zip' (52.1 kB), and 'hcp\_hvf\_hvm.tar.gz' (923.9 MB). Each file has a 'Download' button.

Name	Size	Action
bulk_modulus.tar.gz md5:8c6d6b493addbe2f052145aec269bd7b	117.1 MB	
fcc_hvf_hvm.tar.gz md5:152ba1410738e2acb78c3528803d558	163.4 MB	
figure_excel_files.zip md5:68fafe2a5471e947e9191a528fb4e20	52.1 kB	
hcp_hvf_hvm.tar.gz md5:ae042b4846d58c848cd826f2e0d95c37	923.9 MB	

# MDF Forge



<https://github.com/materials-data-facility>

## Forge

pypi v0.5.1 build passing coverage 85%

Forge is the Materials Data Facility Python package to interface and leverage the MDF Data Discovery service. Forge allows users to perform simple queries and facilitates moving and synthesizing results.

## Installation

```
pip install mdf_forge
```

Compile records for a larger, mixed source, result set

```
elements = ["Al"]
sources = ["khazana_vasp", "sluschi", "ab_initio_solute_database"]
my_ep = "c8ee7e5c-6d04-11e5-ba46-22000b92c6ec"
my_path = "/Users/ben/Desktop/blaiszik-macbookpro/dft_training_set"

mdf = Forge()
res = mdf.search_by_elements(elements=elements, sources=sources, limit=9999)
mdf.get_globus(res, dest=my_path,
               local_ep=my_ep, preserve_dir=True)

Processing records: 100%[██████████] 10/10 [00:00<00:00, 19.05it/s]
Submitting transfers: 100%[██████████] 1/1 [00:00<00:00, 3.60it/s]
All transfers submitted
Submission IDs: 3fbfc637-7181-11e7-a9fd-22000bf2d287
```



Materials Informatics  
Skunkworks (D. Morgan)



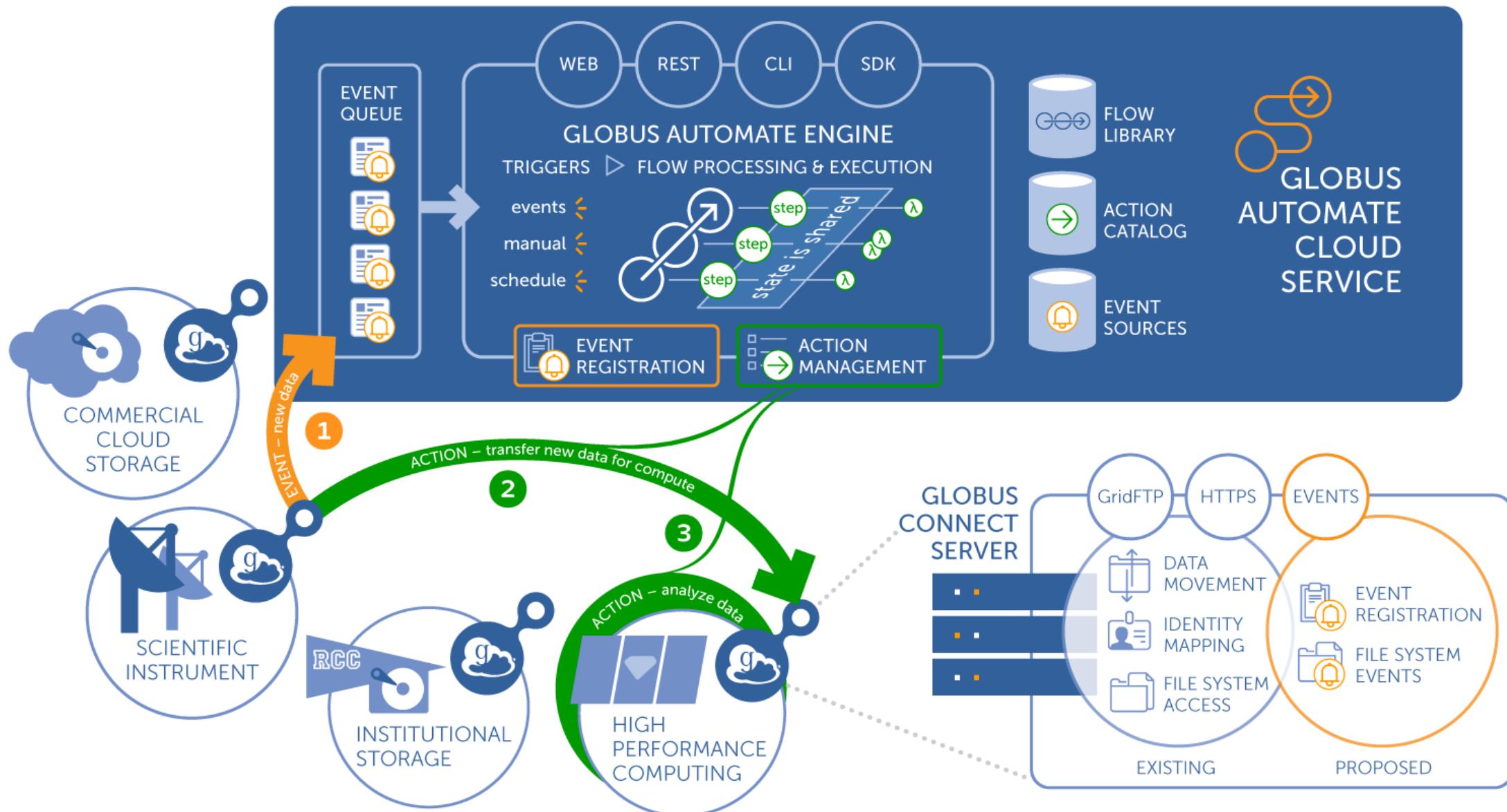
- Data mining and featurizing in materials science
- Integrated to pull data from MDF

(L. Ward, A. Jain)



NDS Labs Workbench

# Automation in Science



# Some Automation Actions

## Auth



User login

Secure service interactions

## Identify



Manage namespace

Mint DOI

## Search



Catalog

Query



## Transfer



Transfer data

Set permission

## Execute



Remote execution

Self optimization

## Describe



Gather metadata

Validate metadata

## Curate



Solicit input

Edit / Approve

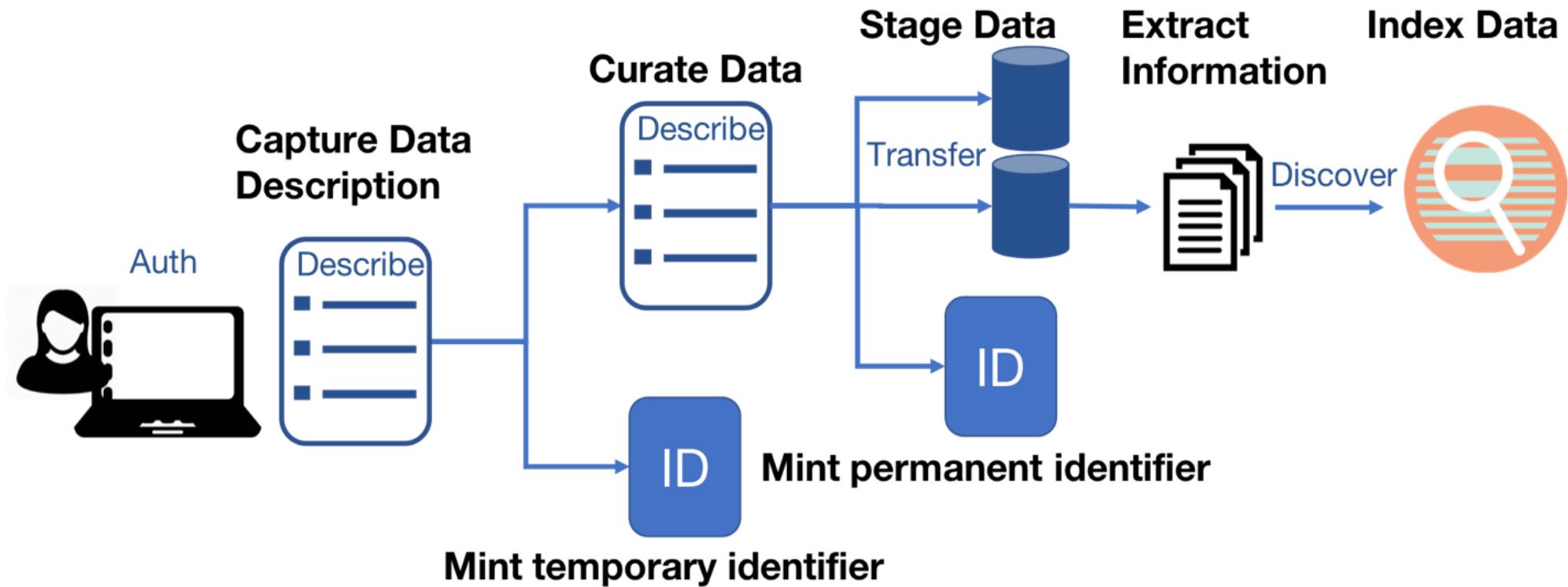
## Extract



Automate metadata extraction

Invoke container extraction

# Data Publication with Automate



# Deep Learning and Materials Design

Logan Ward

Postdoctoral Scholar

Data Science and Learning Division, ANL

Department of Computer Science, University of Chicago

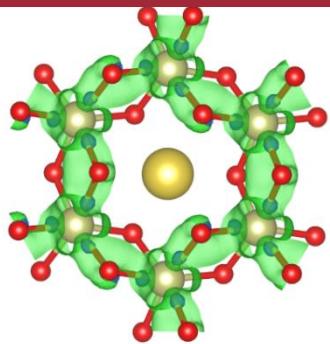
18 October 2018

# I'm a Materials Scientist

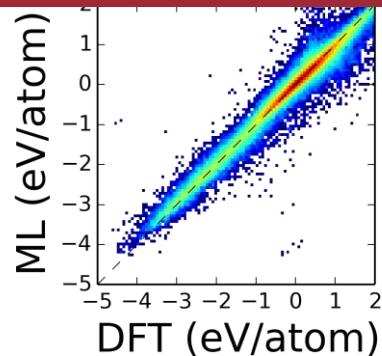
**Affiliation:** Post-doc with Ian Foster since Jan. 2017

**Background:** PhD in Materials Science, Northwestern  
BS/MS in Materials Science, from The Ohio State University

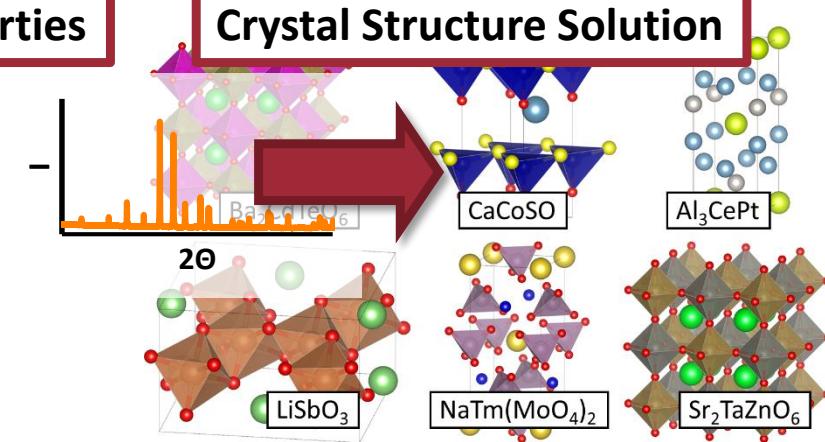
## Density Functional Theory



## ML of Material Properties



## Crystal Structure Solution



## Why UC/ANL?

1. Wanted to see ML in other fields
2. Build tools for others to use



# Materials Engineering: A narrow introduction

3

## Examples of “Aluminum”

### Application



### Key Properties

Conductivity

Low cost  
Shapeable

Castable  
Strong

???

### Alloy

Al + Mg

Al + Mg,Mn

Al + Si,Cu

???

How do we tailor materials for new technologies?

How can we do this quickly?

# Challenges in Engineering Materials

4

**“Engineering”**

Goal/means

0.1 nm

Many “knobs” to adjust

Performance

**Bottom Line:** Understanding or engineering new materials can require **many experiments** and, sometimes, **many years**

Processing



0.1 km

Many effects to decouple  
cause and effect

**“Science”**

Source: Olson. Science. (1997) 1237 [chain], Wikipedia [images]

# Computational Materials Engineering

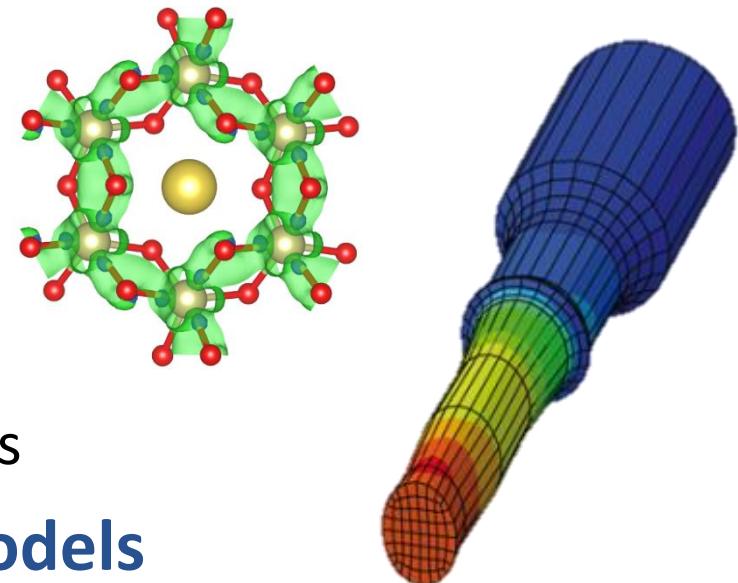
5

**Goal:** Accelerate design of materials

**Method:** Replace experiments with computers

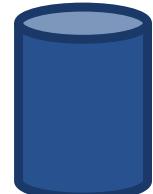
**Many Established Tools:**

- Density Functional Theory
- Phase Field
- Finite Element Analysis
- Computational Thermodynamics



**Emerging Field: Data-driven Models**

Materials  
Data



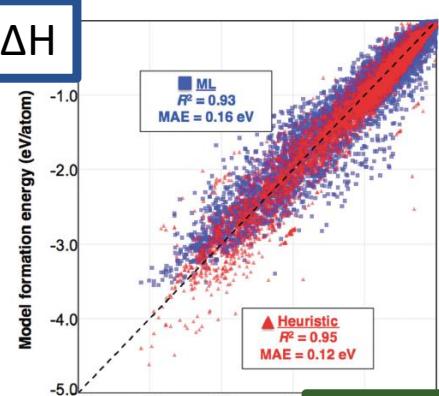
Machine Learning

$$\sigma_Y = f(x)$$

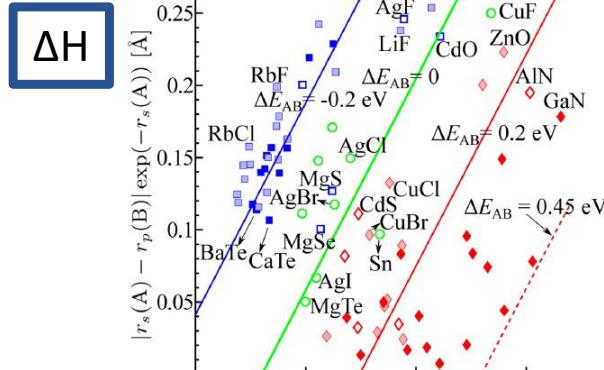
Predictive  
Model

# ML + Materials = “Materials Informatics”

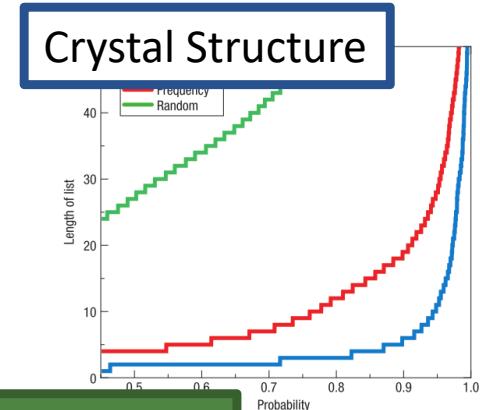
6



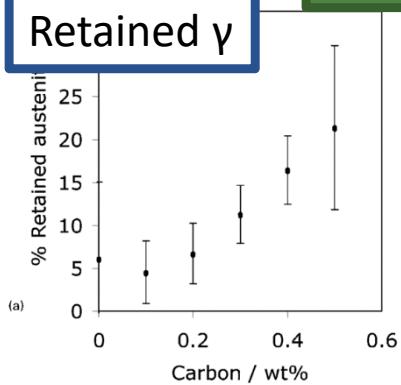
Meredig *et al.* PRB 80, 195101 (2009)



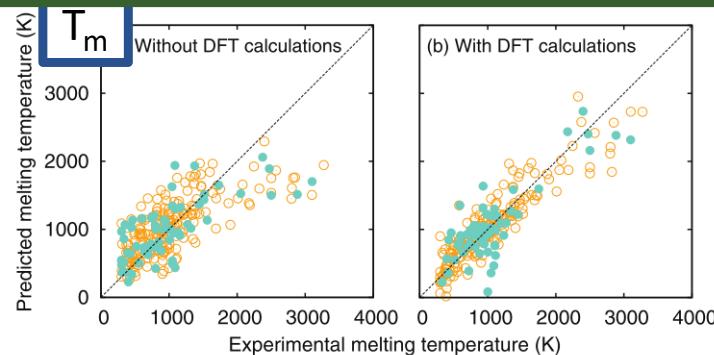
Nat. Mat. (2006), 641



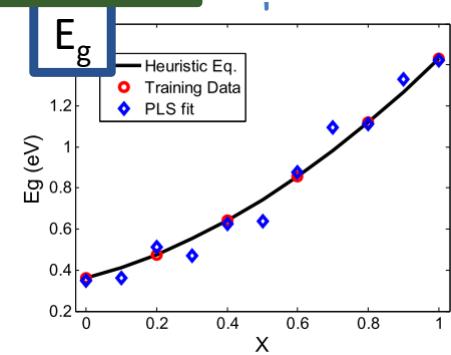
Question for today: Why do this?



Chatterjee *et al.* MS&T (2007), 819



Seko *et al.* PRB (2014), 054303



Srinivasan, Rajan. Materials (2013), 279

Reviews: Ward, Wolverton. COSSMS. (2017), 167.  
 Ward *et al.* MRS Bulletin. (2018), 683.

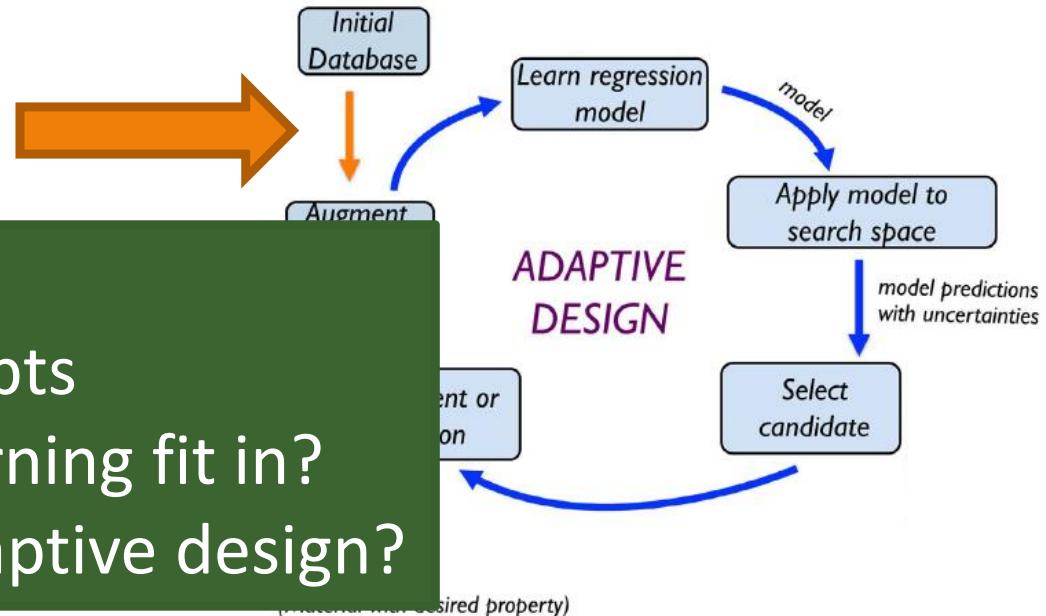
# Why Machine Learning?

7

## Advantages over Physics Models:

- ✓ Fast                                     $10^4\text{-}10^7$  evaluations/CPU/sec
- ✓ Adaptable                            Limited need to know underlying physics
- ✓ Self-correcting                      Improves with more data.
- ✓ Unbiased                              Can lead to unexpected predictions

Combined, this means



## My Talk Today:

- Dive into these concepts
- Where does deep learning fit in?
- How to do “deep” adaptive design?

# ML/DL and Materials Design

# ML models are fast

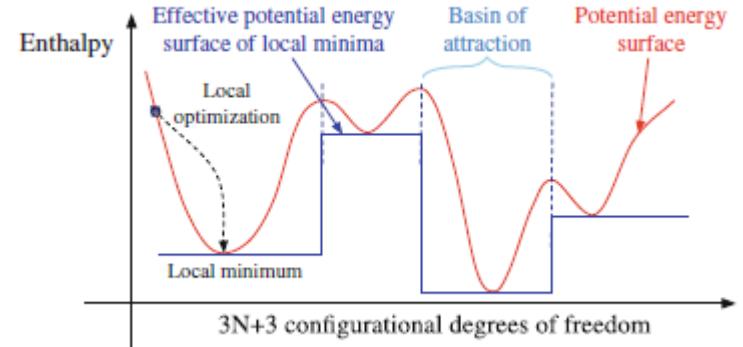
9

*Physics codes tend to be expensive*

**Example:** Density Functional Theory

DFT:  $10^0\text{-}10^3$  CPU-hr/compound

ML:  $\sim 10^{-5}$  CPU-hr/compound

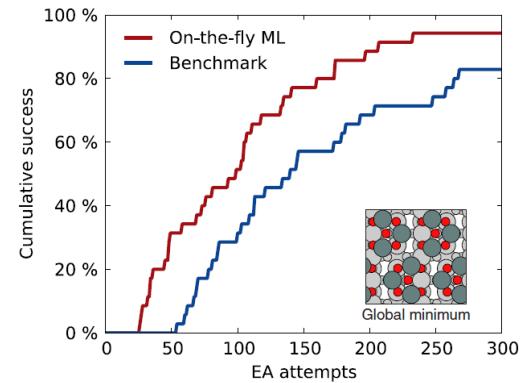


Revard *et al.* Top Curr Chem. (2014), 181–222.

**Design requires many evaluations:**

1. High-throughput parameter sweeps
2. GA-based optimization

...

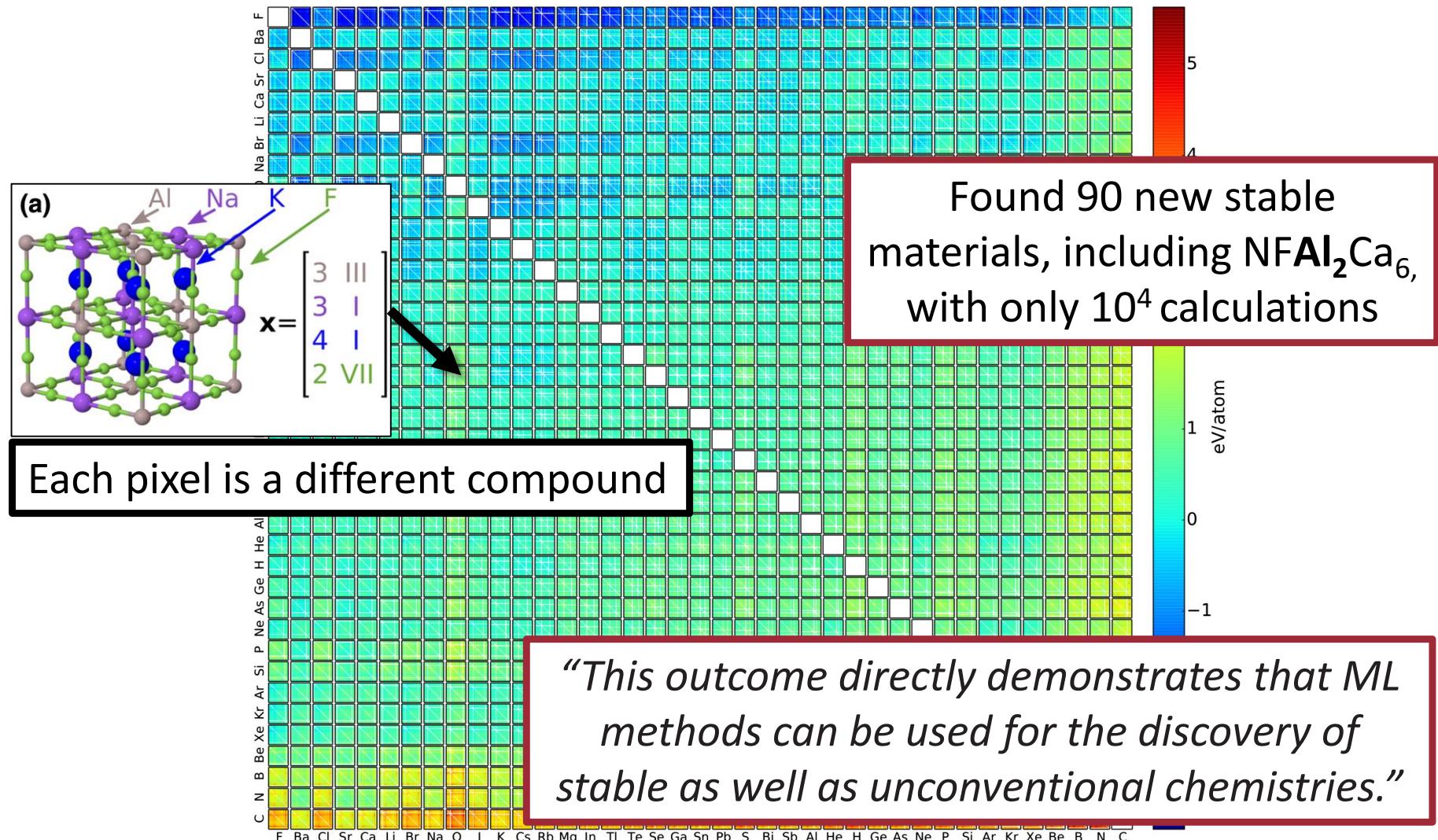


Jacobsen *et al.* PRL (2018), 026102

**Machine learning models are useful as fast surrogates**

# Scanning millions of Elpasolites

10

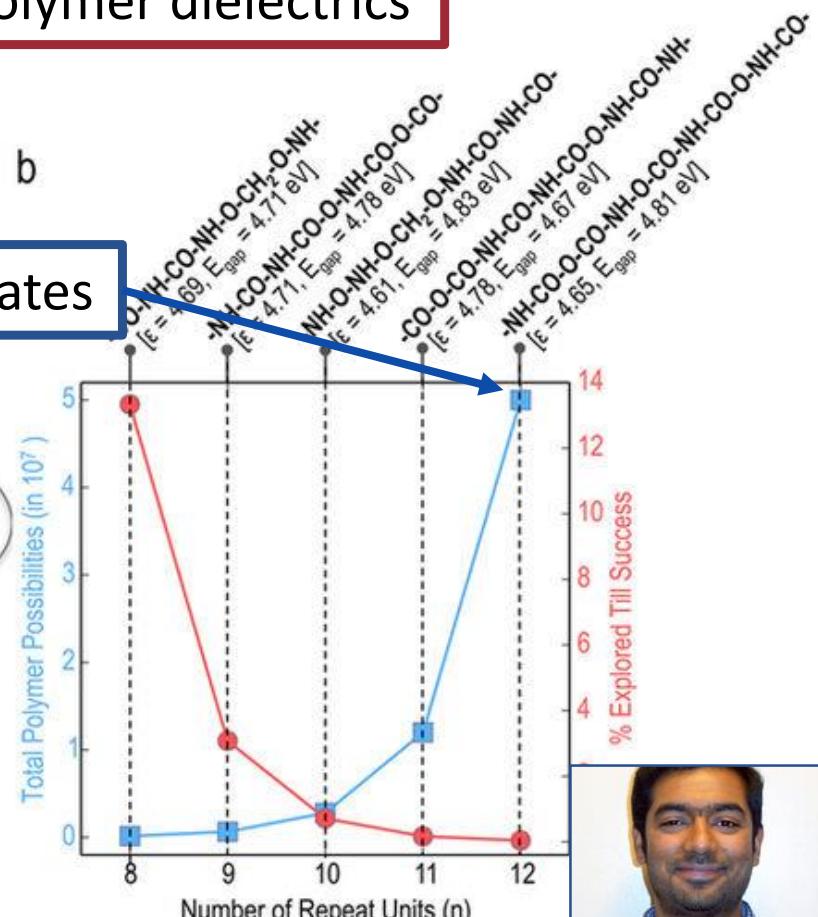
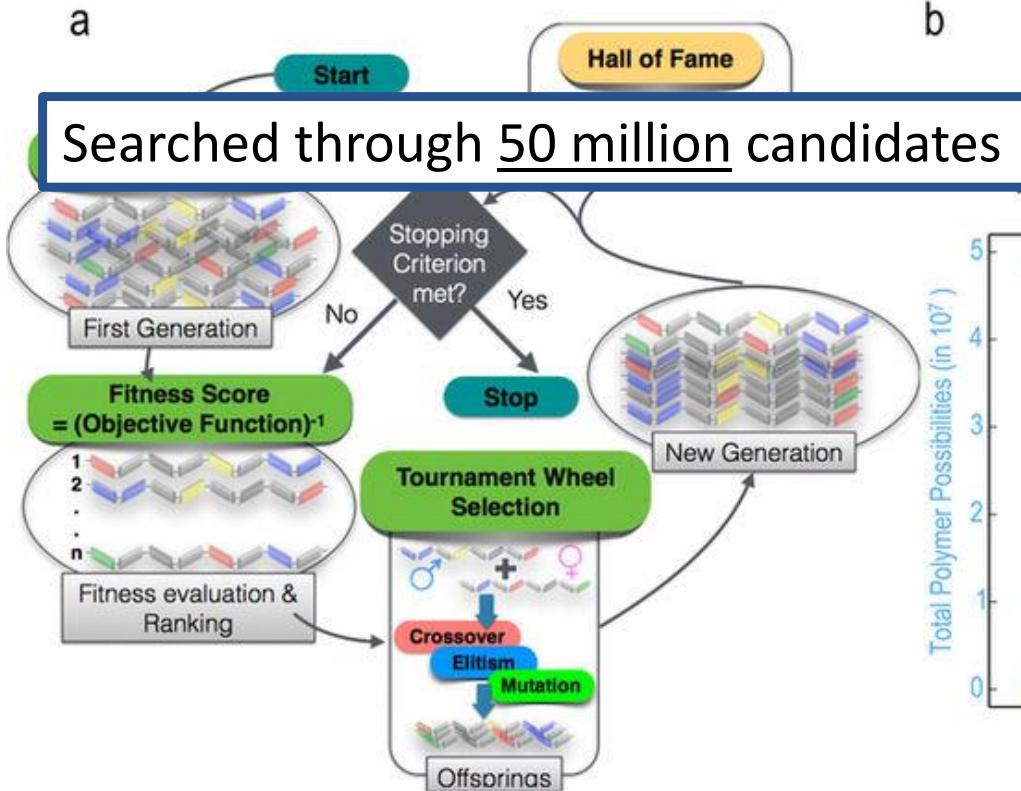


Ref: Faber *et al.* PRL. (2016) 1333502. doi: 10.1103/PhysRevLett.117.135502

# ML design of polymer dielectrics

11

Used ML+Genetic algorithm to find new polymer dielectrics



## *ML makes otherwise impossible searches practical*

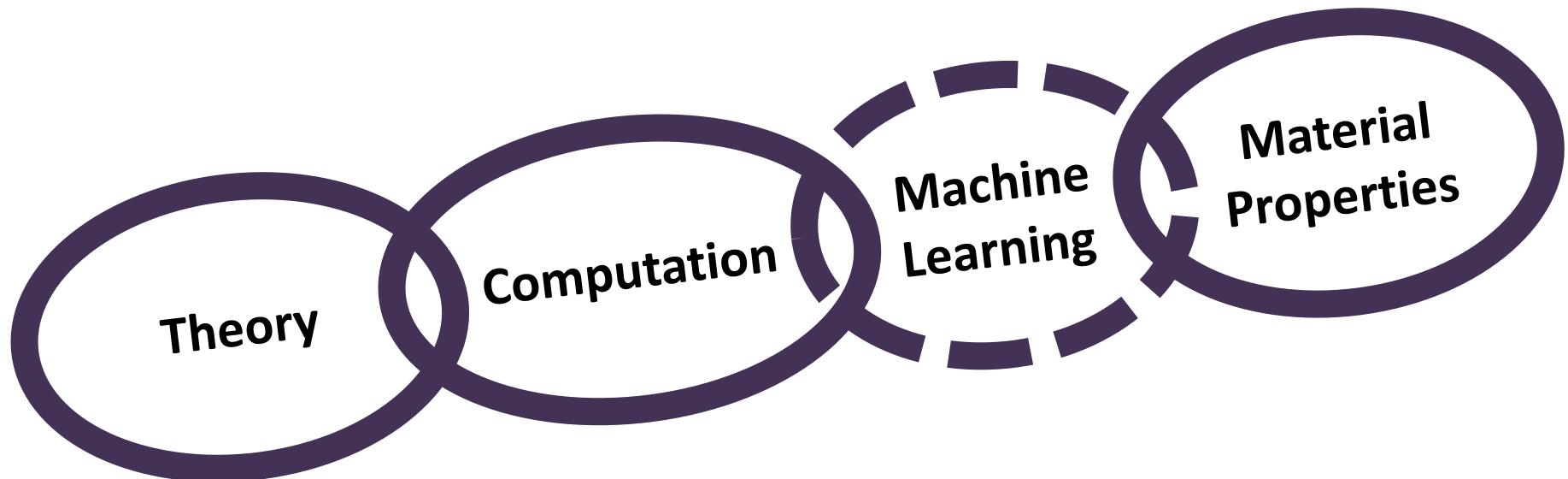
A. Mannodi-Kanakkithodi, CNM

# ML models are adaptable

12

*Not all properties are accurately computable (yet!)*

**Example:** Glass-forming ability of metal alloys

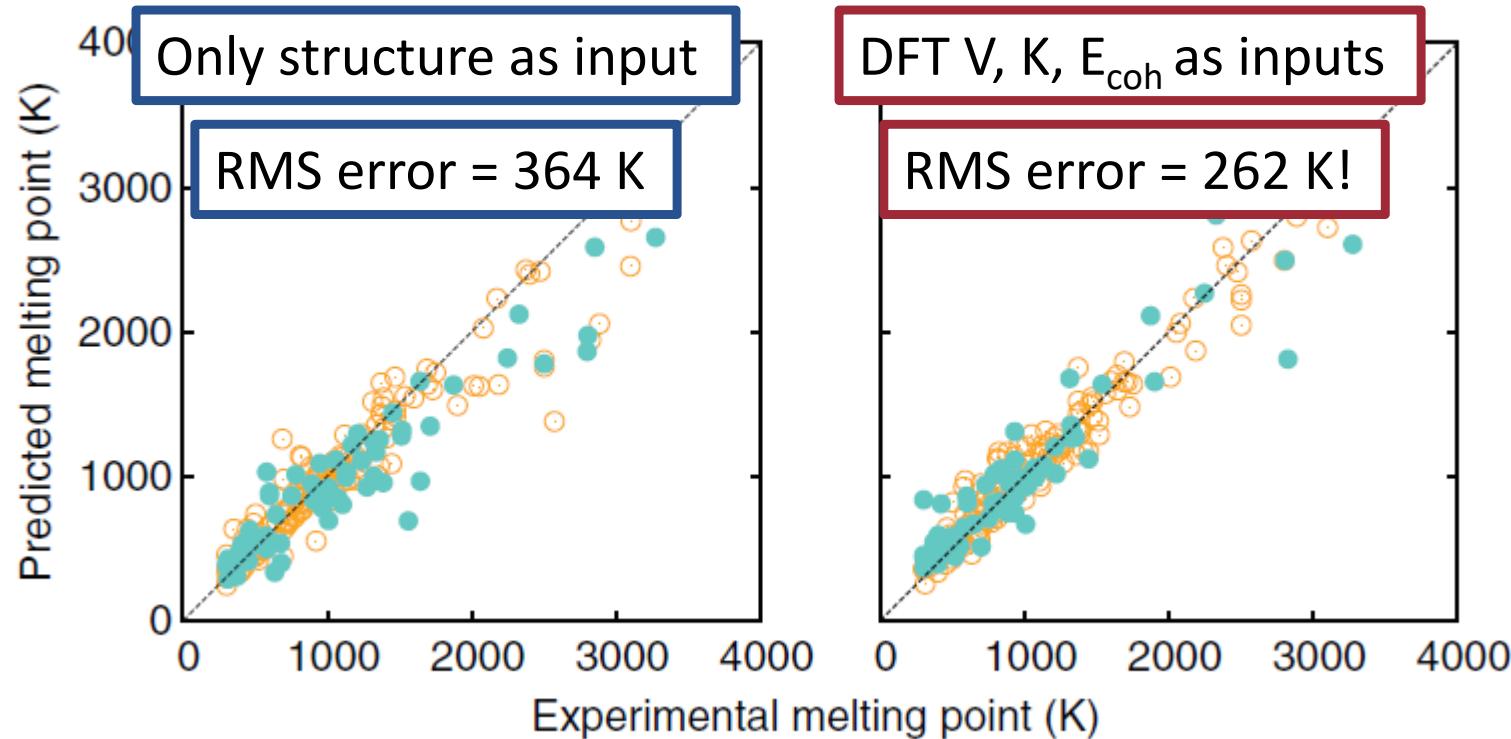


**ML can fill in links between computation and properties**

# Fast melting point calculations

13

*Melting temperature can be computed with DFT, but is very expensive*



**Physics models + ML leads to more accurate, faster predictions**

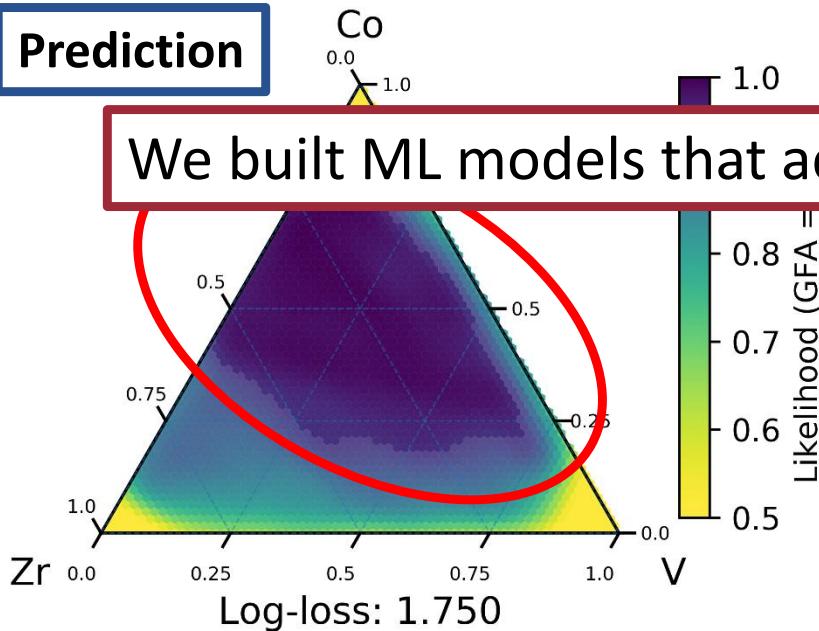
Ref: Seko *et al.*, PRB. (2014), 054303

# Predicting Glass-forming ability

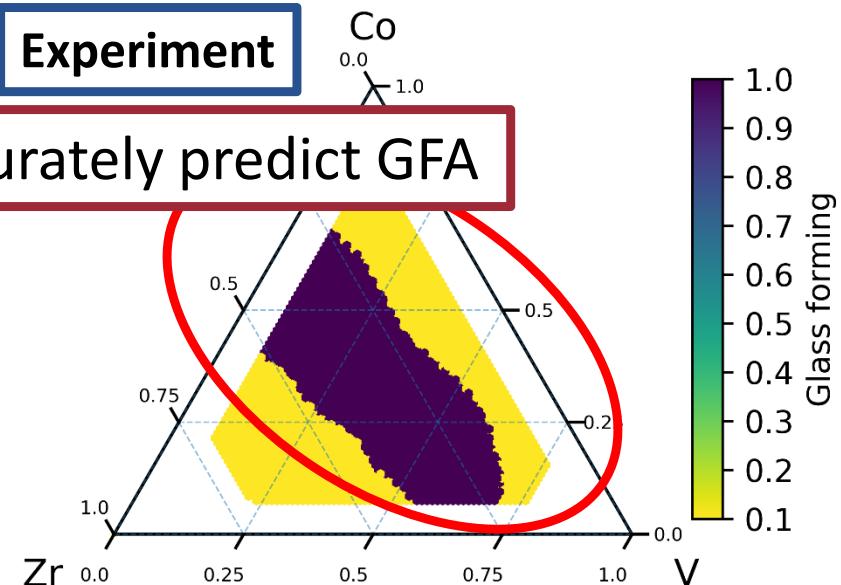
14

*Few models/theories link composition -> glass-forming ability*

**Prediction**



**Experiment**



*Machine learning fills missing links*

Ref: Ren et al. Sci. Adv., eaaq1566



J. Hattrick-Simpers, NIST  
Fang Ren, SLAC (Apple)  
Apurva Mehta, SLAC

# ML models are self-correcting

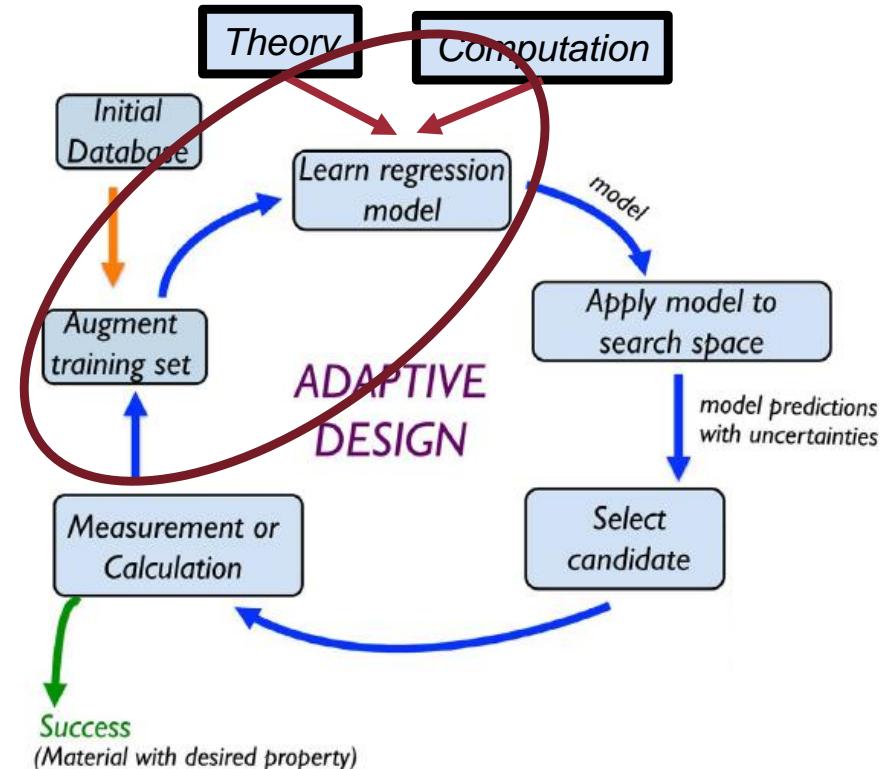
15

*Machine learning models are trivial to update*

Models can [often] be improved  
by just adding new data

**Best case:** No human effort

**Implication:** Updating model  
predictions faster than making  
new measurements



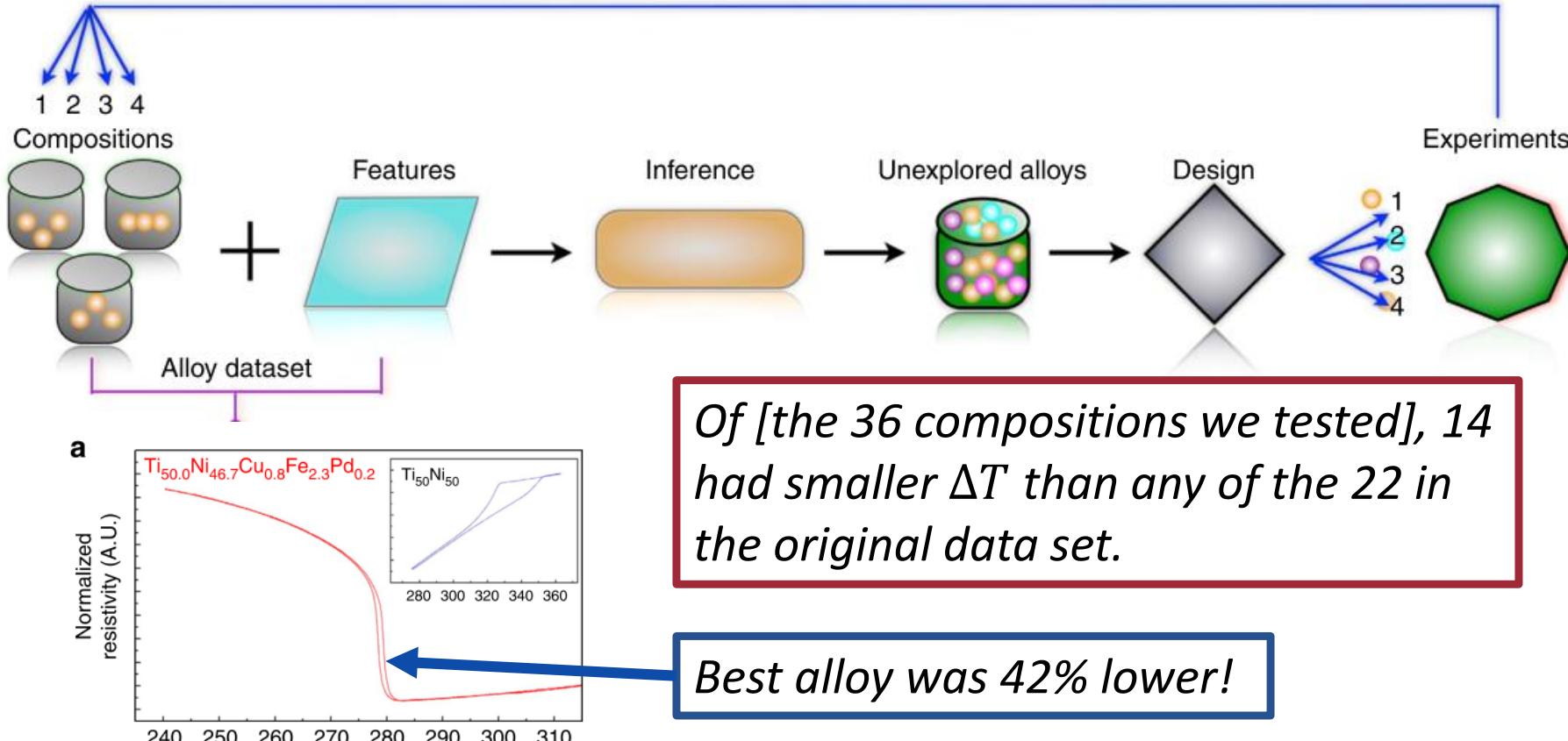
Easy updating allows for adaptive design

# Adaptive design of shape-memory alloys

16

## Example: Shape memory alloys with small $\Delta T$

Feedback from experiments: augmented data set with four new alloys



Adaptive search leads to rapid materials design

Ref: Xue *et al.* *Nat Comm.* (2016), 11241. doi: 10.1038/ncomms11241

# What about deep learning?

# Why Deep Learning?

18

*Conventional ML*



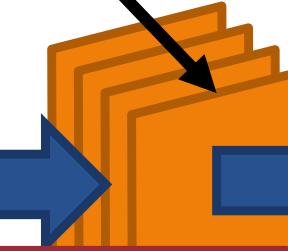
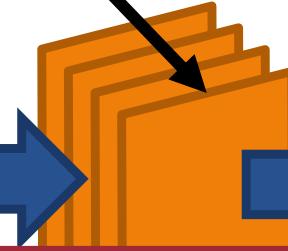
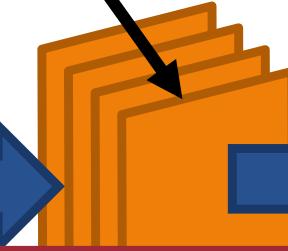
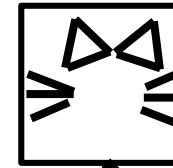
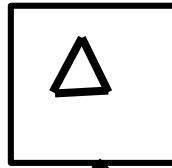
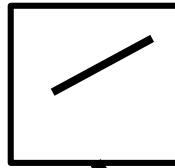
whiskers: True  
Eyes: 2  
Pupil aspect: 4  
...

ML

{ 'cat': 0.9,  
'dog': 0.1 }

Features require knowledge, model learned automatically

*Deep Learning*



{ 'cat': 0.9,  
'dog': 0.1 }

Features and model learned automatically

Question: Can representation learning work for materials?

# Predicting Crystalline Materials

19

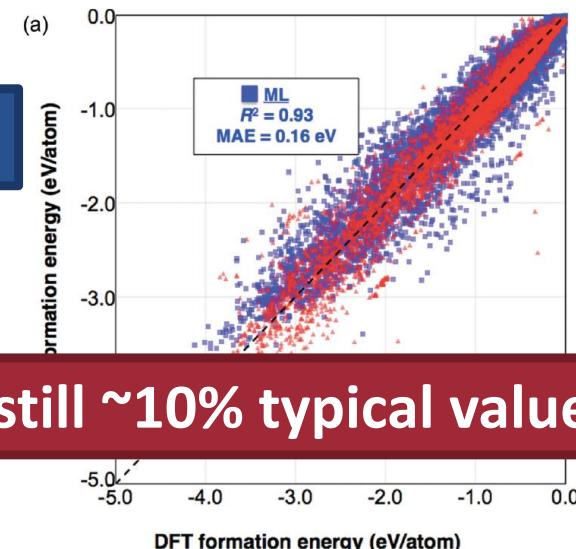
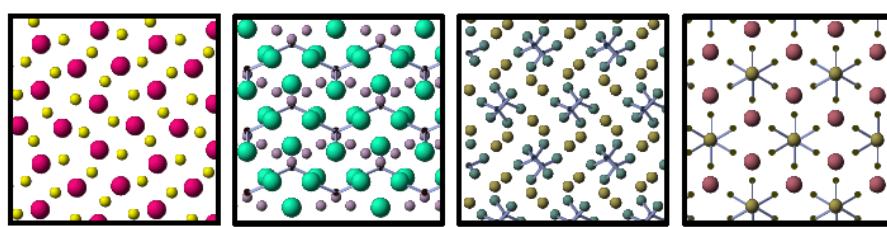
**Goal:** Identify new crystalline materials

**Model:** Given composition predict  $\Delta H_f$

**Initial attempt:** 2014



**Successfully identified 8 new compounds!**



**Can we do even better with deep neural networks?**

Ref: Meredig *et al.* PRB. (2014), 094104. doi: 10.1103/PhysRevB.89.094104

# Training Data

20

Collect

Process

Represent

Learn

**Data Source:** Open Quantum Materials Database

~470k DFT Calculations

## Why OQMD?

- ✓ Open ✓ Large
- ✓ Contains unstable materials

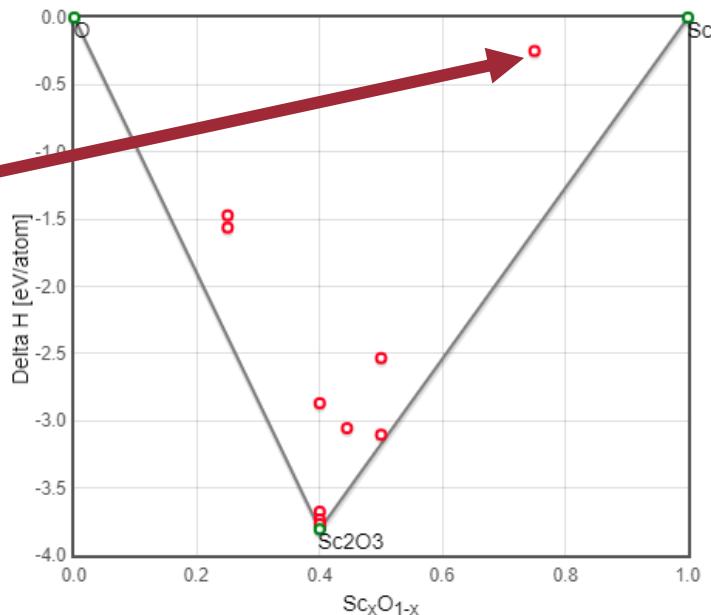
## Training Set:

Subset: Lowest energy structure

Size: 230336

Input: Composition

Output: Formation Enthalpy  $\Delta H_f$



# “Representation”

21

Collect

Process

Represent

Learn

Previous state-of-the-art: Compute physical features

Composition

Element Fractions

Physical Features

$\text{LiFePO}_4$

$x_H$	$x_{He}$	$x_{Li}$	...
0	0	$1/7$	...

Range Electronegativity:  
2.46

Is it charge-neutral:  
Yes

Maximum  $T_m$ :  
1538  
[...]

With Deep Learning: Rely on representation learning

# ElemNet: Architecture

22

Collect

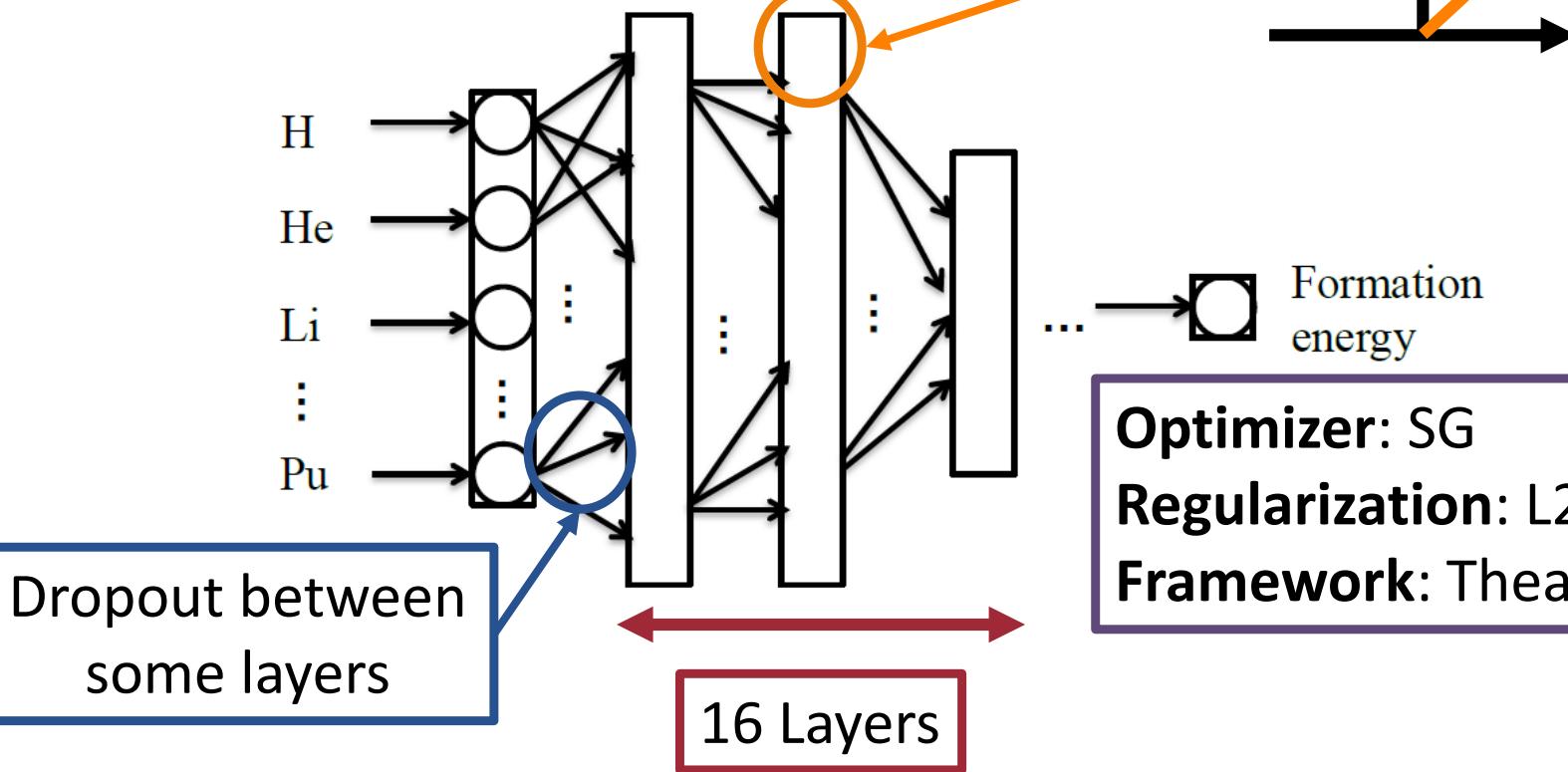
Process

Represent

Learn

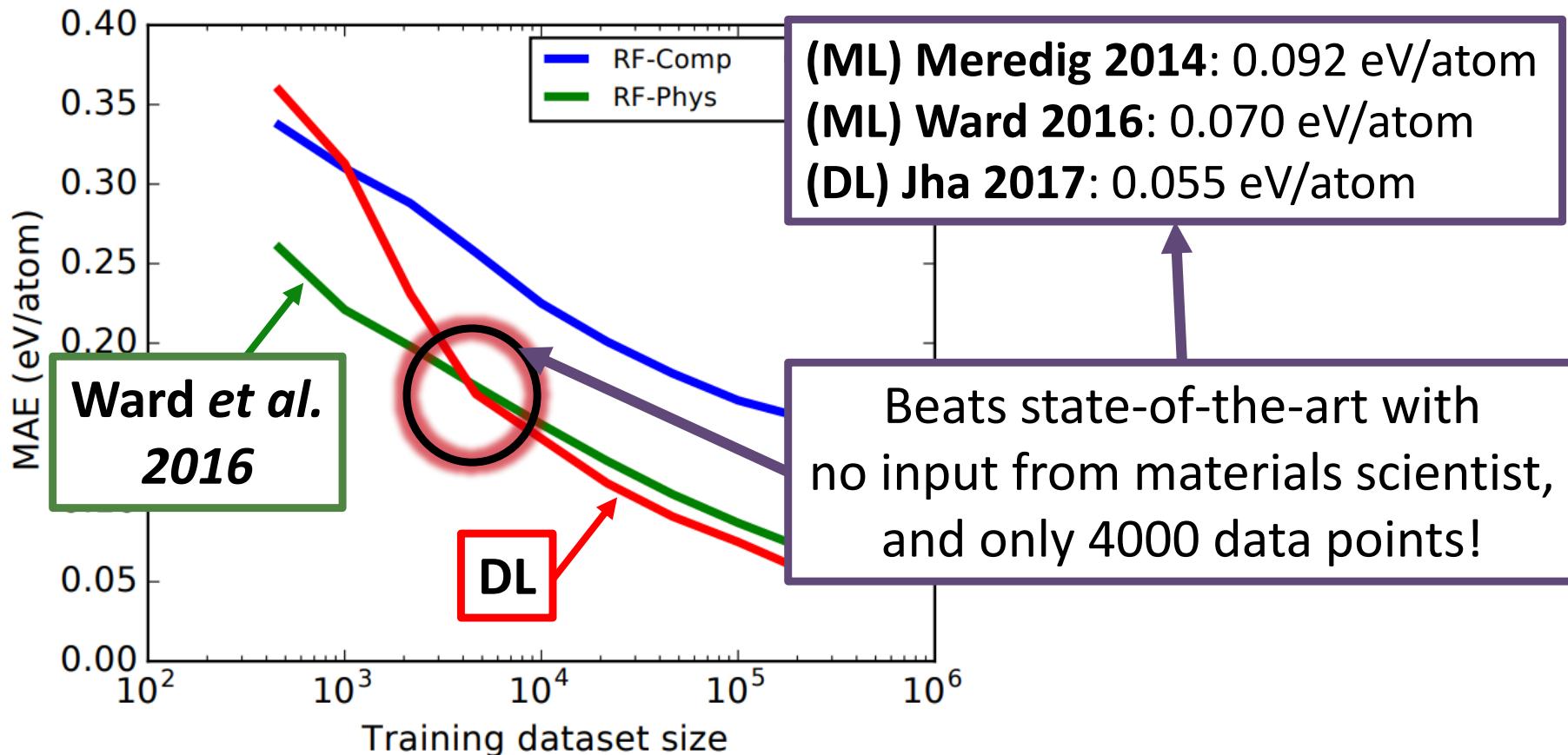
Fully connected MLP

ReLU activation



# Better than conventional learning?

23

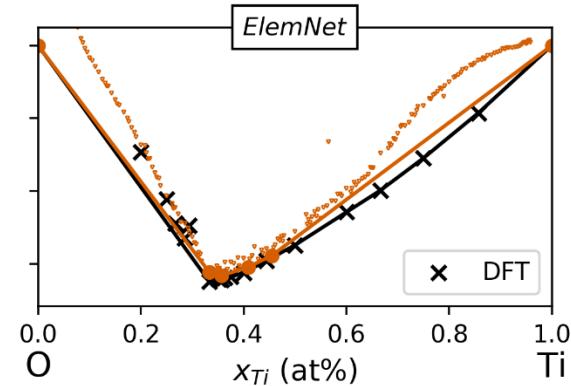
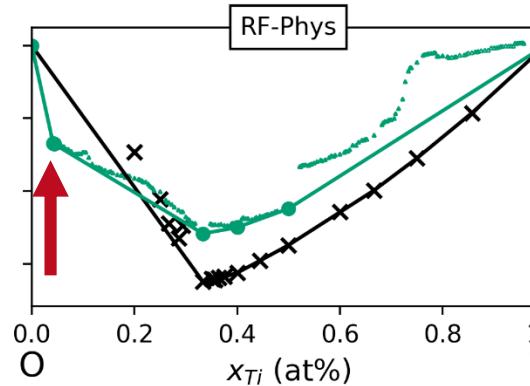
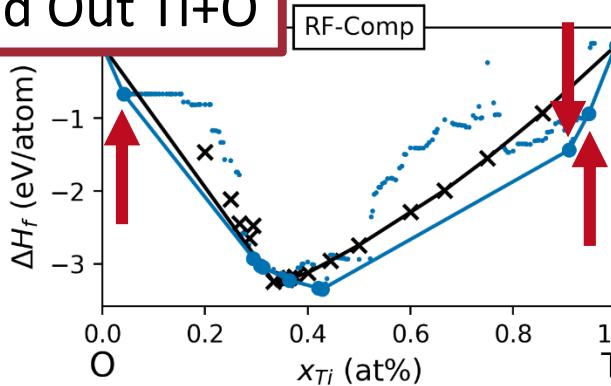


Ref: Jha *et al.*, *in review*

# Can DL interpolate between elements?

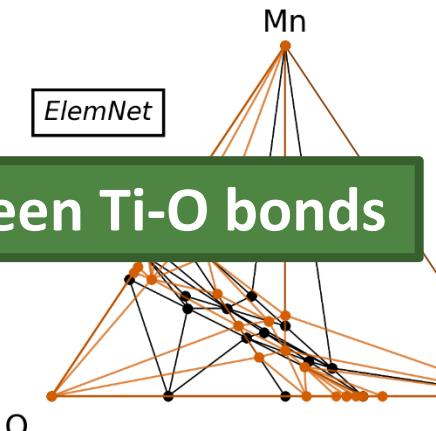
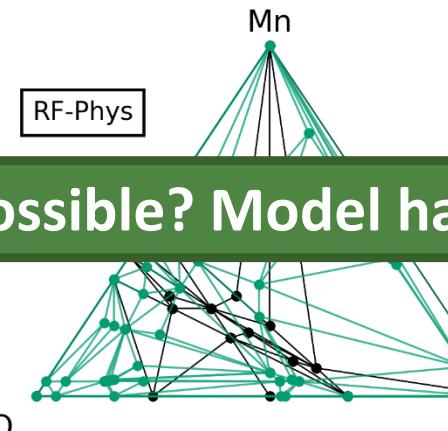
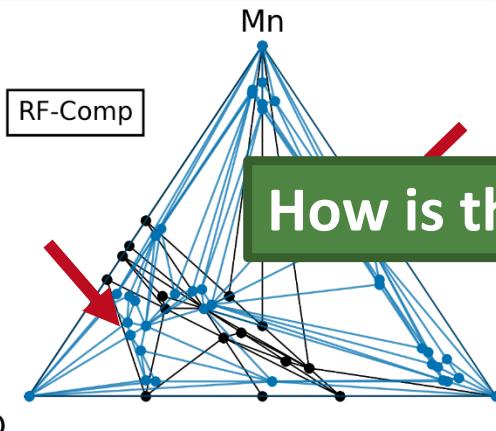
24

Hold Out Ti+O



Deep Learning Yields Fewer Spurious Predictions

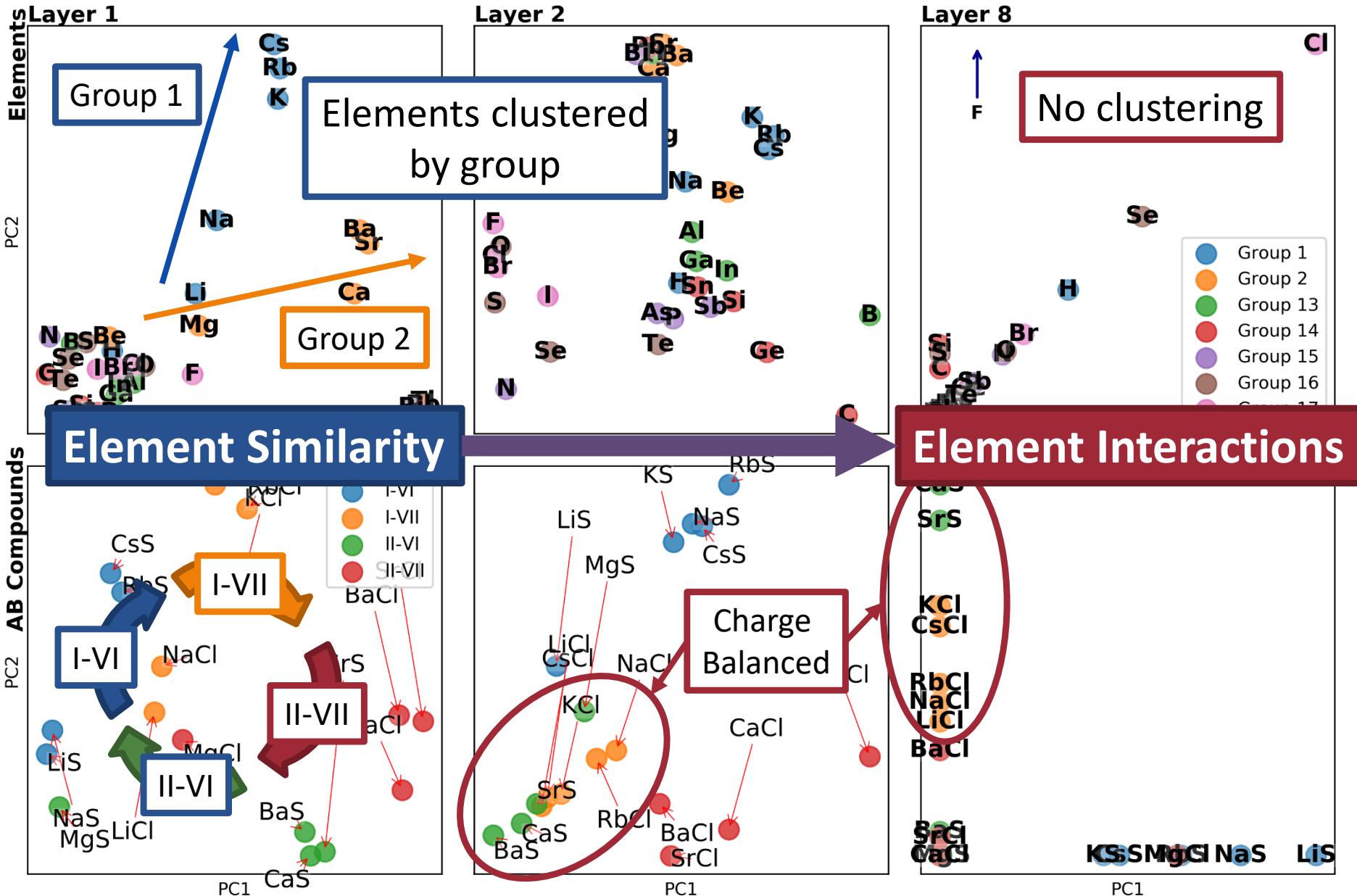
Hold Out Na-Fe-Mn-O



How is this possible? Model hasn't seen Ti-O bonds

# How is it working so well?

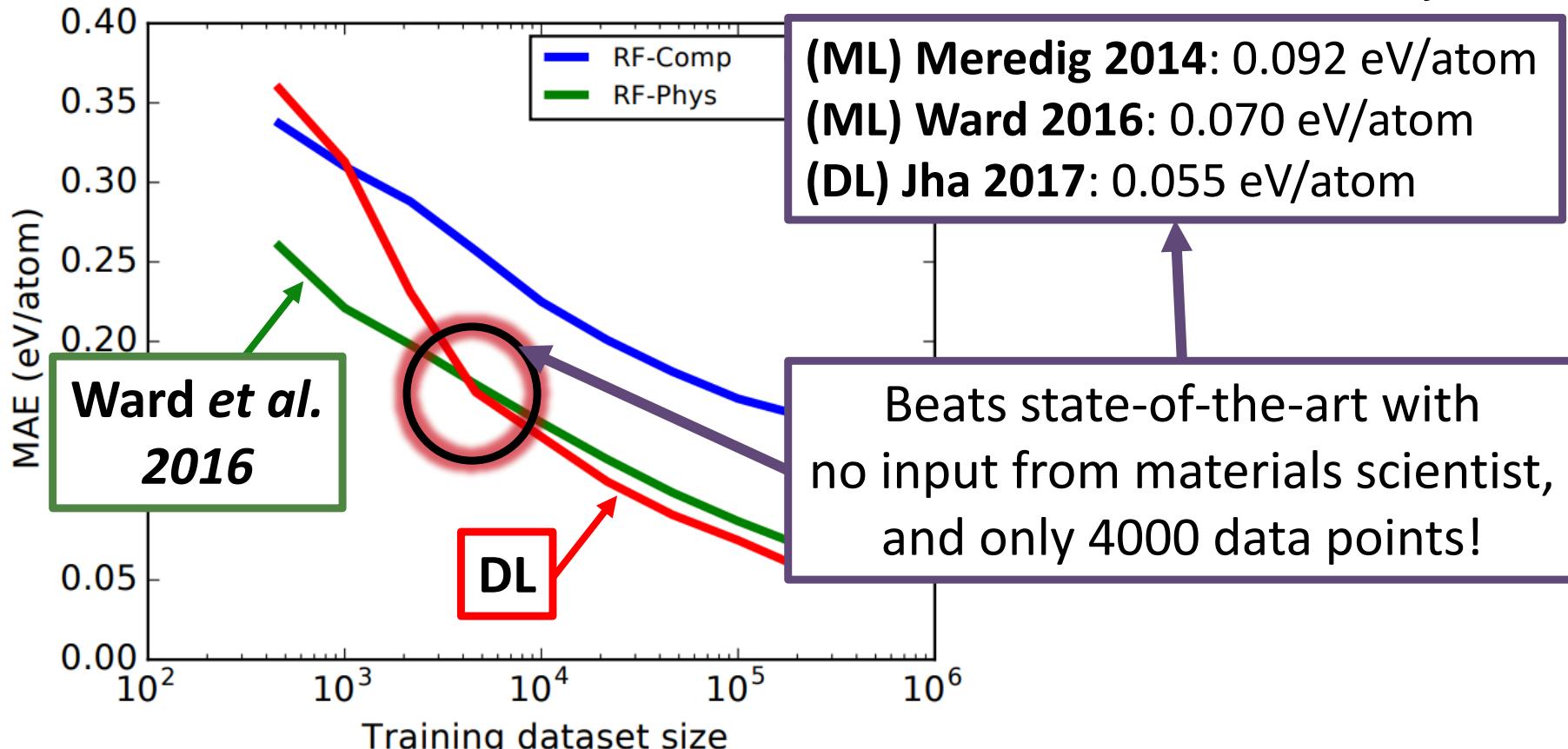
25



# Example #1: Composition -> Property

26

**Challenge:** Given element fractions, predict  $\Delta H_f$



Deep learning could change the way we approach problems

Ref: Jha *et al.*, *in preparation*

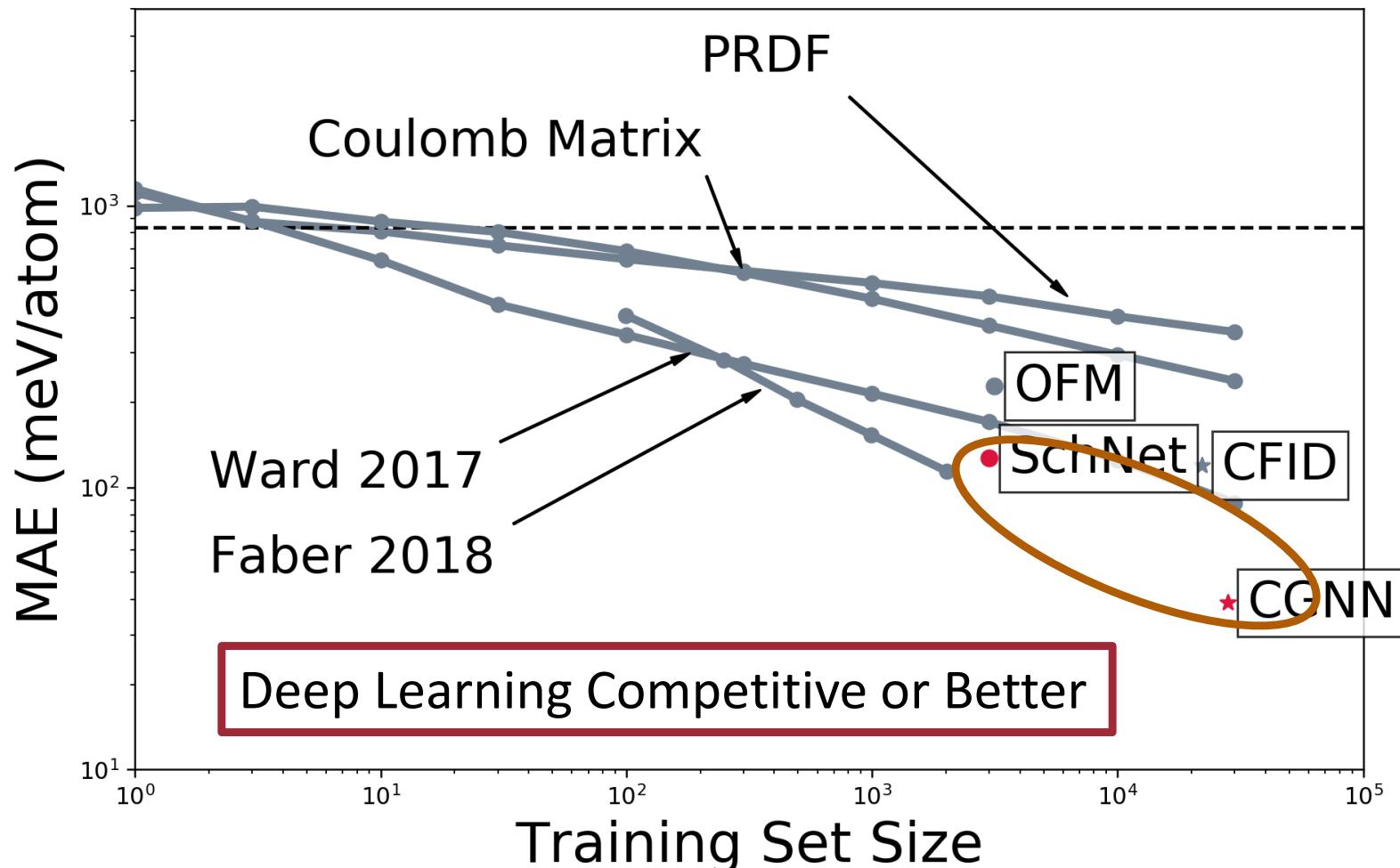
Dipendra Jha, NU

# Example #2: Crystal Structure

27

**Dataset:** 32k DFT  $\Delta H_f$  from the OQMD

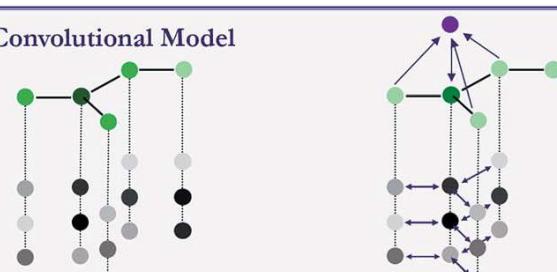
**Test:** Remove 1000, train on  $N$  remaining



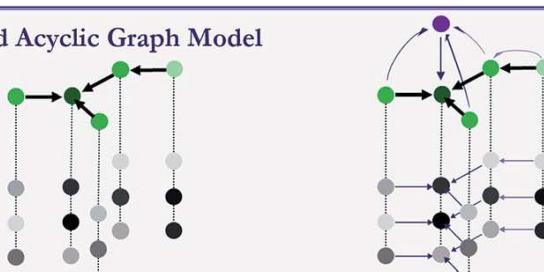
# Example #3: Molecules

28

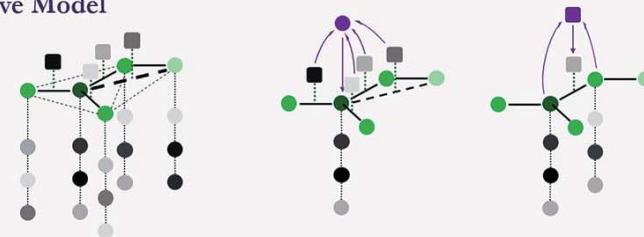
A. Graph Convolutional Model



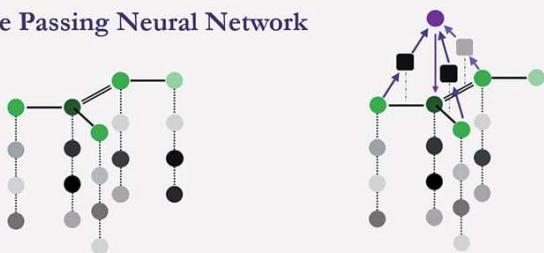
B. Directed Acyclic Graph Model



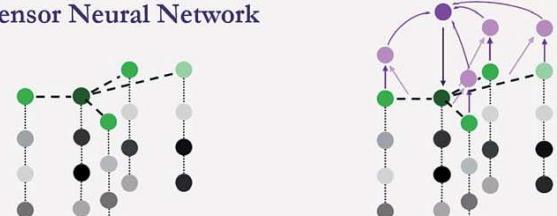
C. Weave Model



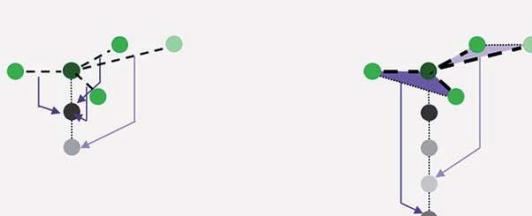
D. Message Passing Neural Network



E. Deep Tensor Neural Network



F. ANI-1



Many methods available for molecular data!

# DL not a cure-all, but taking over

29

Table 3: Summary of performances(test subset): conventional methods versus graph-based methods. Graph-based models outperform conventional methods on 11/17 datasets.

Category	Dataset	Metric	Best performances - conventional methods	Best performances - graph-based methods
Quantum Mechanics	QM7	MAE	KRR(CM): 10.22	<b>DTNN: 8.75</b>
	QM7b	MAE	KRR(CM): 1.05	<b>DTNN: 1.77*</b>
	QM7	Multitask	Multitask: 0.0150	<b>MPNN: 0.0143</b>
	QM7b	Multitask	Multitask(CM): 4.35	<b>DTNN: 2.35</b>
	QM7	MAE	Post: 0.99	<b>MPNN: 0.58</b>
	QM7b	MAE	Post: 1.74	<b>MPNN: 1.15</b>
Biophysics	MOV	AUC-ROC	Post: 0.799	<b>GC: 0.655</b>
	MOV	Multitask	Post: 0.129	<b>GC: 0.136</b>
	HIV	AUC-ROC	Multitask: 0.184	Weave: 0.109
	BACE	AUC-ROC	<b>KernelSVM: 0.792</b>	GC: 0.763
	PDBbind(full)	RMSE	<b>RF: 0.867</b>	Weave: 0.806
	BBBP	AUC-ROC	<b>KernelSVM: 0.729</b>	GC: 1.44
Physiology	Tox21	AUC-ROC	KernelSVM: 0.822	<b>GC: 0.690</b>
	ToxCast	AUC-ROC	Multitask: 0.702	<b>GC: 0.829</b>
	SIDER	AUC-ROC	<b>RF: 0.684</b>	Weave: 0.742
	ClinTox	AUC-ROC	Bypass: 0.827	GC: 0.638
				<b>Weave: 0.832</b>

For benchmark problems of learning  
from molecular data,  
conventional ML better for 6/17 cases

\* As discussed in section 4.4, DTNN outperforms KRR(CM) on 14/16 tasks in QM7b while the mean-MAE is skewed due to different magnitudes of labels.

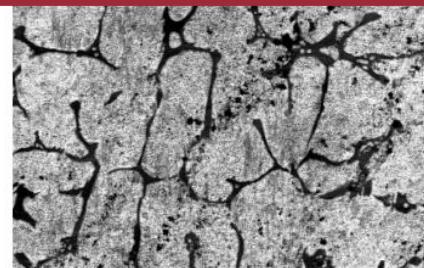
# Example 4: Microstructural Data

30

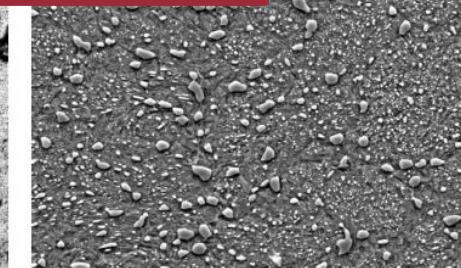
**Major focus area:** DL  
and microstructural data

**Why?** DL techniques  
work well with images

## Microstructure Classification



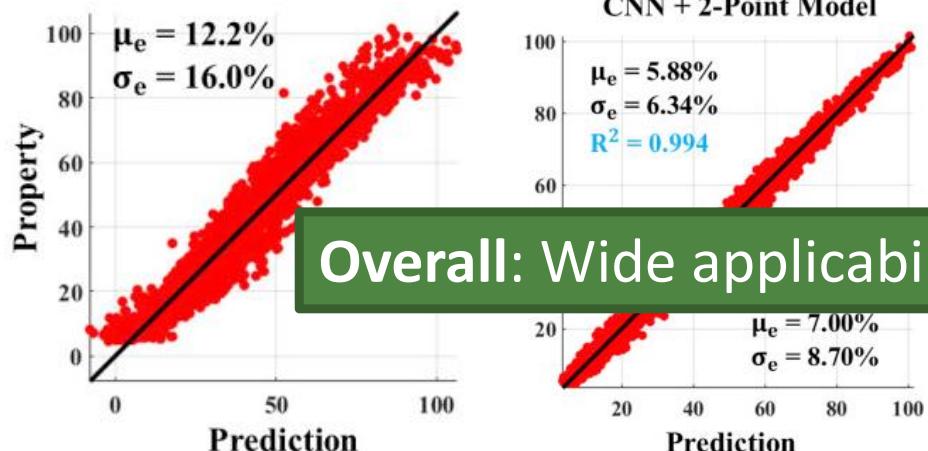
(a) Network



(b) Spheroidite

Ling *et al.* arXiv: 1711.00404v1

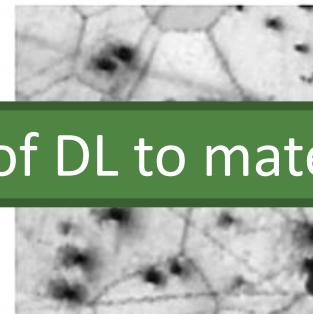
## Surrogate Models for FEM



Cecen *et al.* *Acta Mat.* (2018), 76

## Ionic Conductivity

### Input images



Specific features for  
*low* ionic conductivity



Kondo *et al.* *Acta Mat.* (2017), 29

# Summary

31

## Why use machine learning in materials?

- ✓ Fast
- ✓ Adaptable
- ✓ Self-correcting
- ✓ Unbiased

## Where does deep learning fit in?

- More accurate than conventional ML
  - For many types of data!

## Where's the big opportunity(s)?

- Adaptive design for materials
- Integrating physics into deep learning
- Integrating deep learning into design

# Goal: Tight integration with design

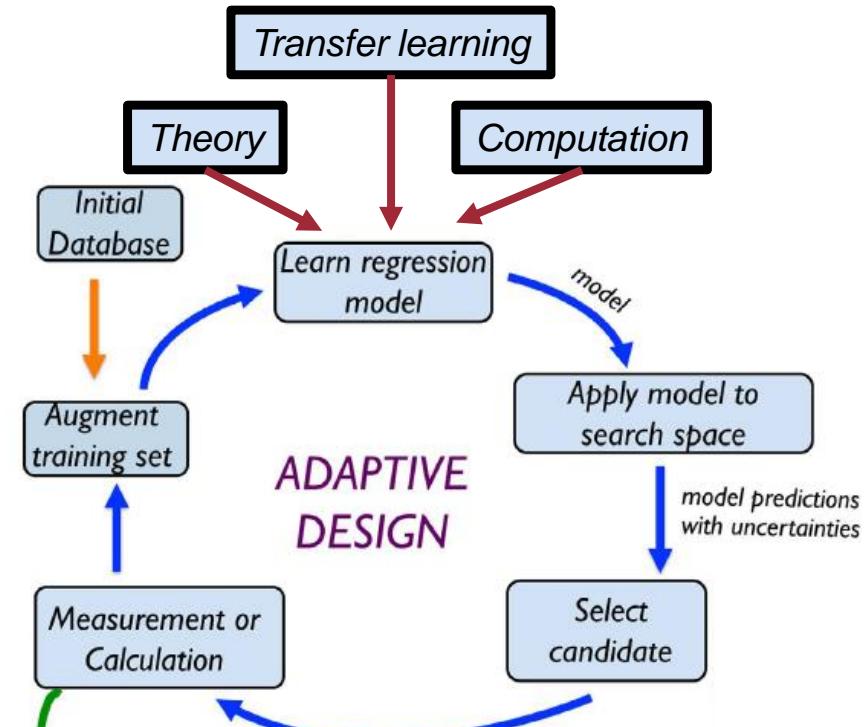
32

Future: Adaptive design with automated experiments

What do we already have?:

- Materials data
- ✓ Representations
- ✓ ML algorithms

Need: Tight integration  
with materials design



Problem: Disconnect between method developers and model users

(Material with desired property)

Figure: Balachandran *et al.* *Sci. Rep.* (2016), 19660. doi: 10.1038/srep19660



ARGONNE LEADERSHIP  
COMPUTING FACILITY



# DLHub

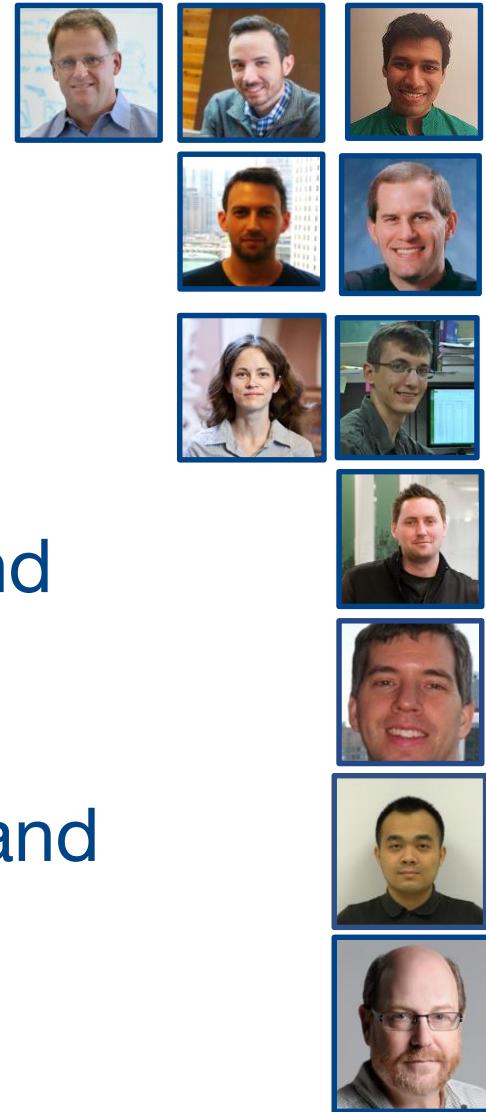
Data and Learning Hub for Science



Ben Blaiszik ([bblaiszik@anl.gov](mailto:bblaiszik@anl.gov)), Ryan Chard, Logan Ward,  
Kyle Chard, Zhuozhao Li, Anna Woodard, Yadu Babuji,  
Steve Tuecke, Mike Franklin, Ian Foster

**Funding:** 2018 Argonne Advanced Computing LDRD

# Data and Learning Hub for Science (DLHub)

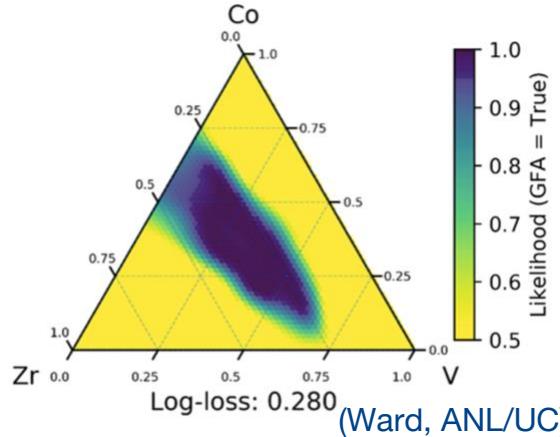


- Collect, publish, categorize models from many disciplines (materials science, physics, chemistry, genomics, etc.)
- Serve models via API to simplify sharing, consumption, and access
- Enable new science through reuse, real-time integration, and synthesis of existing models

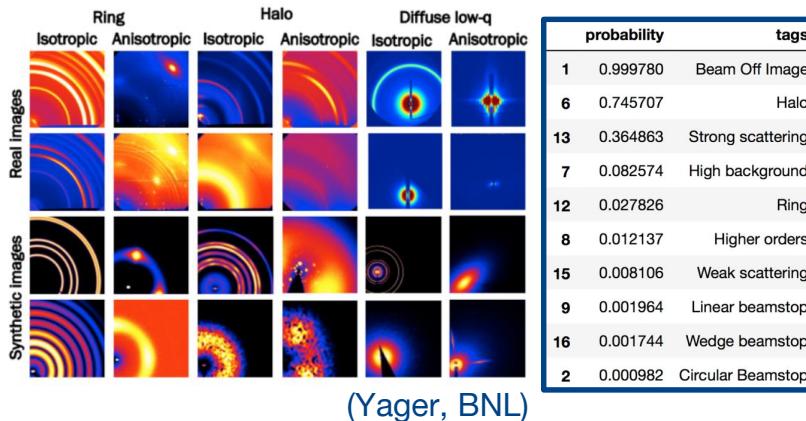
# Select DLHub Use Cases

## Model-driven Experimentation and Data Tagging

- Metallic glass discovery [active learning]

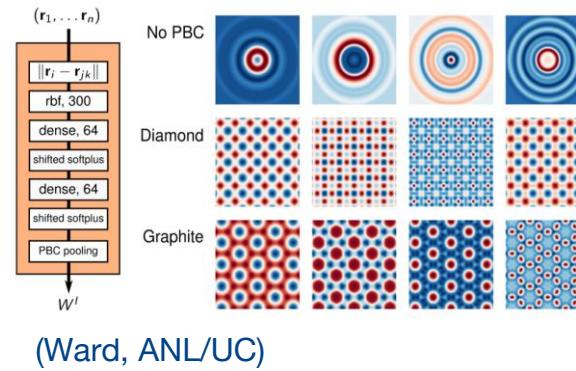


- XRD beamline image tagging

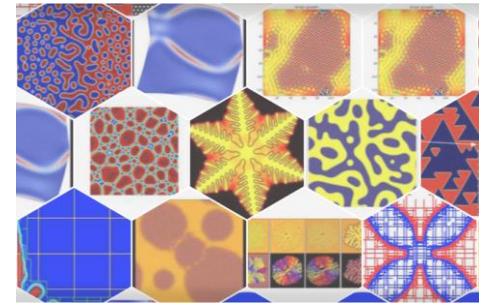


## Community Model Benchmarking

- Crystal structure



- NIST PFHub



(Wheeler, Warren, Heinonen  
NIST/UC/Argonne/NU)

## Automated Model Retraining with New Data

- Models linked to dynamic data sources

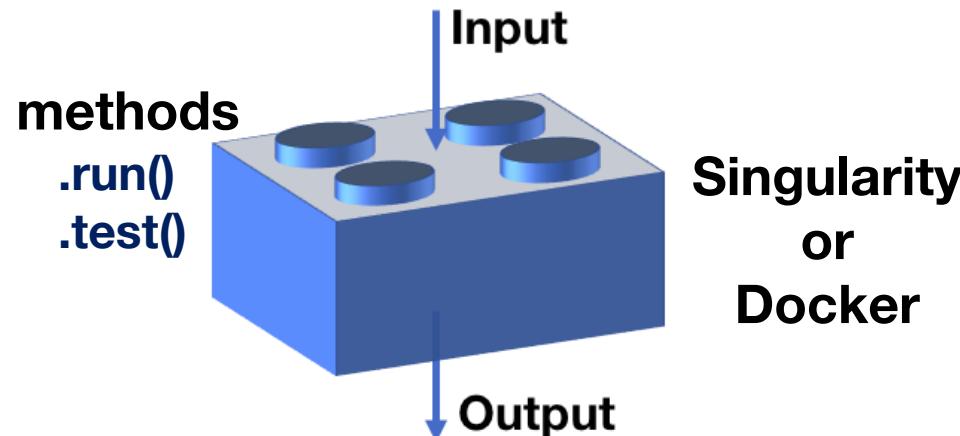


(Center for Hierarchical Materials  
Design NIST/UC/Argonne/NU)

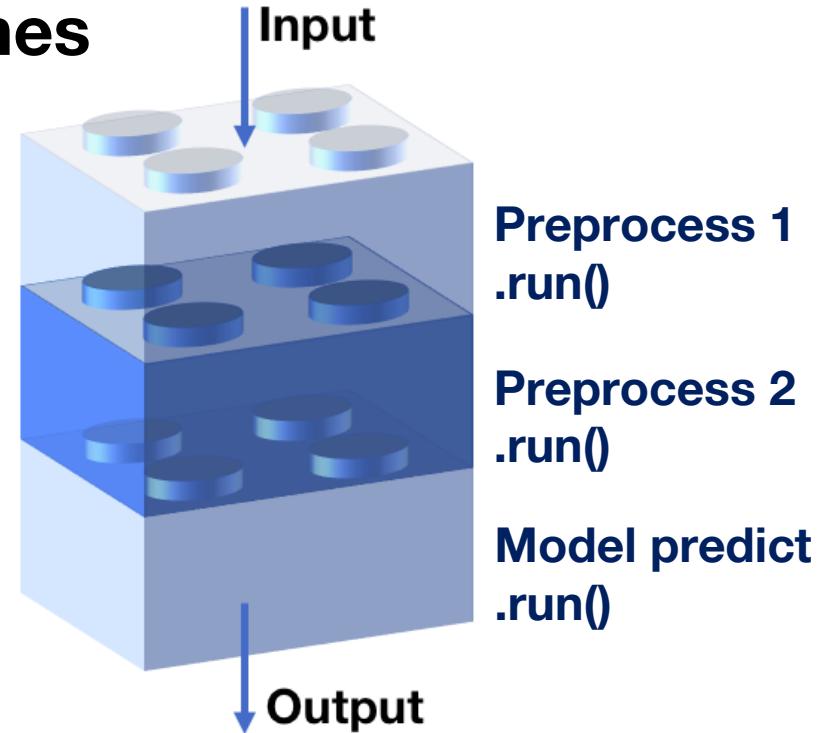
# DLHub Servables and Pipelines



## Servables

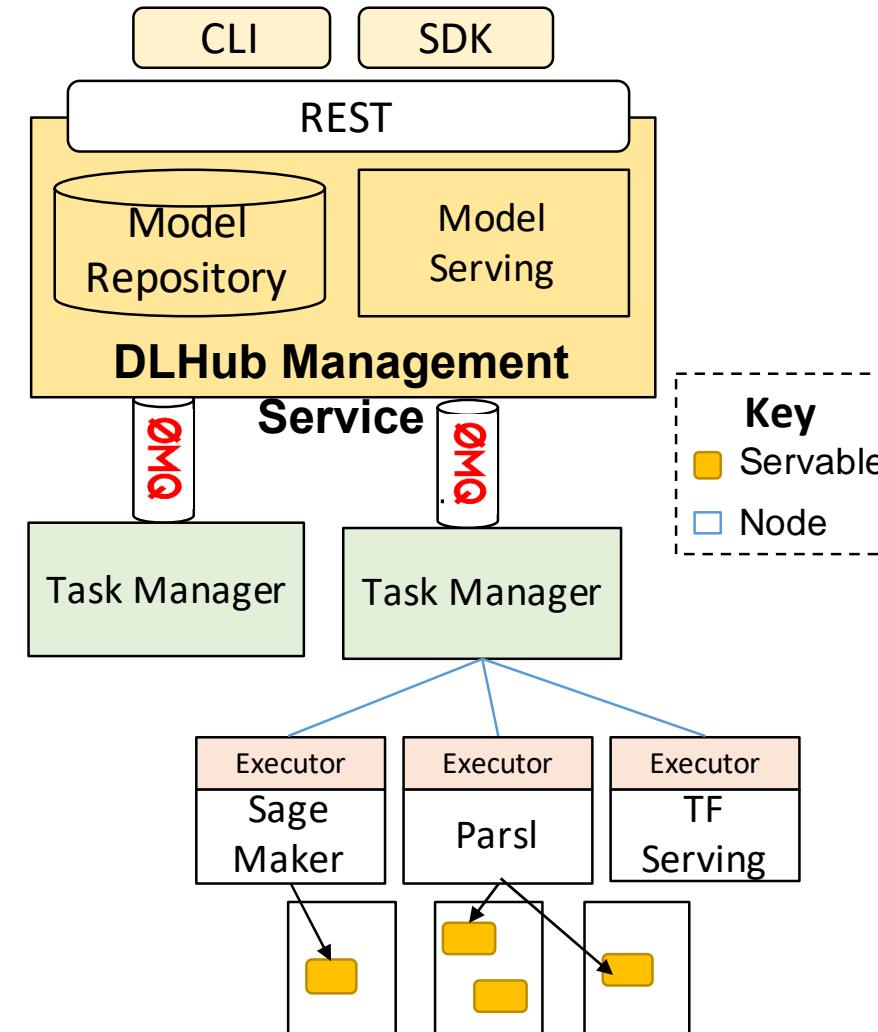


## Pipelines



# DLHub Architecture

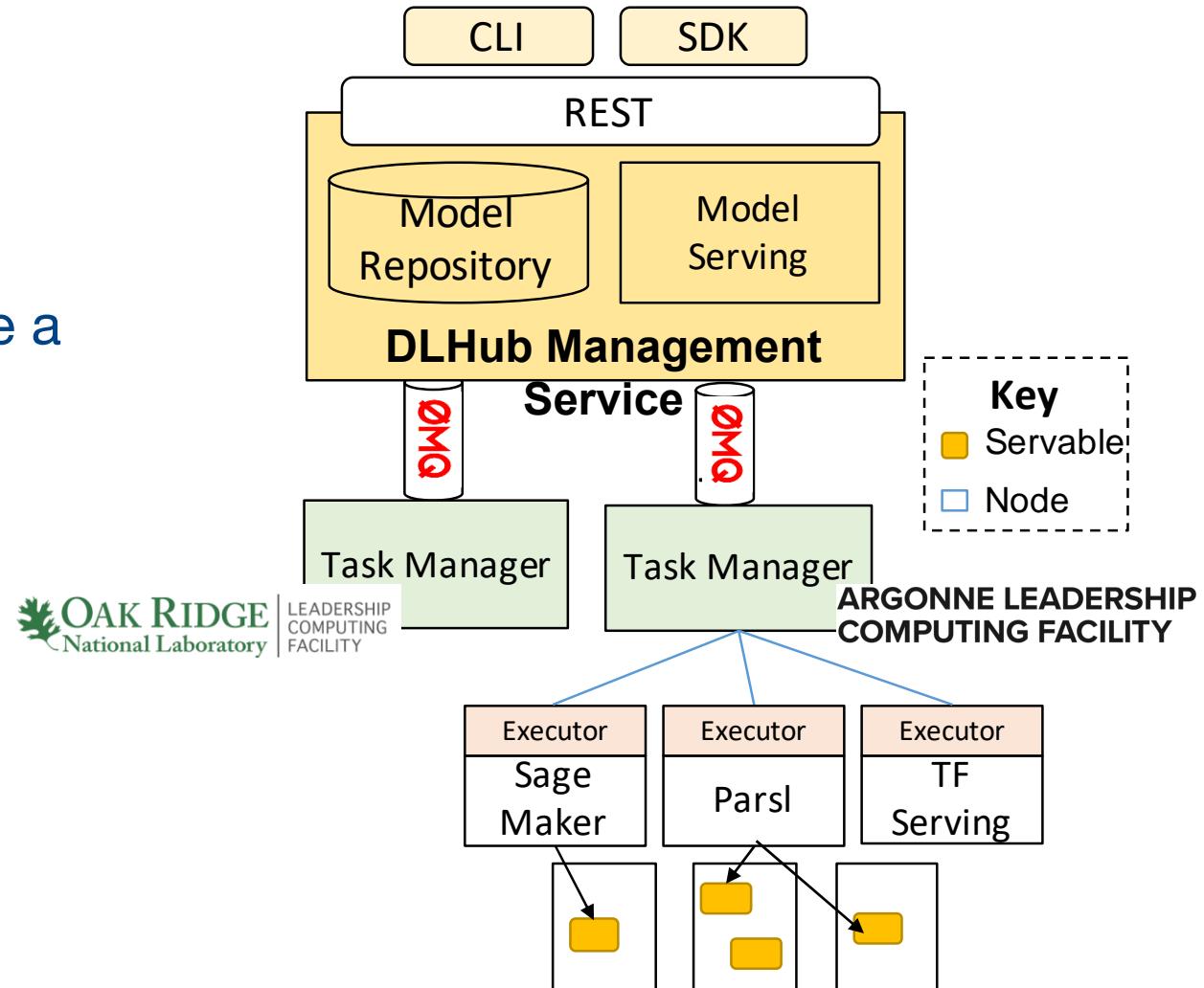
- REST API with Python SDK (available) / CLI (delivery in Nov. 2018 )
  - Support model markup, data staging, registration, and invocation
- Model Repository
  - Container registry
  - Advanced search functions
  - Identifier minting capabilities



<https://github.com/DLHub-Argonne>

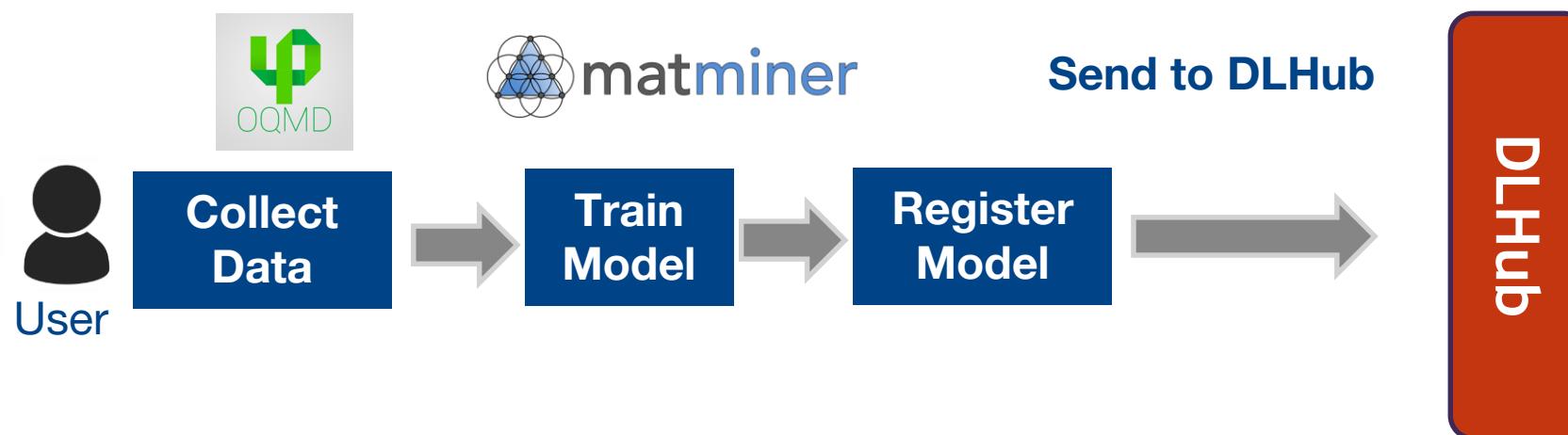
# DLHub Architecture

- Task Managers (TM) to support execution on various compute resources
- Executors chosen by TM to invoke a given servable'
- Caching at TM
- Data staging with Globus
- Batch submissions
- Scalability through deployment of model replicas



# DLHub Model Registration and Publication

- **Register** model metadata, weights, and files to improve discoverability and reusability
- **Containerize** model to enhance interoperability
- **Identify** model with a permanent identifier (e.g., DOI, minid)
- **Version** model and data pre/post processing steps
- **Deploy** model with simplified interfaces for users
- **Control** access to model metadata and usage
- (future) **Automate** retraining and testing when new data are available



# Marking up a Model - Python SDK

Existing Model

User Mark Up with SDK

SDK Extracts Metadata  
for Known Model  
Types

Send to DLHub  
(via Globus or HTTPS)

DLHub  
Containerization

Populate Search  
Index / Mint  
Identifiers

```
from dlhub_toolbox.models.servables.keras import KerasModel
import pickle as pkl
import json

# Describe the keras model
model = KerasModel('model.hd5', list(map(str, range(10)))))

# Describe the model
model.set_title("MNIST Digit Classifier")
model.set_name("mnist_tiny_example")
model.set_domains(["general","digit recognition"])

# Add link to paper describing the dataset
model.add_related_identifier("10.1109/CVPR.2007.383157", "DOI",
| | | | | | | | "IsDescribedBy")

model.set_authors(["Lecunn, Yann", "Cortes, Corinna"])

# Describe the outputs in more detail
model.output['description'] = 'Probabilities of being 0-9'
model.input['description'] = 'Image of a digit'
```

# Python SDK - Automated Metadata Generation

## Citation Metadata

```
"datacite": {  
    "creators": [{  
        "givenName": "Yann",  
        "familyName": "Lecunn",  
        "affiliations": []  
    },  
    {  
        "givenName": "Corinna",  
        "familyName": "Cortes",  
        "affiliations": []  
    }  
],  
    "titles": [  
        {"title": "MNIST Digit Classifier"}],  
    "publisher": "DLHub",  
    "publicationYear": "2018",  
    "relatedIdentifiers": [  
        {"relatedIdentifier": "10.1109/CVPR.2007.383157",  
        "relatedIdentifierType": "DOI",  
        "relationType": "IsDescribedBy"}],  
    "identifier": {  
        "identifier": "10.YET/UNASSIGNED",  
        "identifierType": "DOI"}  
},  
    "resourceType": {  
        "resourceTypeGeneral": "InteractiveResource"  
    }  
},
```

## DLHub Metadata

```
"dlhub": {  
    "version": "0.1",  
    "domains": [  
        "general",  
        "digit recognition"  
    ],  
    "visible_to": [  
        "public"  
    ],  
    "id": null,  
    "name": "mnist_tiny_example",  
    "files": {  
        "other": [],  
        "model": "model.hd5"  
    }  
},
```

### Access Control

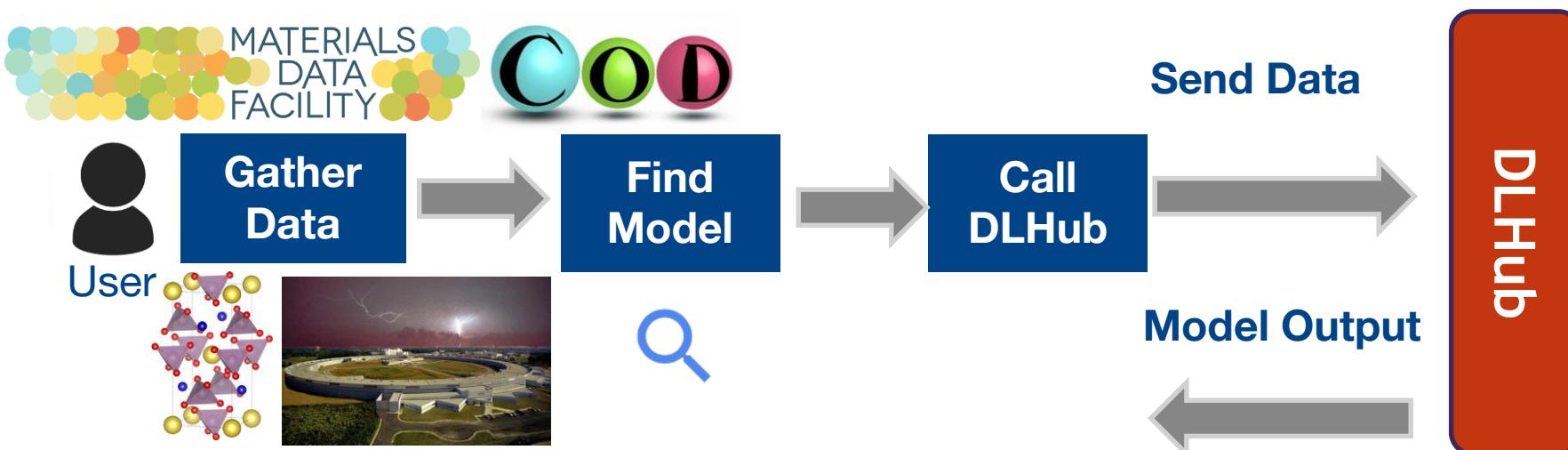
- **Public**
- **Globus users**
- **Globus groups**

## Servable Metadata

```
"servable": {  
    "methods": {  
        "run": {  
            "input": {  
                "type": "ndarray",  
                "description": "Image of a digit",  
                "shape": [null,28,28,1]  
            },  
            "output": {  
                "type": "ndarray",  
                "description": "Probabilities of being 0-9",  
                "shape": [null,10]  
            },  
            "parameters": {},  
            "method_details": {  
                "method_name": "predict",  
                "classes": ["0","1","2","3","4",  
                           "5","6","7","8","9"]  
            }  
        },  
        "shim": "keras.KerasServable",  
        "language": "python",  
        "dependencies": {  
            "python": {  
                "keras": "2.2.4",  
                "h5py": "2.8.0"  
            }  
        },  
        "type": "Keras Model",  
        "model_type": "Deep NN"  
    }  
},
```

# DLHub Model Discovery and Usage

- Find curated and tested models
- Use models through simple interfaces

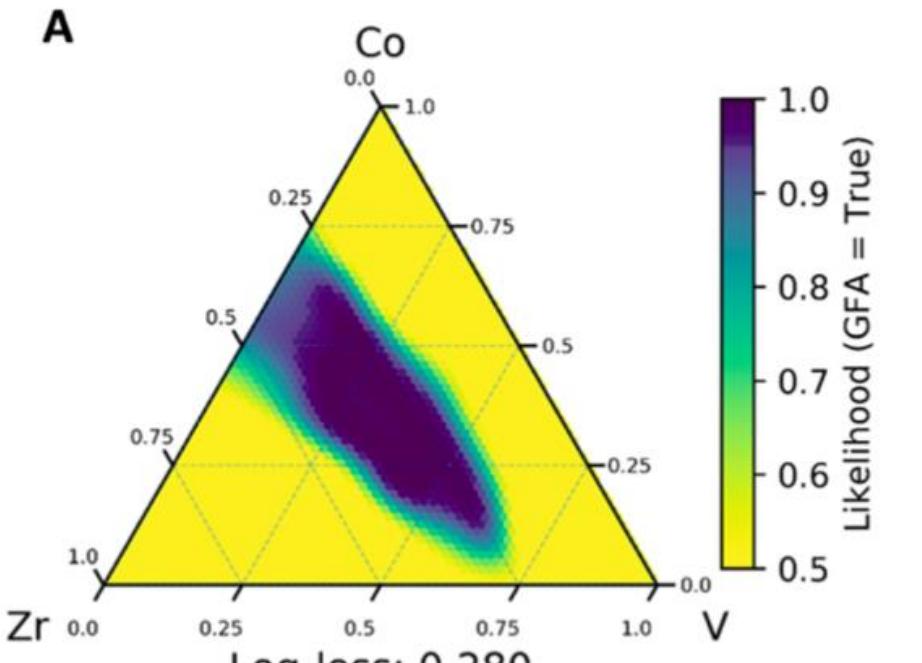


# Predicting Glass-forming Ability

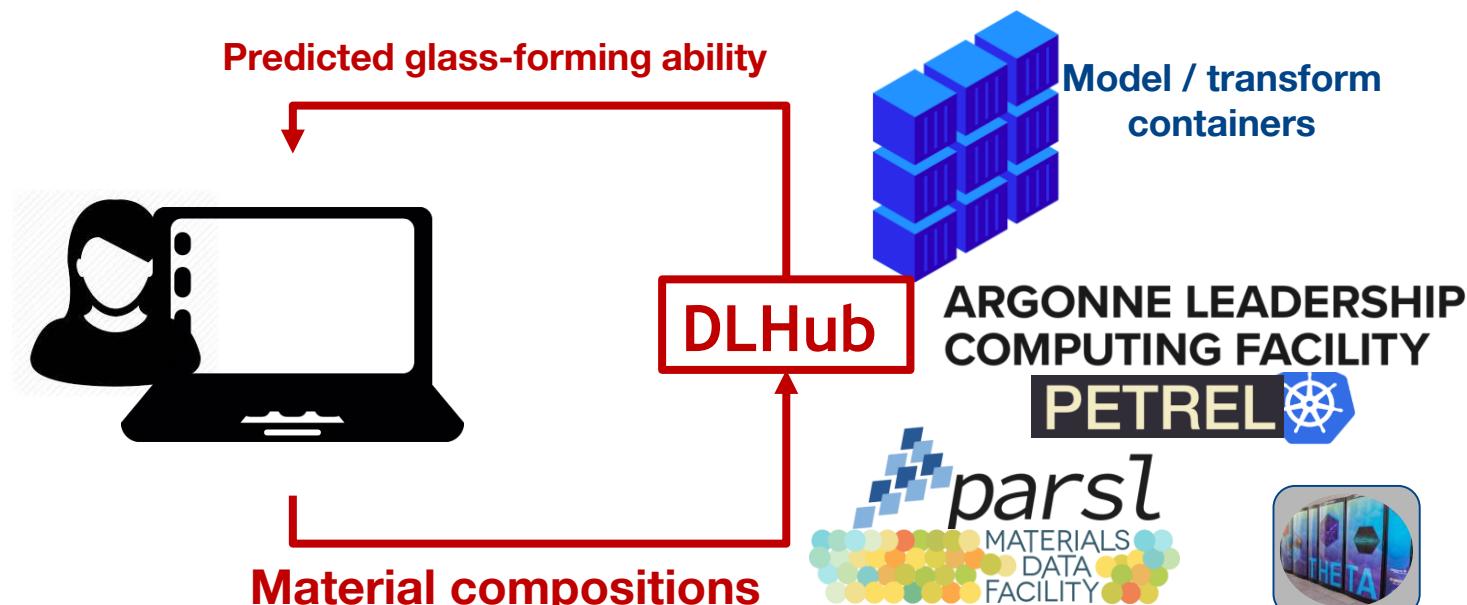
Accelerated discovery of metallic glasses through iteration of machine learning and high-throughput experiments

Fang Ren<sup>1,\*</sup>, Logan Ward<sup>2,3,\*</sup>, Travis Williams<sup>4</sup>, Kevin J. Laws<sup>5</sup>, Christopher Wolverton<sup>2</sup>, Jason Hattrick-Simpers<sup>6</sup> and Apurva Mehta<sup>1,†</sup>

10.1126/sciadv.aaq1566



- Where are the model and trained weights?
- How do I run the model on my data?
- How can I retrain the model on new data?
- How can I build on this work?



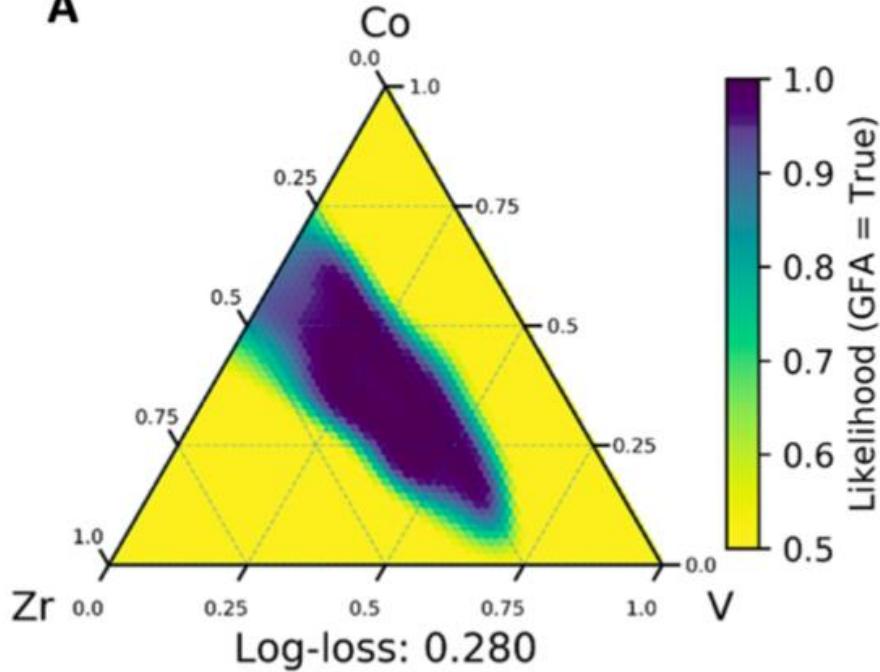
# Predicting Glass-forming Ability

Accelerated discovery of metallic glasses through iteration of machine learning and high-throughput experiments

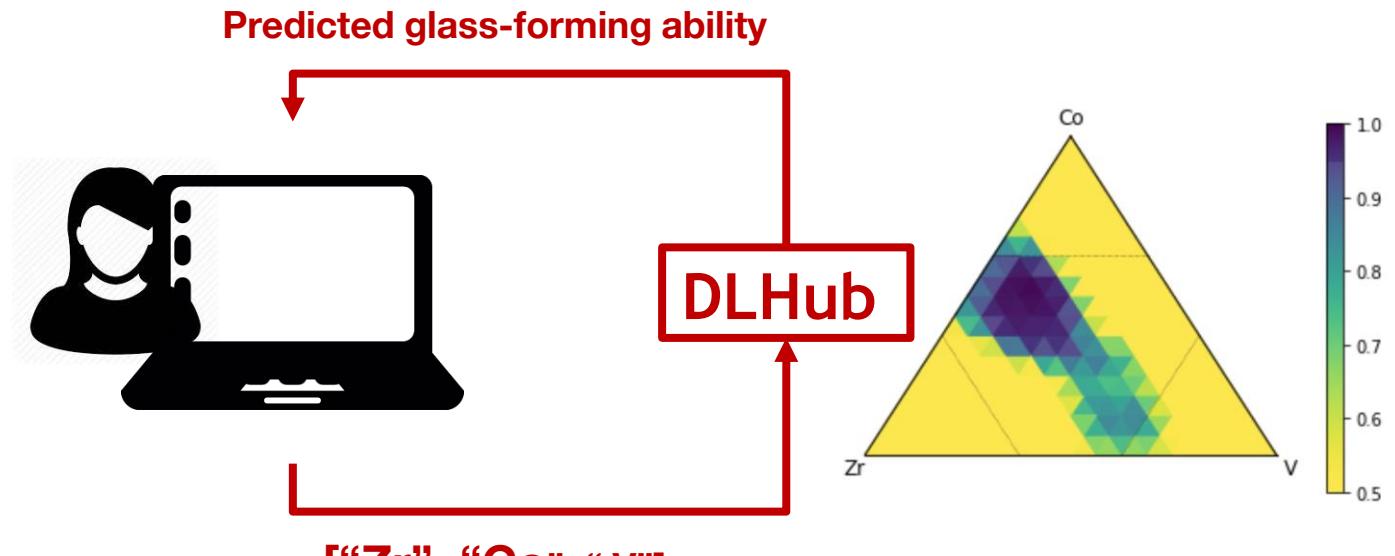
Fang Ren<sup>1,\*</sup>, Logan Ward<sup>2,3,\*</sup>, Travis Williams<sup>4</sup>, Kevin J. Laws<sup>5</sup>, Christopher Wolverton<sup>2</sup>, Jason Hattrick-Simpers<sup>6</sup> and Apurva Mehta<sup>1,†</sup>

10.1126/sciadv.aaq1566

A



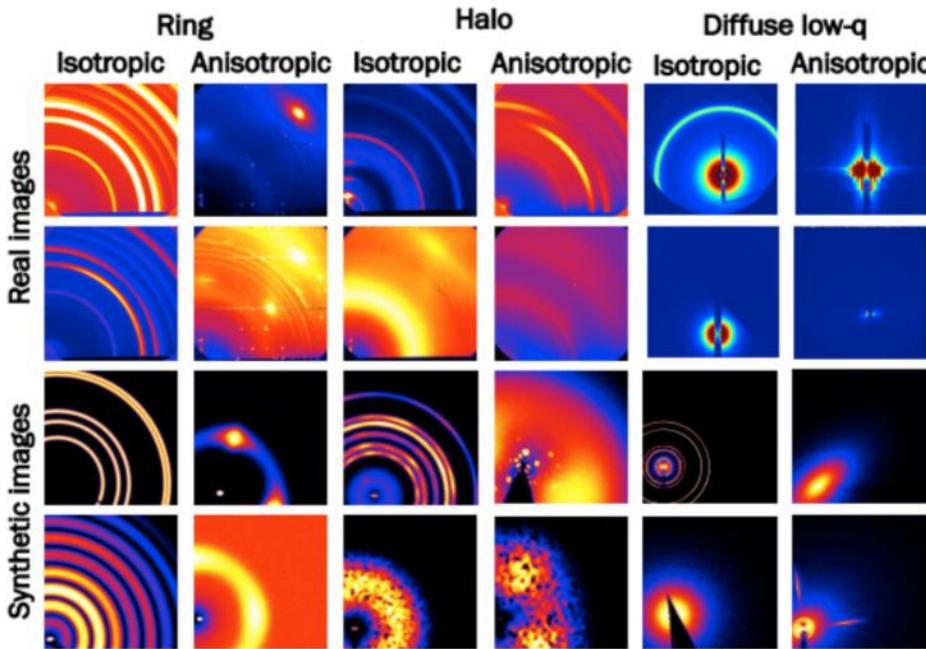
```
servable_name = "metallic_glass"  
servable_id = dl.get_id_by_name(servable_name)  
elems = ["V", "Co", "Zr"]  
  
res = dl.run(servable_id, {"data":elems})
```



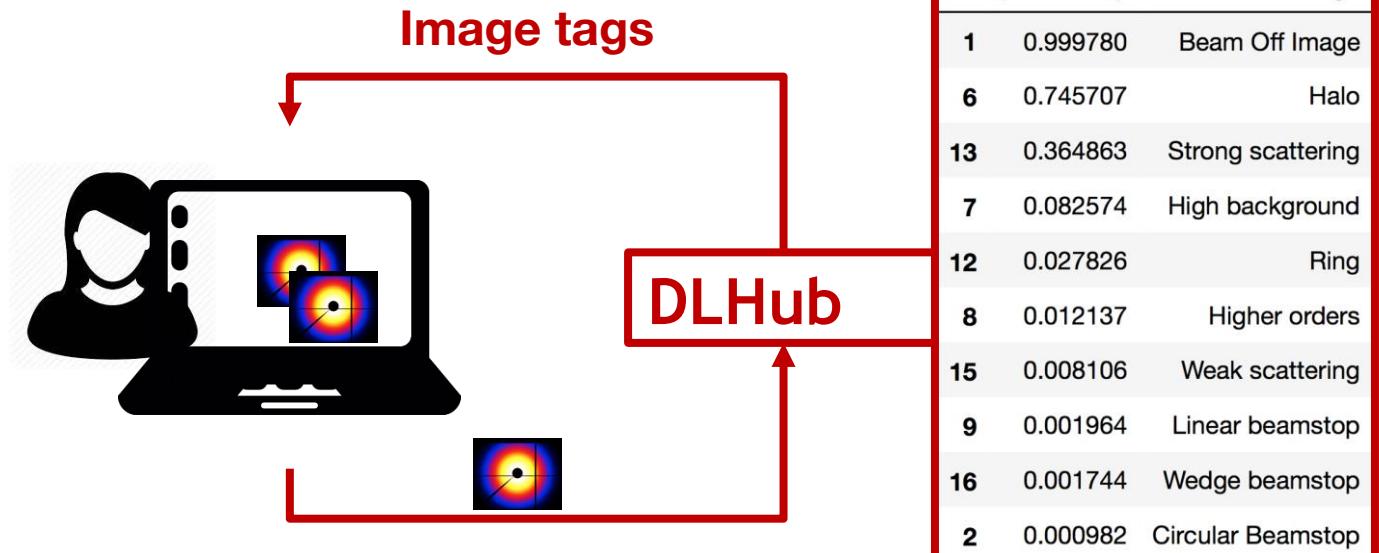
# Analyzing Beamline Images

Robust and Scalable Deep Learning for X-ray Synchrotron Image Analysis

Nicole Meister<sup>1\*</sup>, Ziqiao Guan<sup>2\*</sup>, Jinzhen Wang<sup>3</sup>, Ronald Lashley<sup>4</sup>,  
Jiliang Liu<sup>5</sup>, Julien Lhermitte<sup>5</sup>, Kevin Yager<sup>5</sup>, Hong Qin<sup>2</sup>, Bo Sun<sup>6</sup>, Dantong Yu<sup>3</sup>



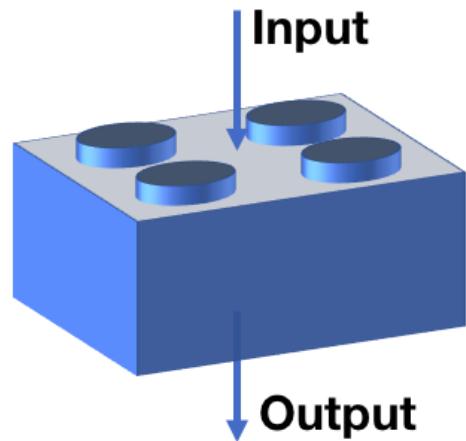
- Stage data into containers via Globus HTTPS
- Pass valid token and data location



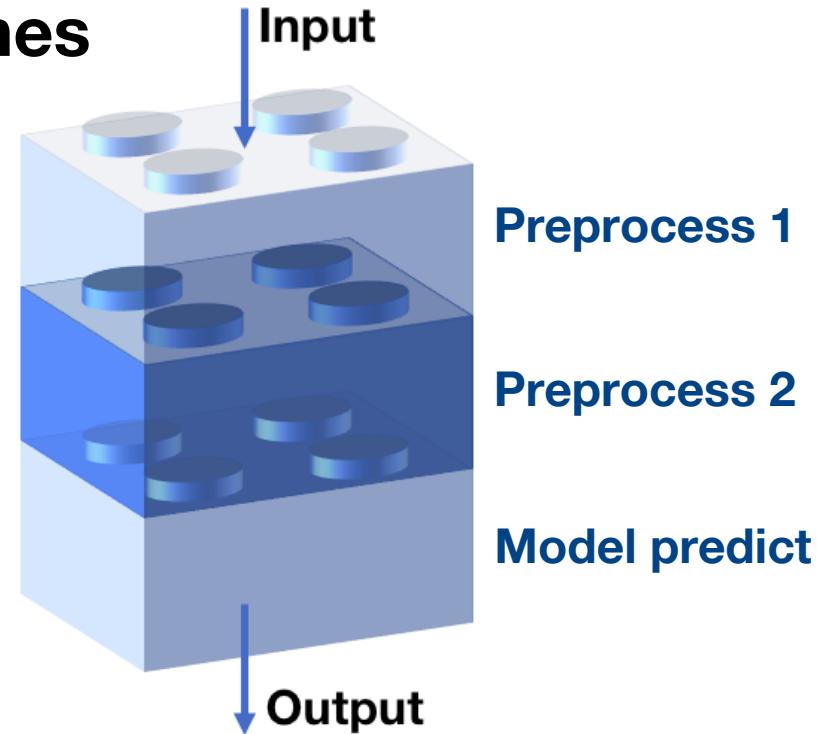
# Data and Learning Hub (DLHub): Pipelines



**Servables**

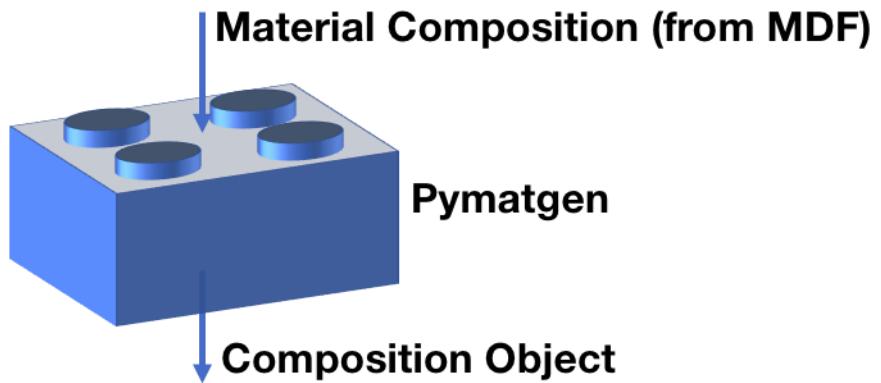


**Pipelines**

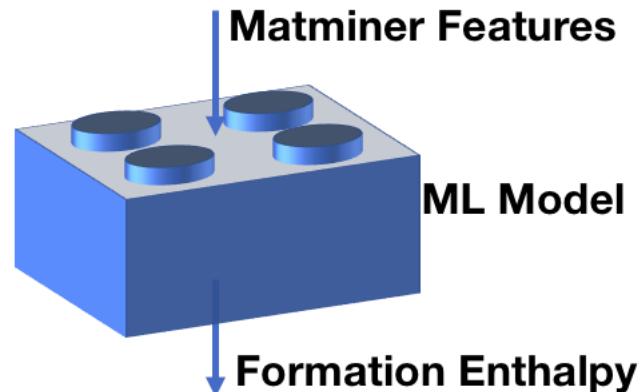


# Pipelines: Predicting Formation Enthalpy

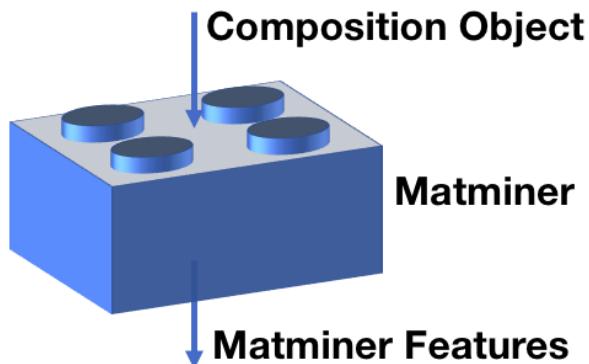
## Step 1



## Step 3

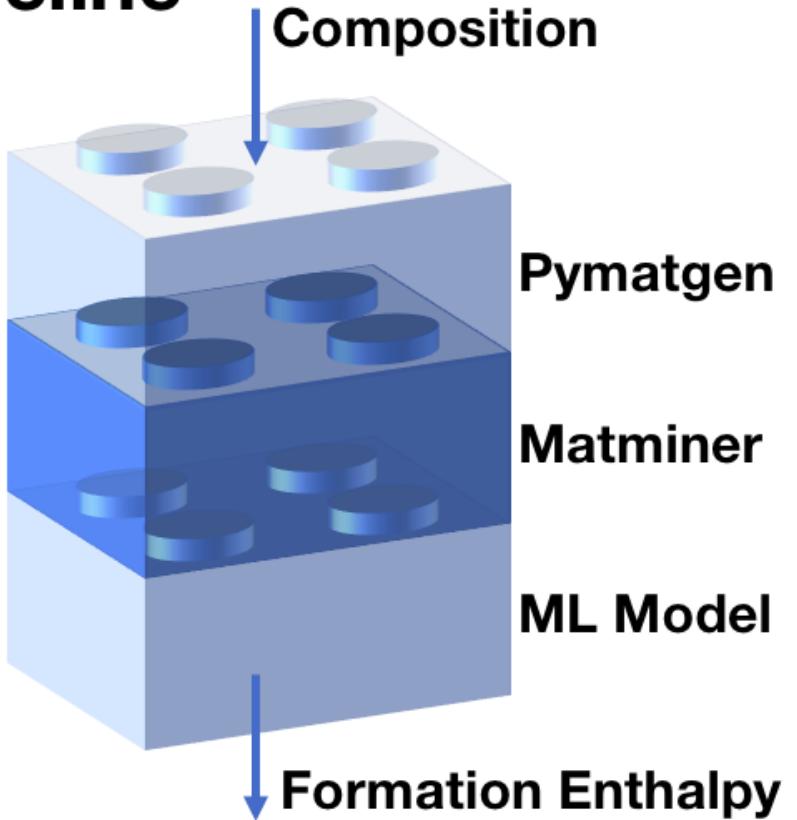


## Step 2



# Predicting Formation Enthalpy

**Pipeline**

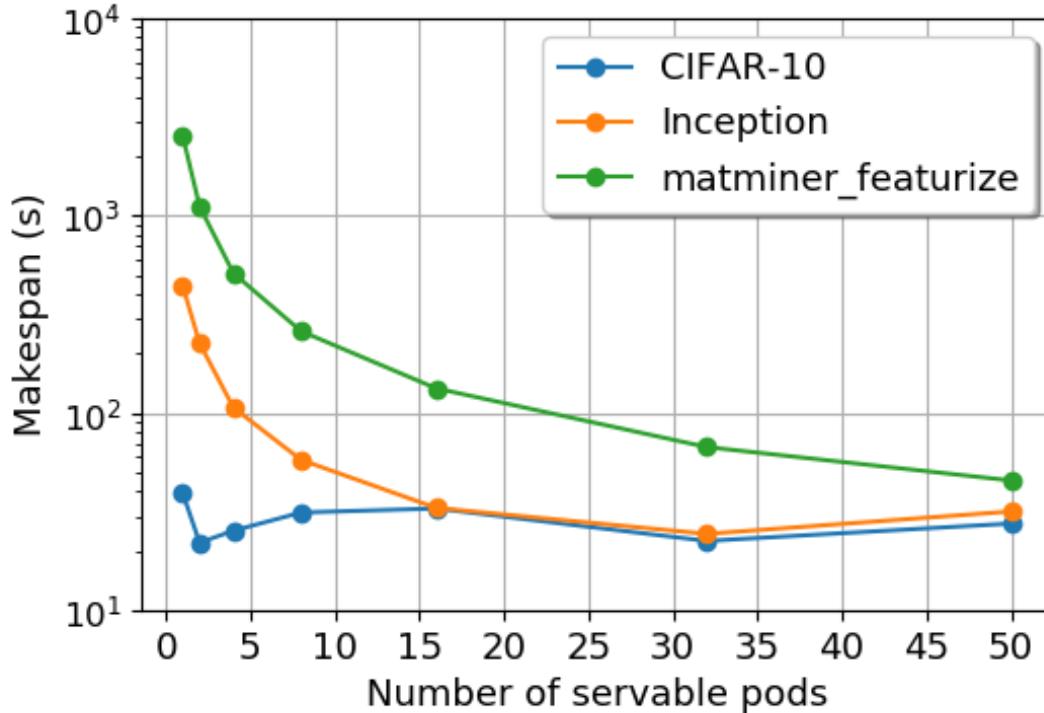


# DLHub

# Performance

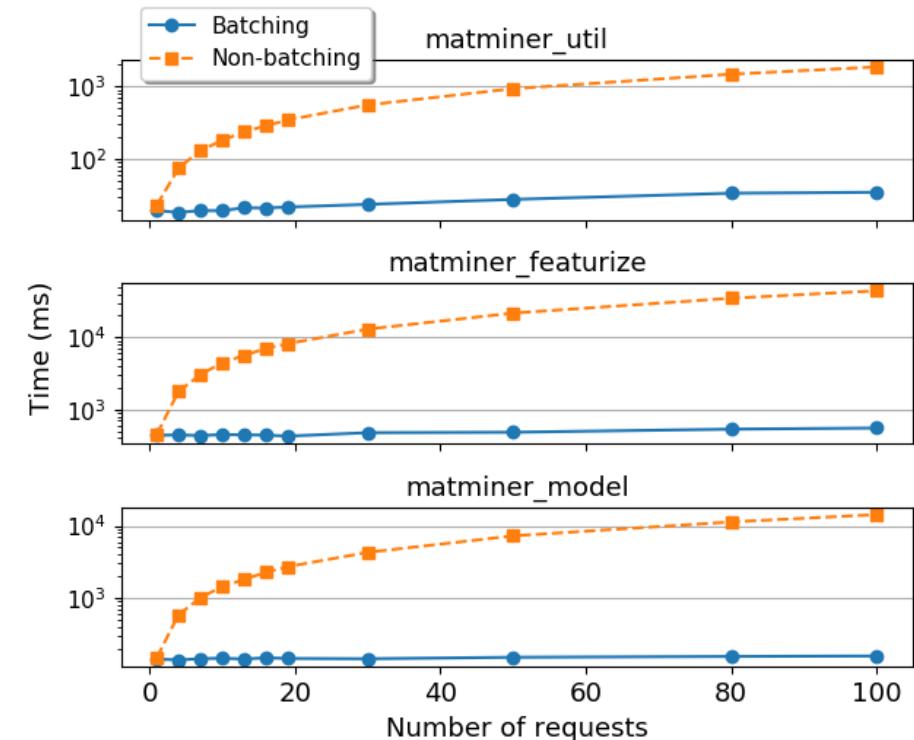
# DLHub Performance

## Scale Testing



The time required for the Inception, CIFAR10, and Matminer-featurize models to process 5000 inferences with varying numbers of replicas.

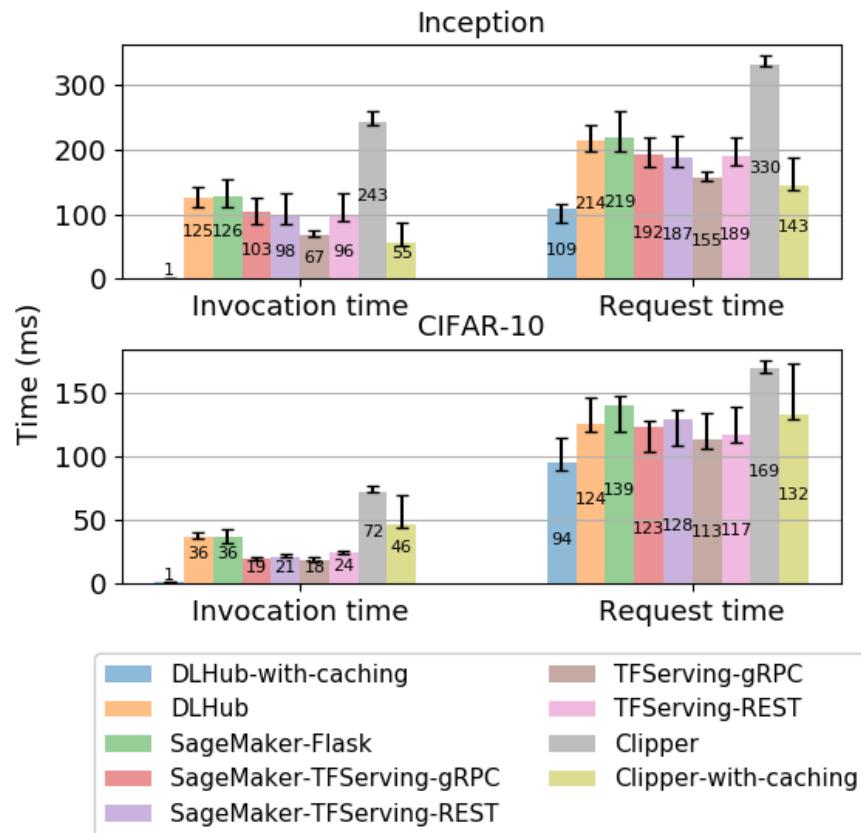
## Batching



Servable invocation time, with and without batching.

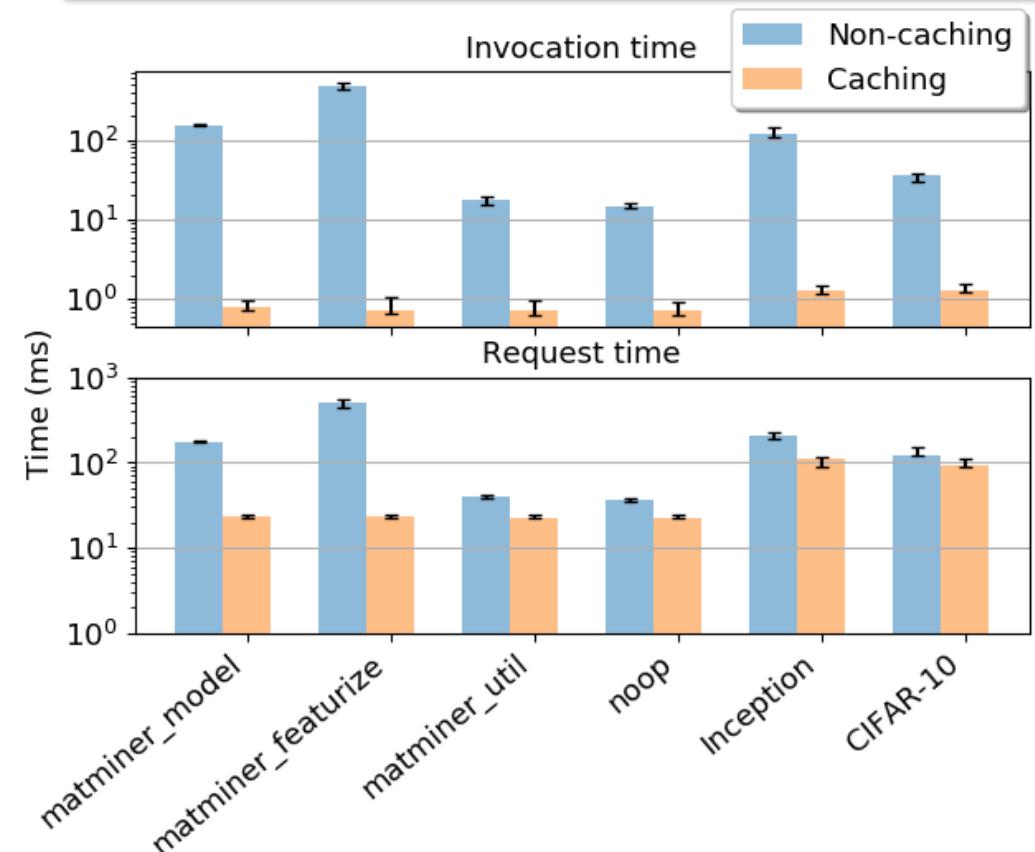
# DLHub Performance

## Serving General Models



Performance of different serving systems on the Inception and CIFAR-10 problems.

## Caching



Performance impact of caching in DLHub. Bars and error bars show median and 5th/95th percentiles

# DLHub Summary

---

## Model deposit and discovery

- Developed a model schema to promote discovery
- Implemented advanced search and filtering
- Built ingest flow: models are dynamically staged, packaged, dockerized, published, and indexed

## Model serving

- Deployed capabilities for users to run inference with SDK and CLI
- Automated testing of containers
- Implemented caching and batching

## Support for multiple execution sites

- PetrelKube: Parsl, TF serving, Sagemaker
- Other: AWS, OSG

## Authentication

- Protected model metadata and inference with GlobusAuth
- Secured data staging

## Monitoring and statistics

- Request, invocation, data staging

## Future work

- Dynamic scaling by load
- Build Web UI to create pipelines and invoke models
- Cache at the servable level within pipelines
- Couple DLHub to data sources (MDF, etc.)
- Integrate with ML frontend tools (DeepForge), optimization tools (DeepHyper), and more
- Create interface for training and retraining of models

# Thanks to our sponsors!

## ARGONNE LEADERSHIP COMPUTING FACILITY

ALCF DF

PETREL

Data Management and Sharing Pilot

Parsl



Globus



U. S. DEPARTMENT OF  
ENERGY

NIST

NATIONAL INSTITUTES  
OF HEALTH



MATERIALS  
DATA  
FACILITY

NIST

IMaD



CHIMaD  
Center for Hierarchical Materials Design

DLHub



U. S. DEPARTMENT OF  
ENERGY

Argonne  
LDRD



U. S. DEPARTMENT OF  
ENERGY

Argonne  
NATIONAL LABORATORY



THE UNIVERSITY OF  
CHICAGO