



TOYOTA
RESEARCH INSTITUTE

Extracting Insights from Atomistic and Spectroscopic Materials Data

Steven Torrisi
Sr. Research Scientist
7/10/25, NIST

Talk Outline

1

TRI + AMDD Background

2

Challenges in AI-guided materials design

3

Characterization & bridging experiment & theory

Toyota Research Institute



**Energy &
Materials**

**Human
Centered AI**

**Human
Interactive
Driving**

Robotics

HQ
Los Altos, CA

CAM
Cambridge, MA



 **TOYOTA**

 **LEXUS**

**Accelerated
Materials Design
& Discovery**

**Energy &
Materials**

**Carbon
Neutrality
Strategy**



Leadership



Brian
Storey



Joey
Montoya



Linda
Hung



Jens
Bakander

Highlighted in this Talk



Steven
Torrisi



Weike Ye

Accelerated Materials Design & Discovery

AMDD



Hisatugu
Yamasaki



Leena
Sansguiri



Michaela
Burke-
Stevens



Daniel
Schweigert



Amanda
Volk



Jith
Subramanian



Amalie
Trewartha



Kevin Tran



Santosh
Suram



Koki
Nakano



Materials on all length scales

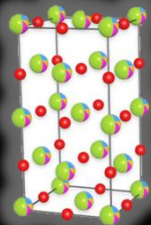
Atomistic

Molecular

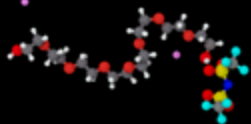
Periodic
(Spectra)

Device-Level

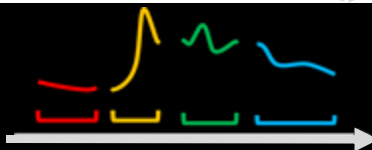
Disordered Rock-Salt
Structures



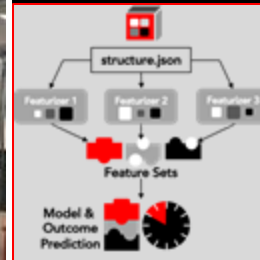
Transition
Paths



Polymer
Trajectories



XANES
Spectra
Featurization &
Inverse Modeling



Battery Charge &
Discharge Features

[1] Liu... **Torrìsì***, Wolverton* *et al*,
Submitted 2025, arXiv

[2] Sheriff, Freitas, Trewartha, **Torrìsì**, NeurIPS AI4Mat 2024

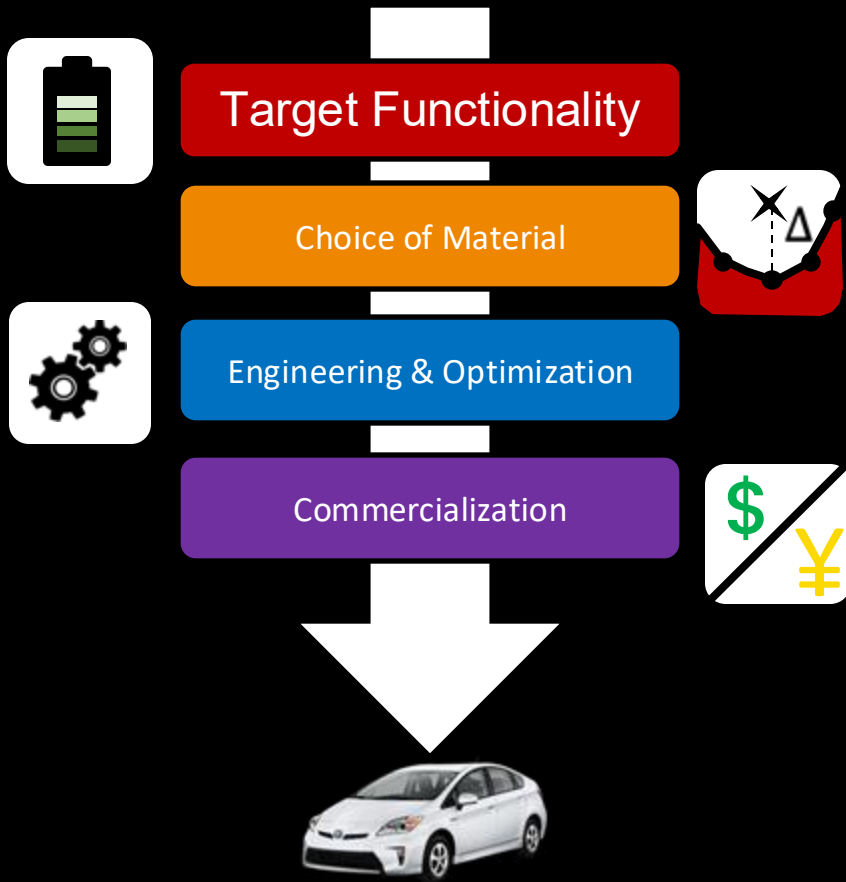
[3] Khajeh, Schweigert, **Torrìsì et al**. Macromolecules (2023)

[4] **Torrìsì et al**. NPJ Computational Materials, 2020

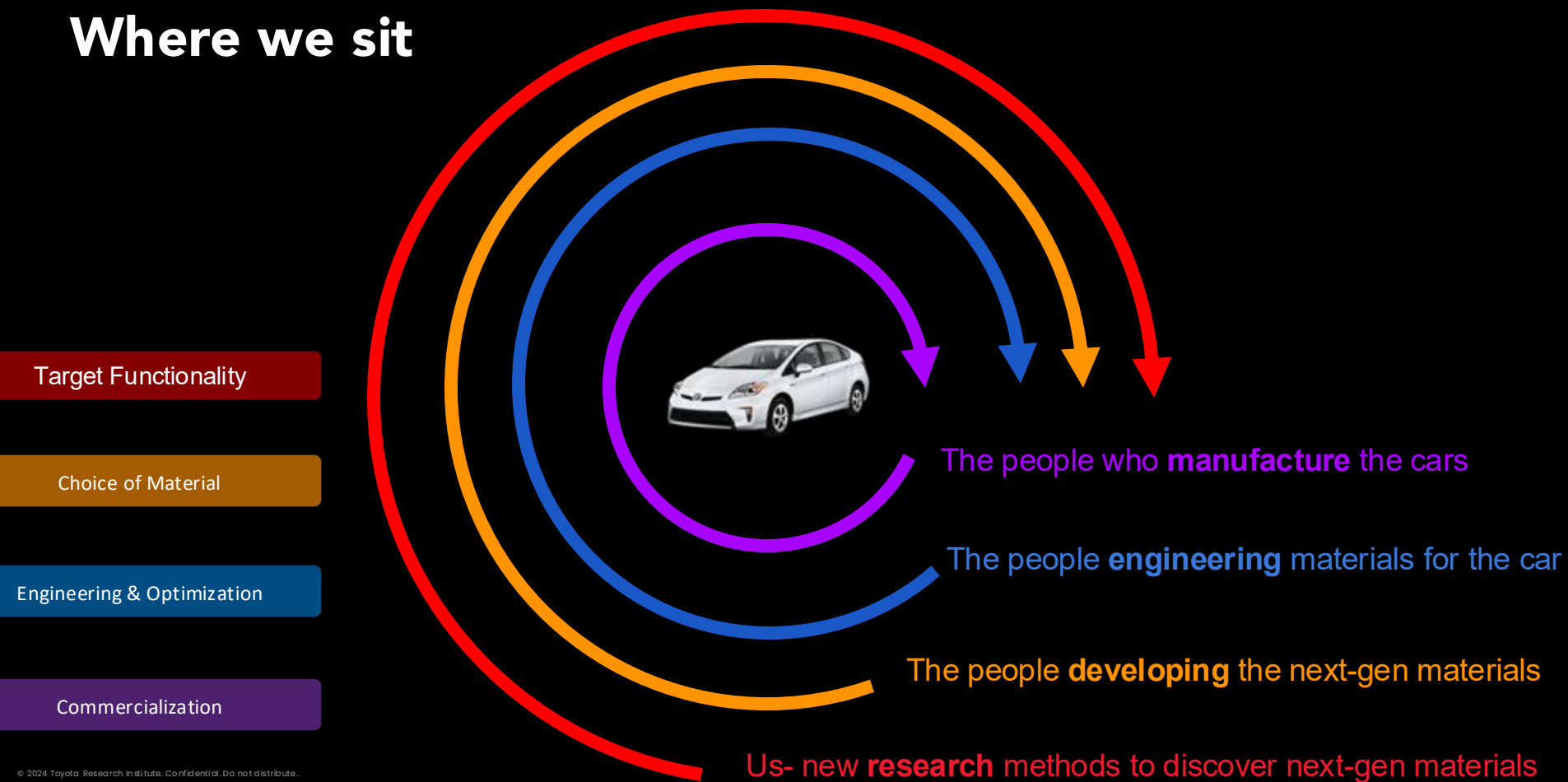
[5] Montoya, Aykol... **Torrìsì**, Trewartha, Storey, Appl. Phys. Rev 2022

[6] Ansari, **Torrìsì**, Trewartha, Sun, J. Energy Storage 2024

Materials Discovery in Context



Where we sit



Talk Outline

1

TRI + AMDD Background

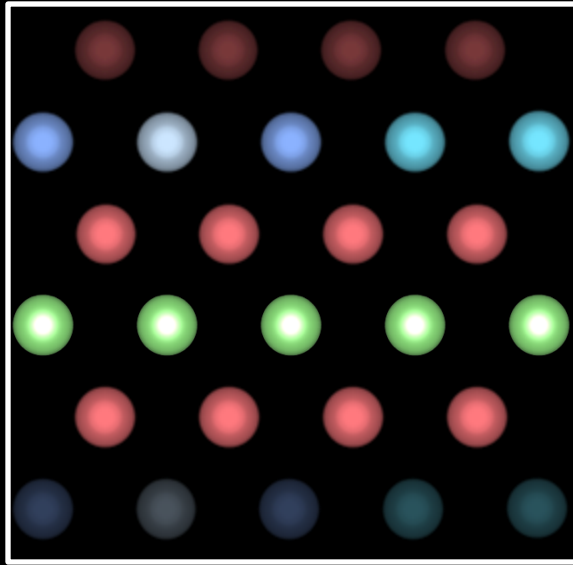
2

Challenges in AI-guided materials design

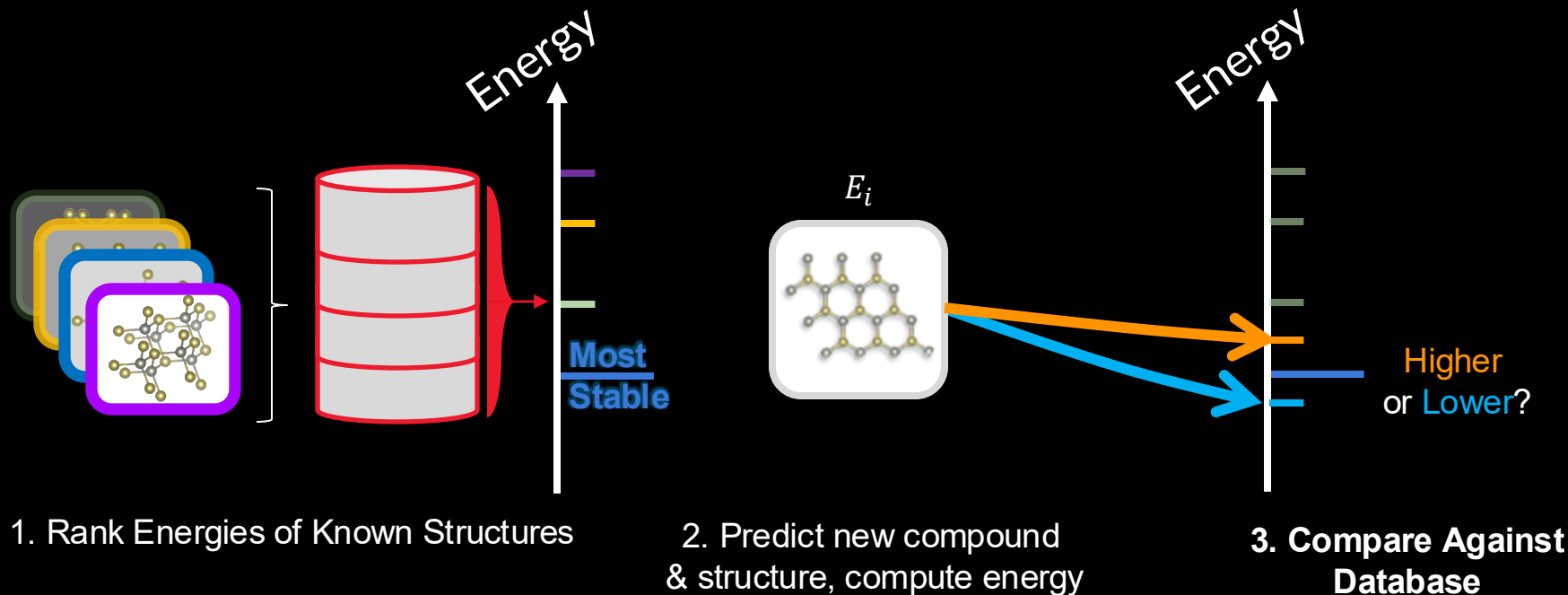
3

Characterization & bridging experiment & theory

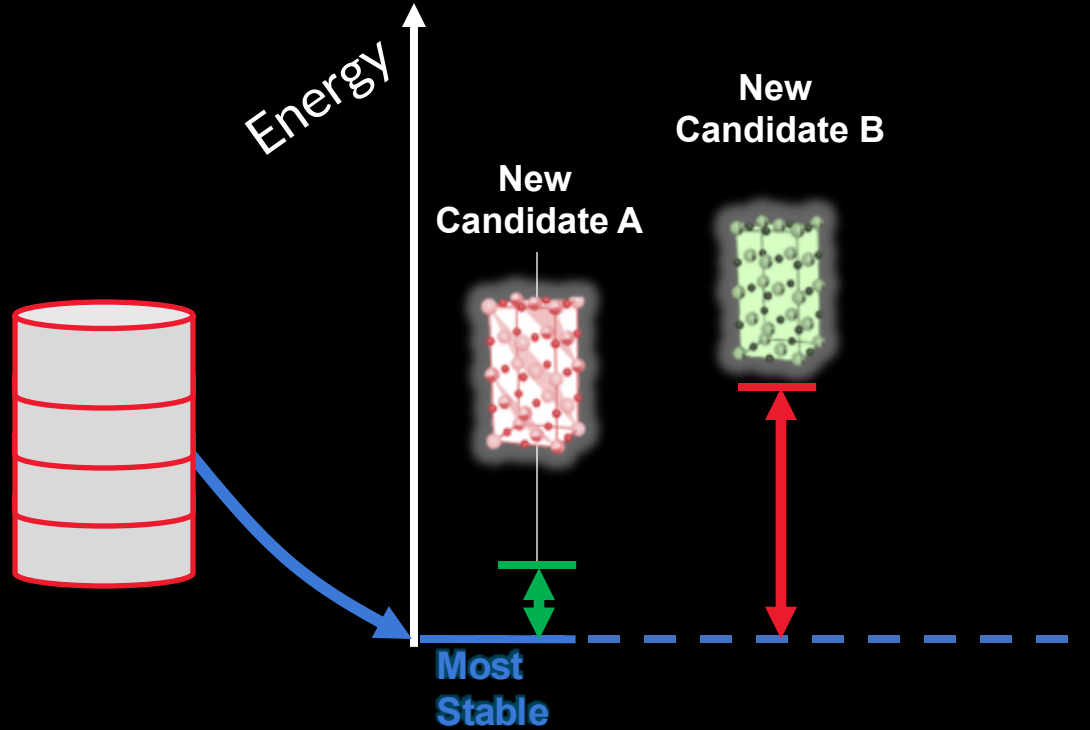
What goes into computational materials discovery?



The Current Paradigm of (inorganic, crystalline) Computational Materials Discovery



How we use energy-ranking



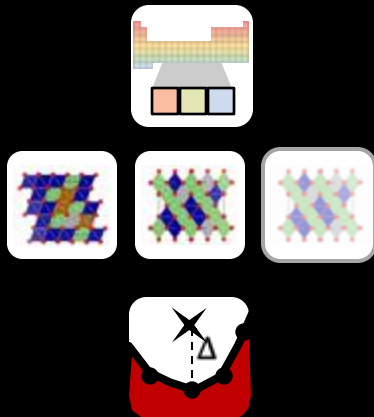
In general: The lower the **energy difference**, the better

Choice of Material

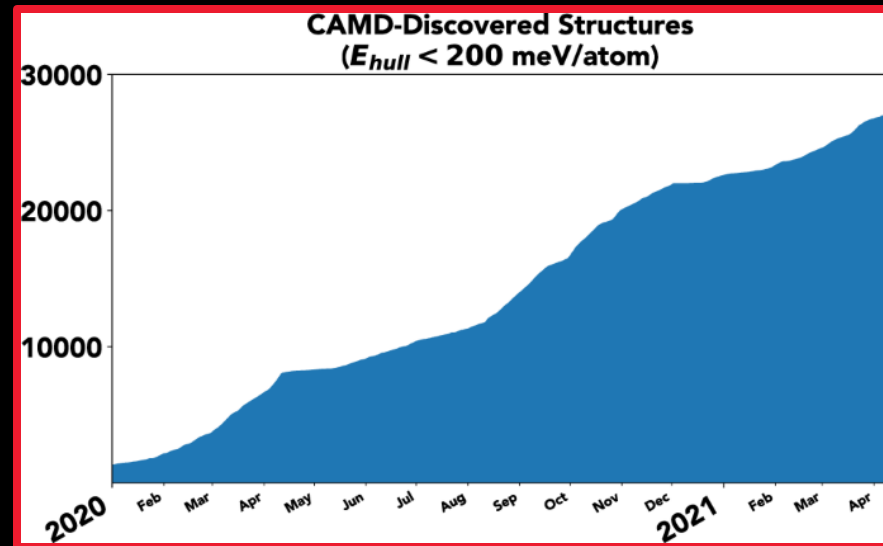
Choice of
Chemical System

Candidate Phases

Feasibility Screening
(e.g. Thermodynamics)



CAMD (2020, TRI)



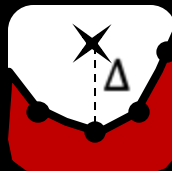
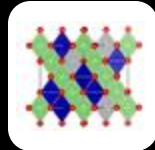
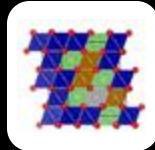
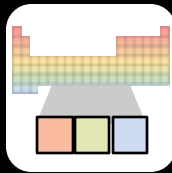
Weike Ye, Ray Lei,
Joseph Montoya, Murat Aykol

Choice of Material

Choice of
Chemical System

Candidate Phases

Feasibility Screening
(e.g. Thermodynamics)



nature computational science

Article | Published: 28 November 2022

A universal graph deep learning interatomic potential for the periodic table

Chi Chen & Shyue Ping Ong

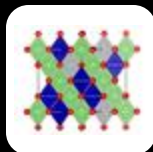
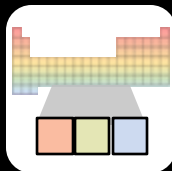
“About **1.8 million materials** were identified from a screening of 31 million hypothetical crystal structures to be potentially stable against existing Materials Project crystals based on M3GNet energies. “

Choice of Material

Choice of
Chemical System

Candidate Phases

Feasibility Screening
(e.g. Thermodynamics)



Article

Scaling deep learning for materials discovery

<https://doi.org/10.1038/s41586-023-06735-9>

Received: 8 May 2023

Accepted: 10 October 2023

Published online: 29 November 2023

Amil Merchant^{1,3,6,7}, Simon Batzner^{1,3}, Samuel S. Schoenholz^{1,3}, Muratahan Aykol¹,
Gwooon Cheon² & Ekin Dogus Cubuk^{1,3,6,7}

Novel functional materials enable fundamental breakthroughs across technological applications from clean energy to information processing^{1–11}. From microchips to

**“2.2 million structures below
the current convex hull...”**

From the perspective:

*“Systems such as GNoME can make many more computational predictions than even an autonomous lab can keep up with, says Andy Cooper, **“What we really need is computation that tells us what to make,”** Cooper says. For that, AI systems will have to accurately calculate a lot more of the predicted materials’ chemical and physical properties.”*

Choice of Material

Finding enthalpically stable ordered structures is now well trailblazed!

Candidate Phases

Feasibility Screening
(e.g. Thermodynamics)



Article

Scaling deep learning for materials discovery

<https://doi.org/10.1038/s41586-023-06735-9>

Received: 8 May 2023

Accepted: 10 October 2023

Amil Merchant^{1,2,3}, Simon Batzner^{1,4}, Samuel S. Schoenholz^{1,2}, Muratahan Aykol¹,
Gwooon Cheon² & Ekin Dogus Cubuk^{1,2,3}

Novel functional materials enable fundamental breakthroughs across technological domains, from information processing^{1–10}. From microchips to

From the perspective:

“Still, it’s clear that systems such as GNoME can make many more computational predictions than even an autonomous lab can keep up with, says Andy Cooper, academic director of the Materials Innovation Factory at the University of Liverpool, UK. **“What we really need is computation that tells us what to make,”** Cooper says. For that, AI systems will have to accurately calculate a lot more of the predicted materials’ chemical and physical properties.”

Choice of Material

But, enthalpically stable
structures are

not the end of the story.

There's some fine print...

Article

Scaling deep learning for materials discovery

<https://doi.org/10.1038/s41586-023-06735-9>

Received: 8 May 2023

Accepted: 10 October 2023

Amil Merchant^{1,2,3}, Simon Batzner^{1,4}, Samuel S. Schoenholz^{1,2}, Muratahan Aykol¹,
Gwooon Cheon² & Ekin Dogus Cubuk^{1,2,3}

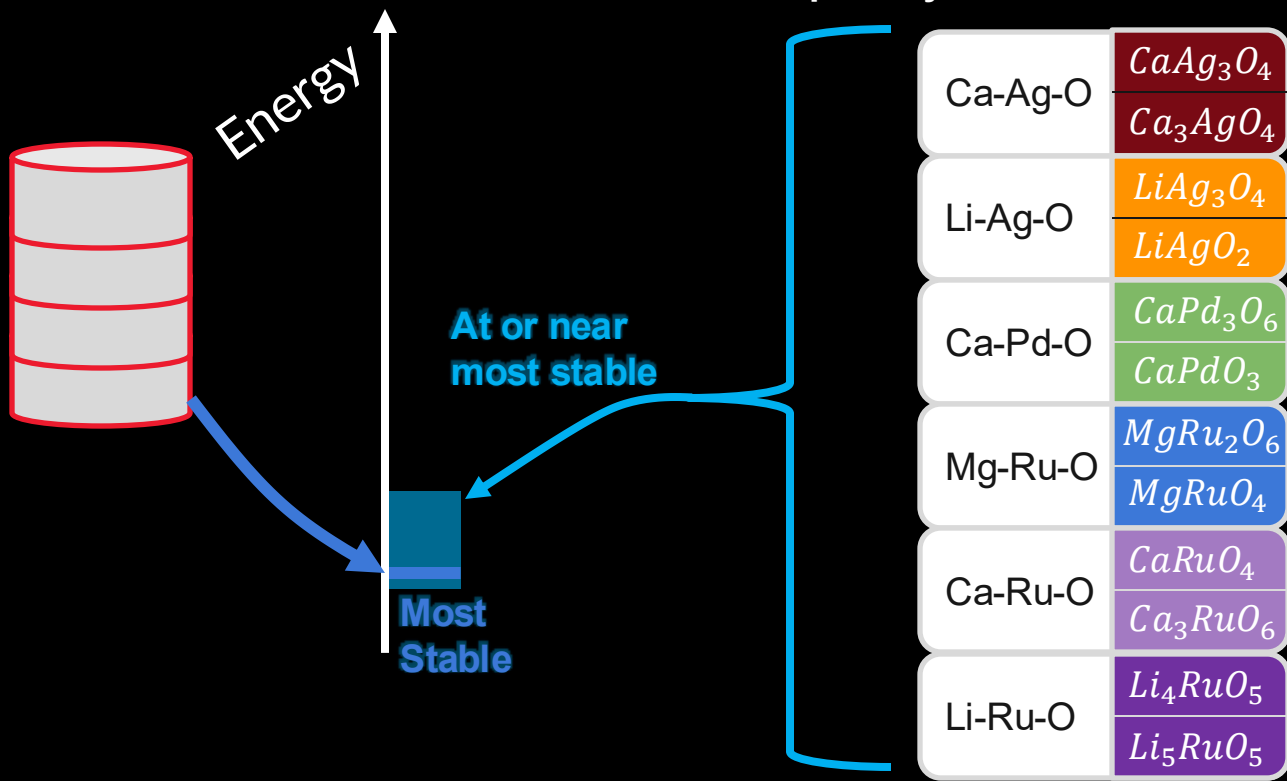
Deep learning has enabled fundamental breakthroughs across technological domains, from microchips to

crystals identified in
7, improved efficiency
ery of 2.2 million
urrent convex hull..."

GNOME can make many
an even an autonomous
oper, academic director
at the University of
ed is computation
oper says. For that, AI
systems will have to accurately calculate a lot more of the
predicted materials' chemical and physical properties."

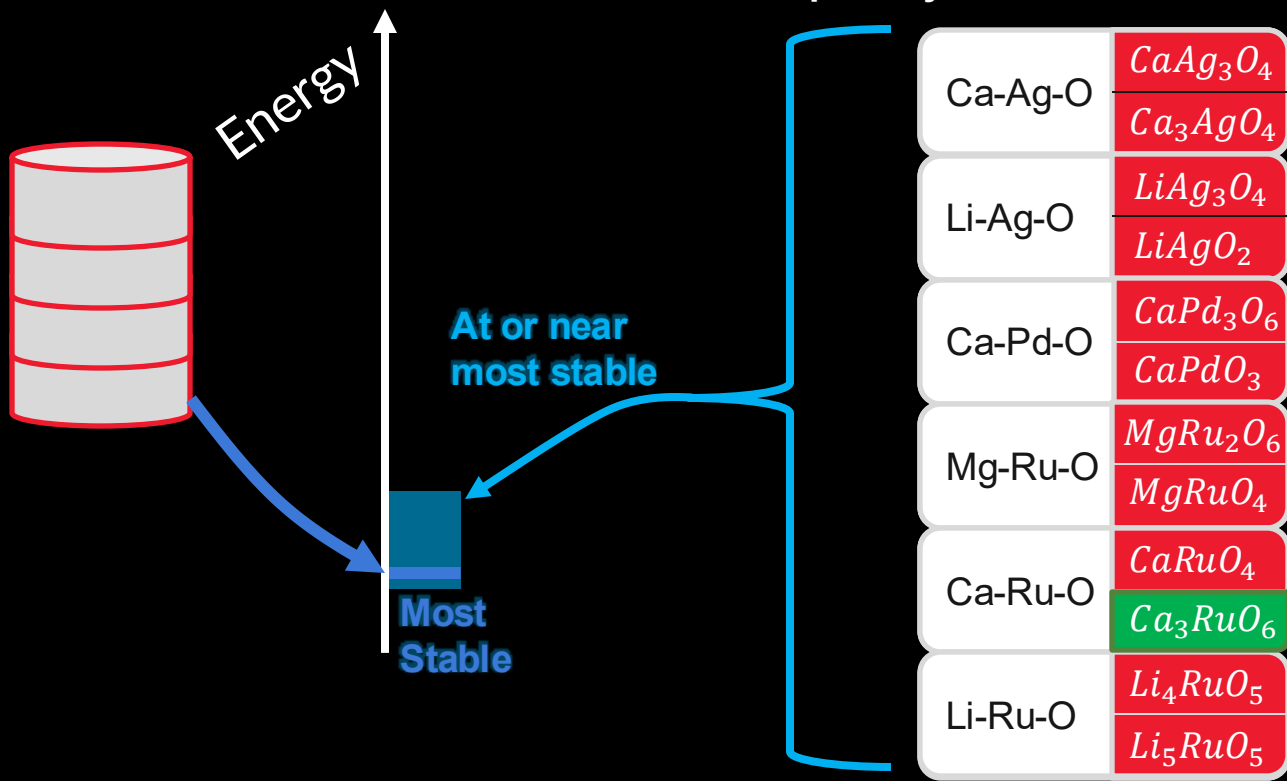
Previous work highlighted synthesis challenges

12 Attempted Syntheses based on Energy Metric



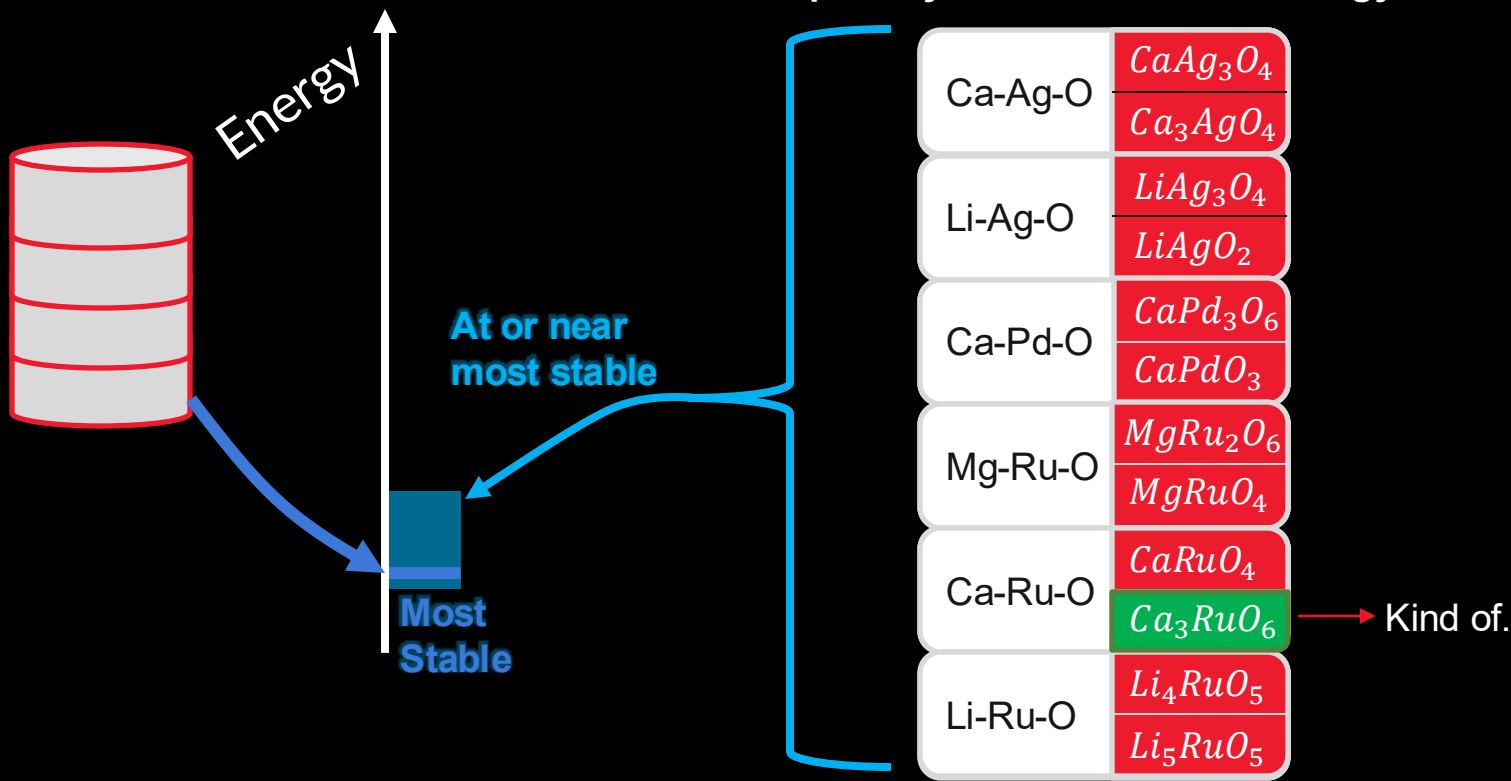
Previous work highlighted synthesis challenges

12 Attempted Syntheses based on Energy Metric



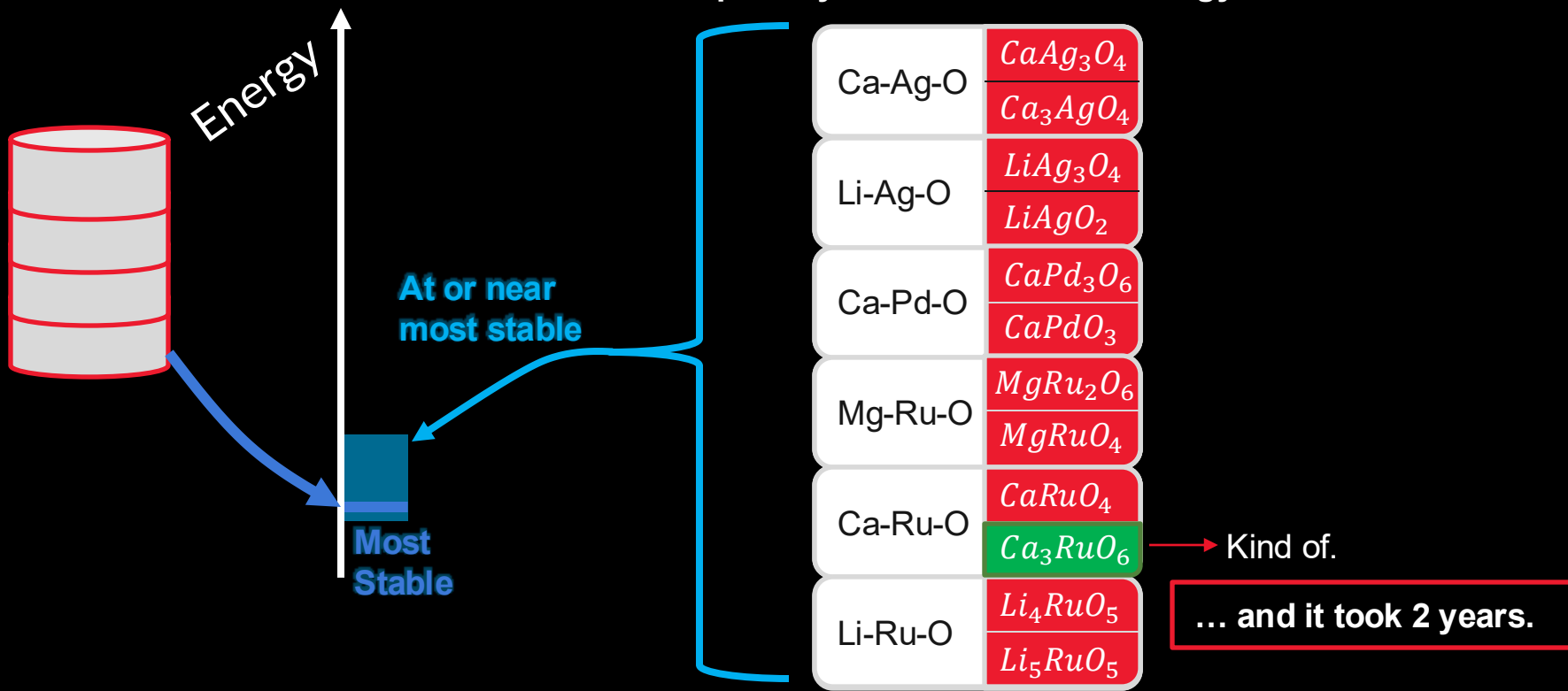
Previous work highlighted synthesis challenges

12 Attempted Syntheses based on Energy Metric



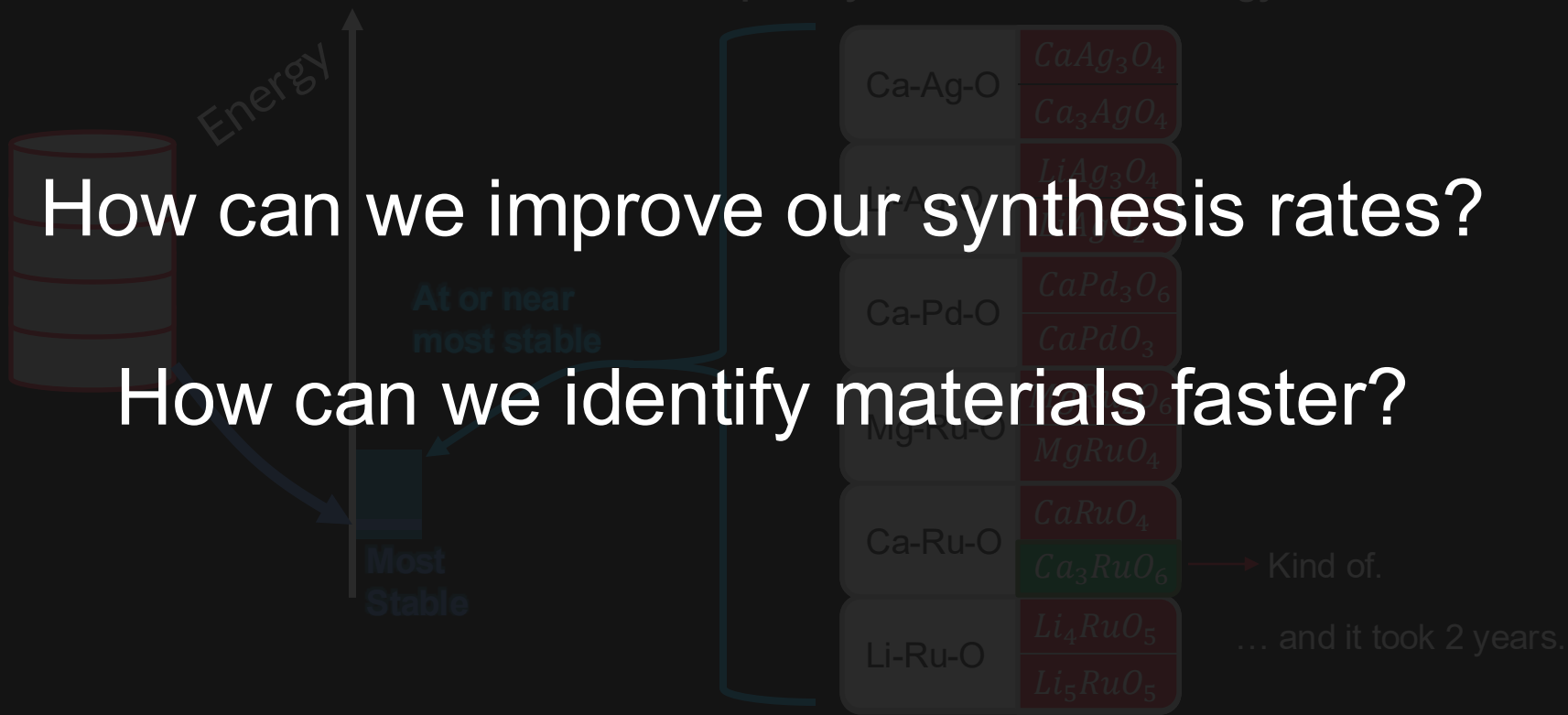
Previous work highlighted synthesis challenges

12 Attempted Syntheses based on Energy Metric



Previous work highlighted synthesis challenges

12 Attempted Syntheses based on Energy Metric



3 challenges for our team

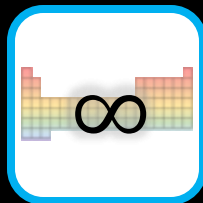
We don't know how to make new candidates.



Identifying materials is hard & involves diverse data.



We have ∞ candidates.



3 challenges for our team

We don't know how to make new candidates.



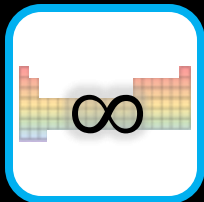
Develop theories of inorganic synthesis

Identifying materials is hard & involves diverse data.



Tools that make it easier to interface with experiment

We have ∞ candidates.



Tools that accelerate atomistic modeling

Talk Outline

1

TRI Background &
Challenges in AI-guided materials design

2

Challenges in AI-guided materials design

3

Characterization & bridging experiment & theory

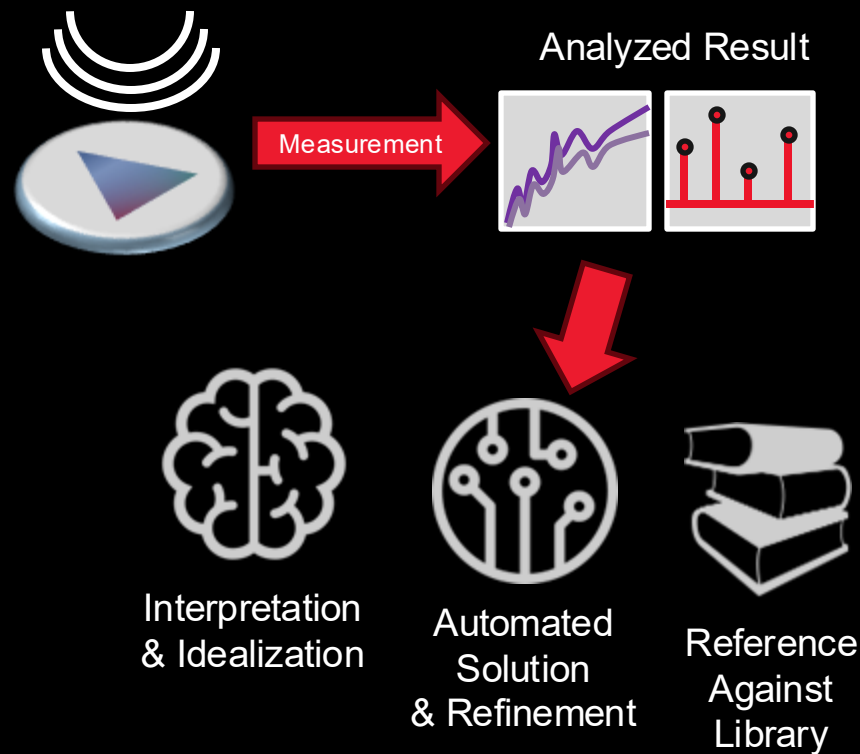
Working with Spectral Data

Multi-modal Analysis & Spectra-scope



Characterization is a Bottleneck for High Throughput Materials Discovery

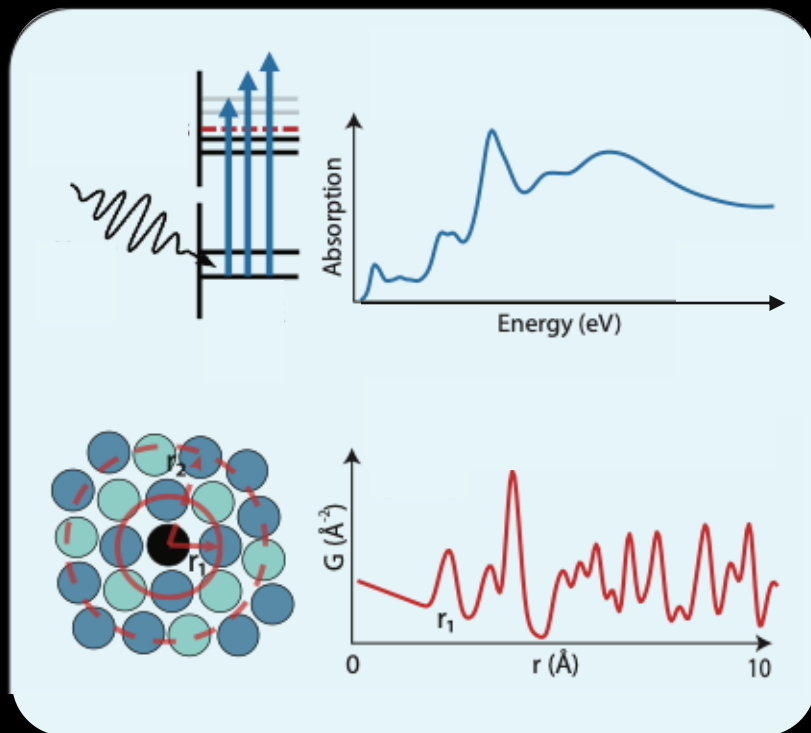
- Low-dimensional forms of characterization (e.g. XRD) can be accelerated
 - but may still require human verification
- Our research program assumes that humans will remain in the loop for the near future
- **Goal is to study interoperability & interpretability of spectra**



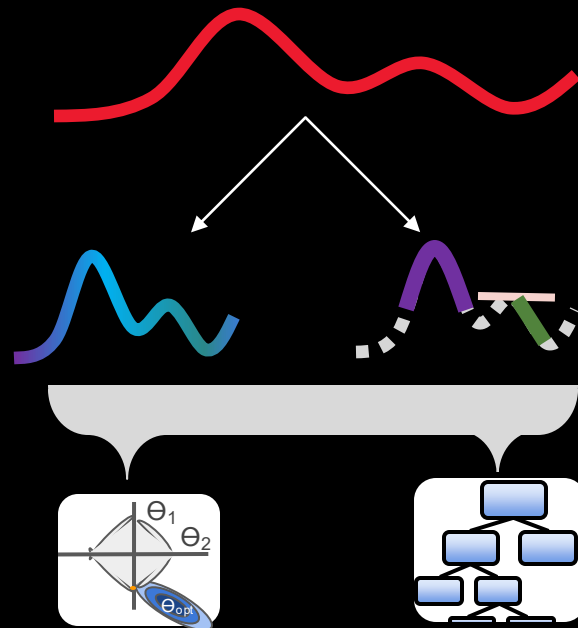
Guiding Questions

- How can we make it easier to utilize and combine diverse modes of spectra?
- How can we make it easier to use ML/statistical methods in the limited-data regime?
- How can extant databases of materials & spectra be utilized for optimal design of experiments?

Multimodal Studies of XANES vs. PDFs



Spectra-scope: A tool for
feature discovery &
downselection of spectrum-
property relationships

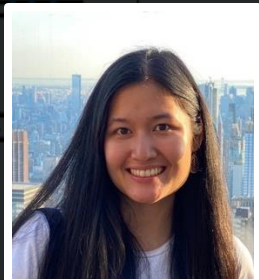


Multimodal Studies of XANES vs. PDFs

Featuring work with Columbia University



Zoe Zachko



Dr. Tina Na Narong



Prof. Simon Billinge

NPJ Computational Materials, 2025
Na Narong, Zachko, **Torrissi***, Billinge*

Spectra-scope: A tool for feature discovery & downselection of spectrum- property relationships

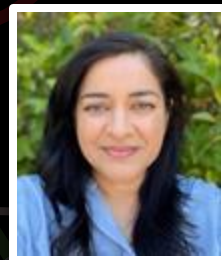
And with TRI intern + scientist



Amalya Johnson
(Stanford)



Dr. Weike Ye

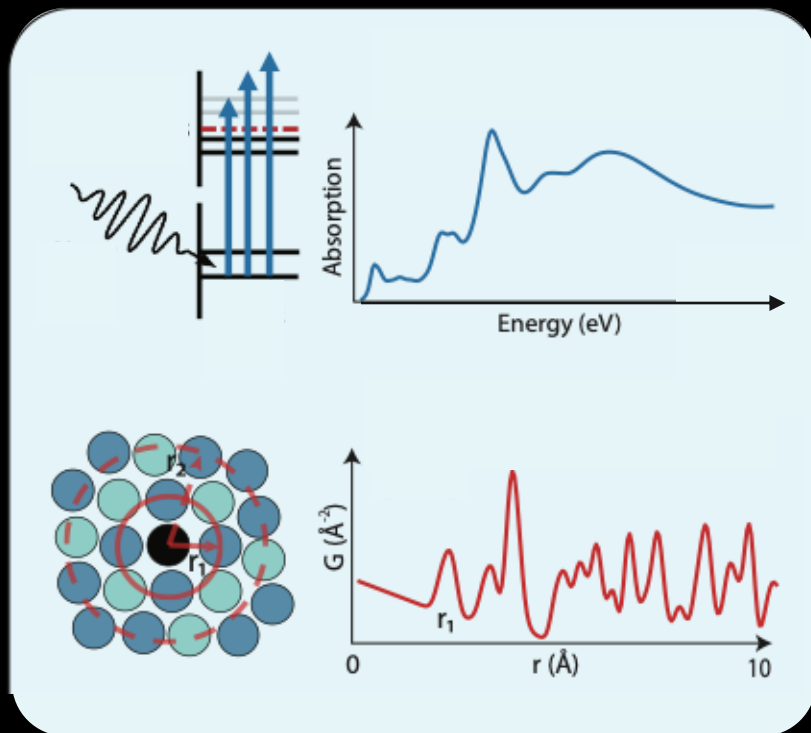


Leena
Sansguiri

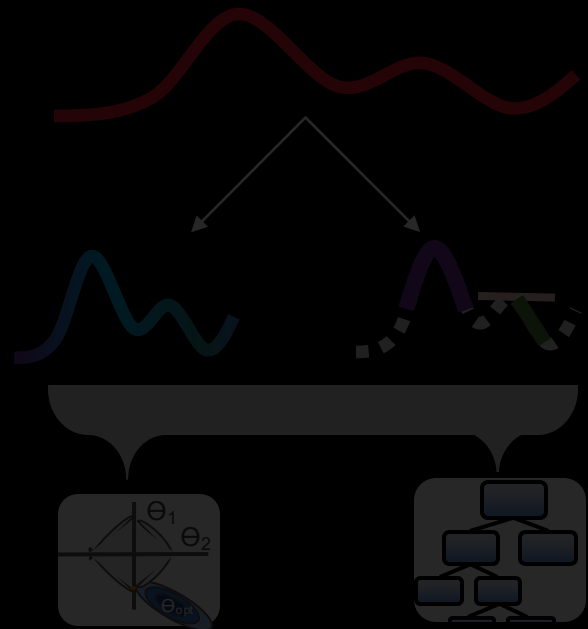


In prep. Johnson, Ye, **Torrissi**

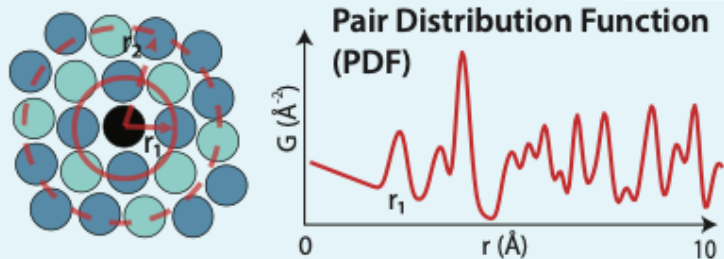
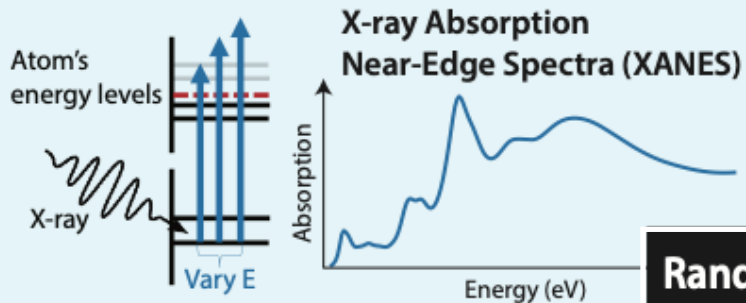
Multimodal Studies of XANES vs. PDF



Spectra-scope: A tool for
feature discovery &
downselection of spectrum-
property relationships



Data



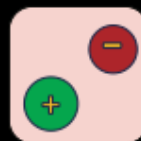
Structures: Transition metal oxides
from Materials Project database



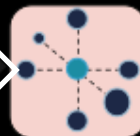
Models

Prediction targets:

Metal's local atomic environments



Charge (oxidation) state
Ti, Mn, Fe: 2+ / 3+ / 4+
Cu: 1+/2+/3+



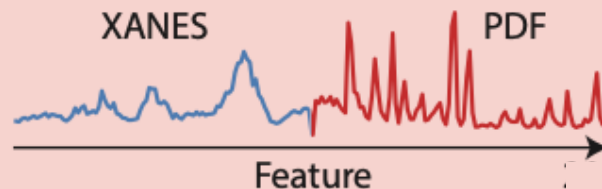
Coordination Number
4/5/6

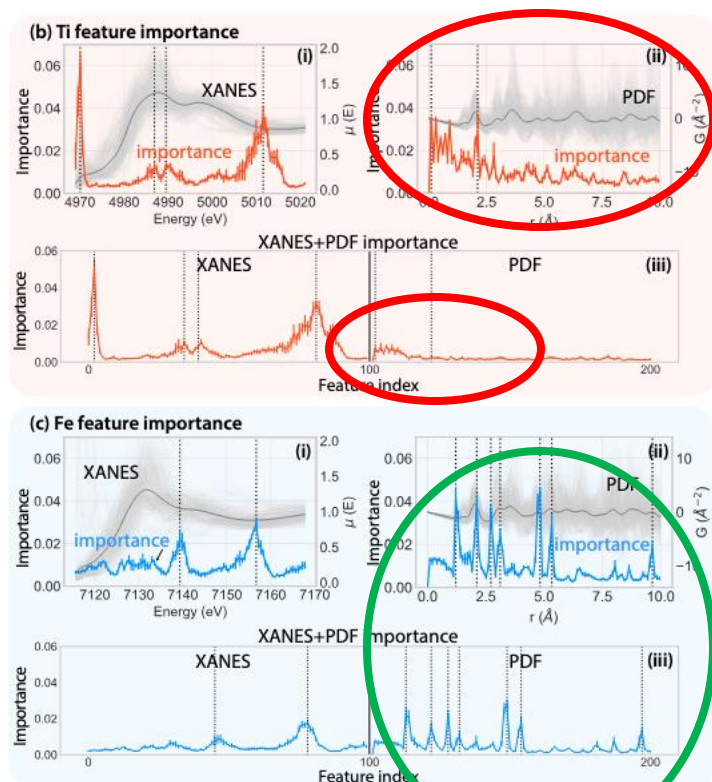
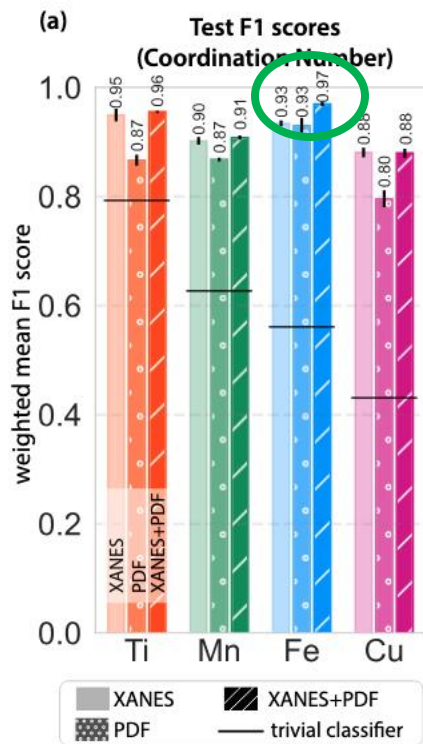


Mean Bond Length
(metal to nearest neighbors)

Random Forest
Models

Feature importance





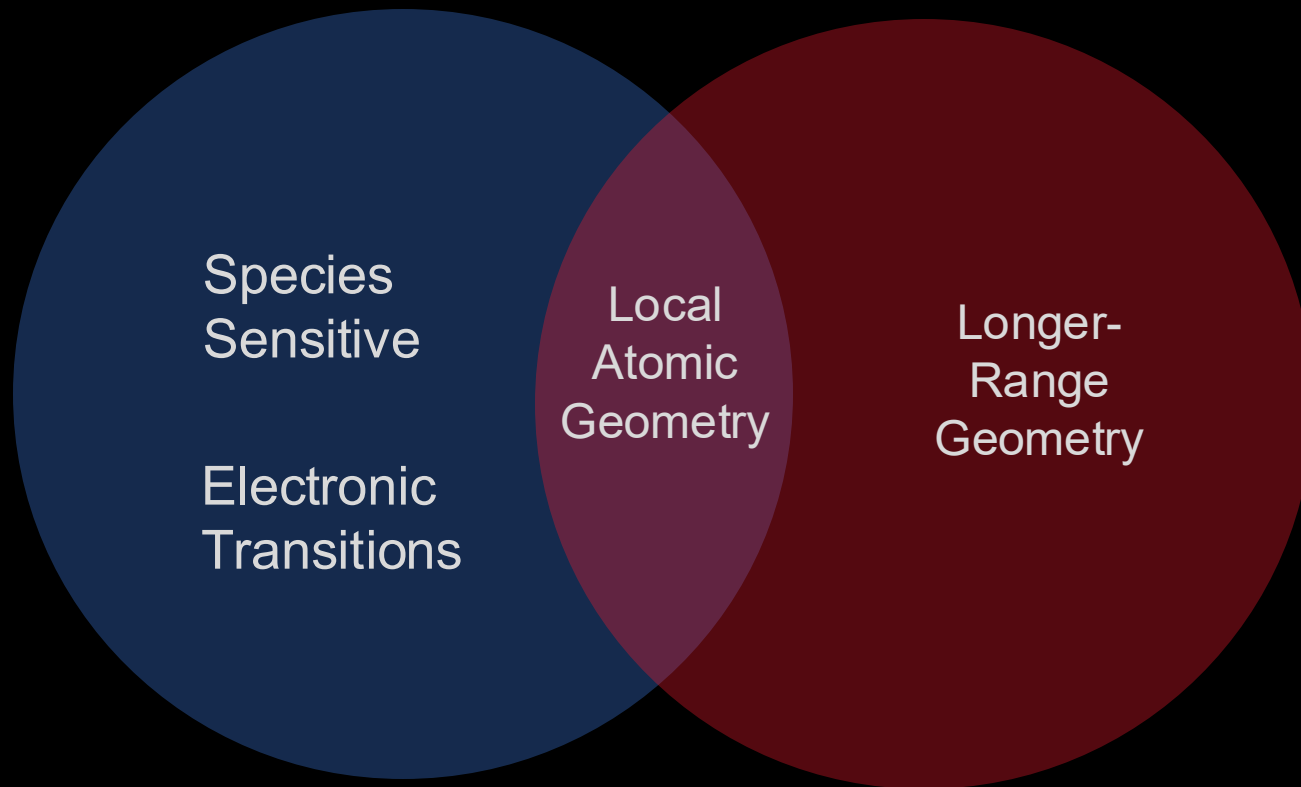
PDF doesn't contribute
When XANES + PDF present
For Titanium

But XANES + PDF improves
performance for Iron

Contributions are system
& (data-set) specific, but can be
measured

XANES

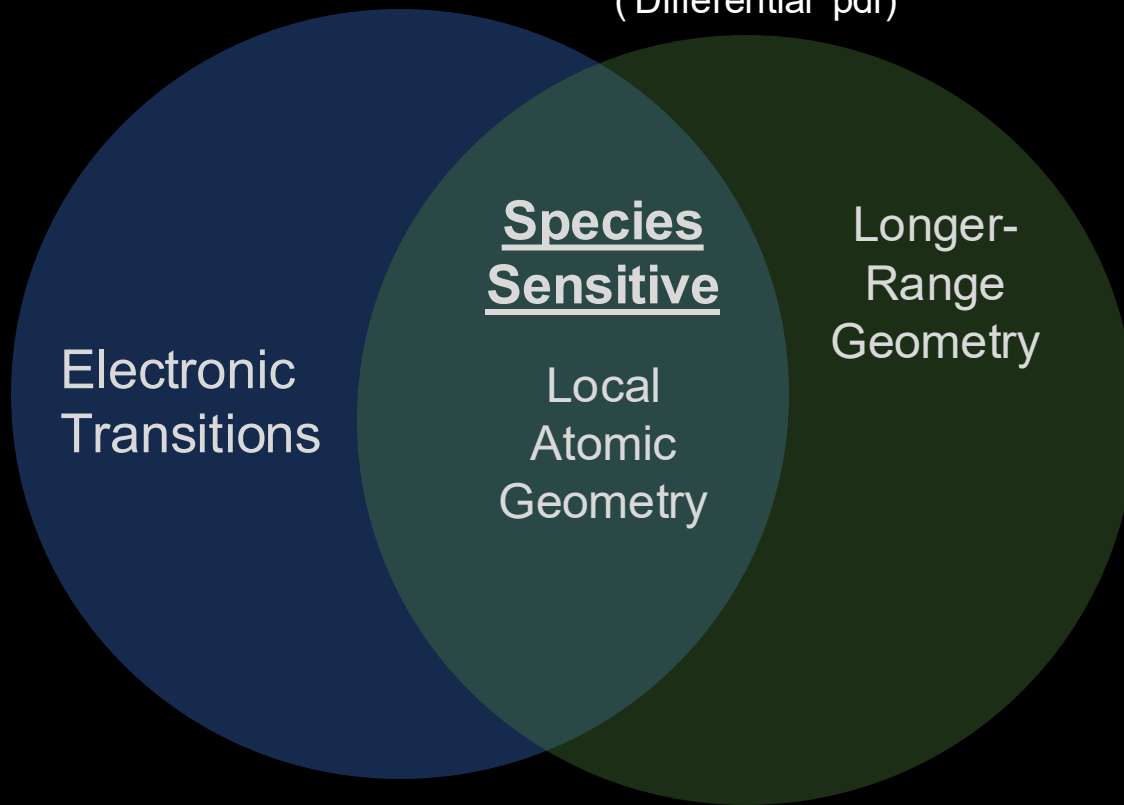
PDF



XANES

dPDF

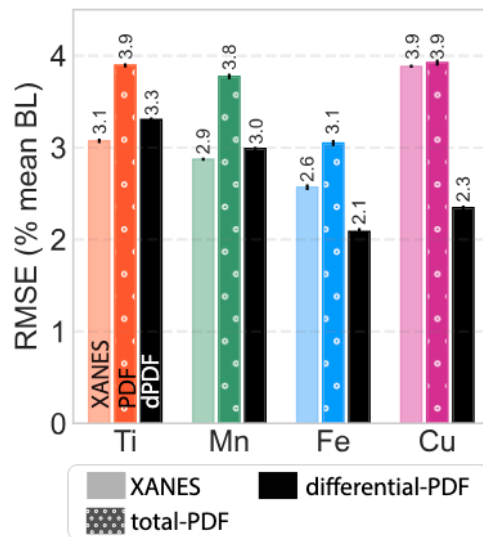
('Differential' pdf)



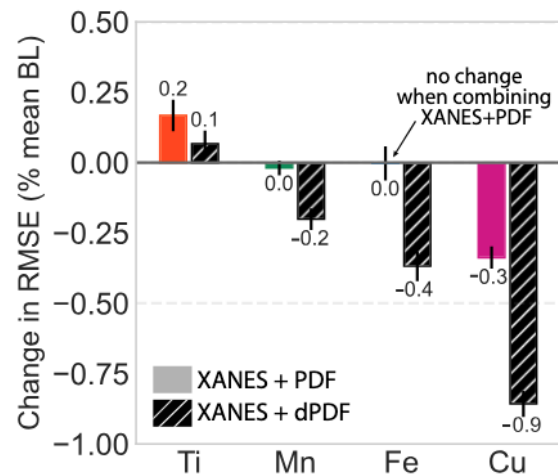
dPDF reduces error significantly;

Closes the gap between XANES and PDF-only models

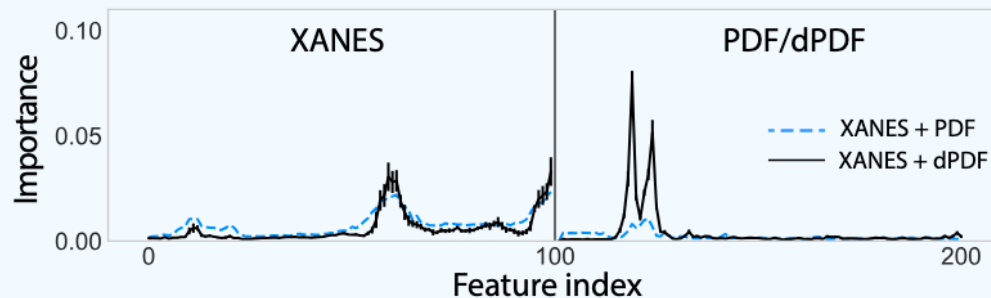
(a) Test RMSEs (Bond Length)
XANES / PDF / dPDF



(b) Multimodal improvement
(difference from XANES-only RMSEs)

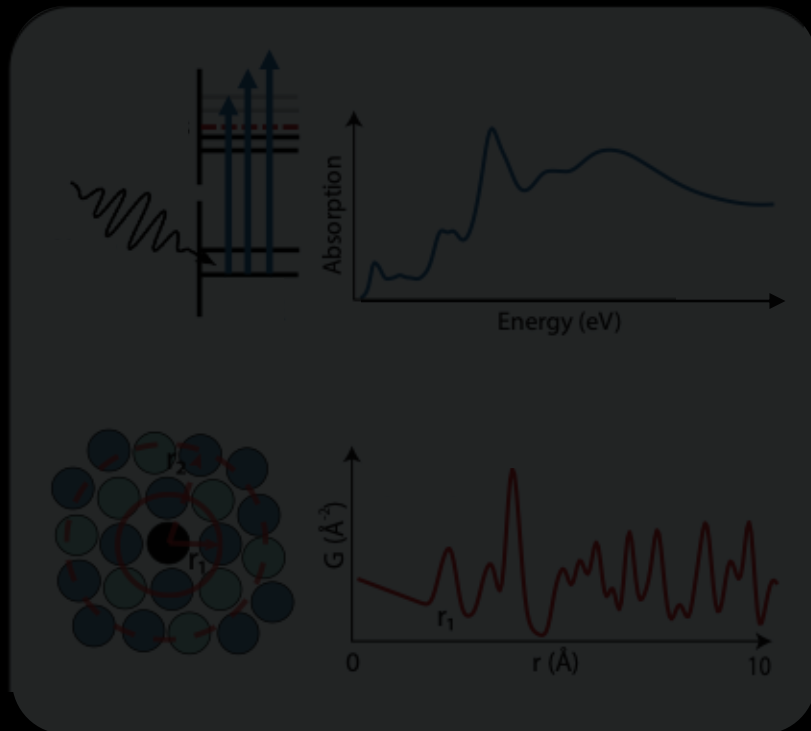


(c) Fe feature importance (multimodal)

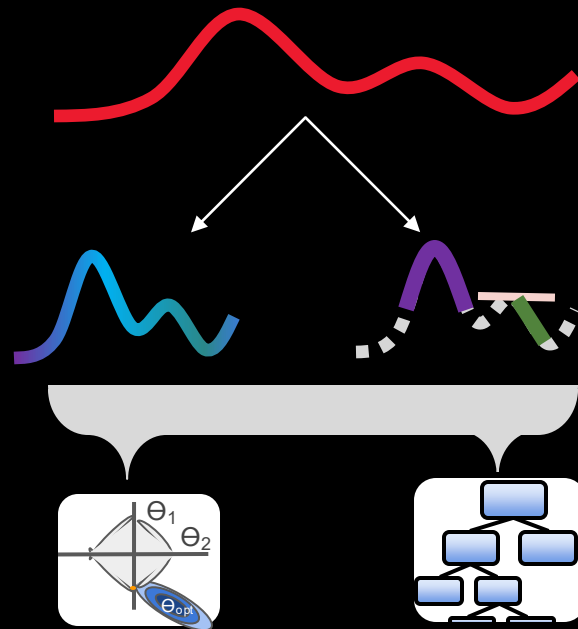


A database approach can help us understand the relative information of different modalities

Multimodal Studies of XANES vs. PDF Spectra



**Spectra-scope: A tool for
feature discovery &
downselection of spectrum-
property relationships**



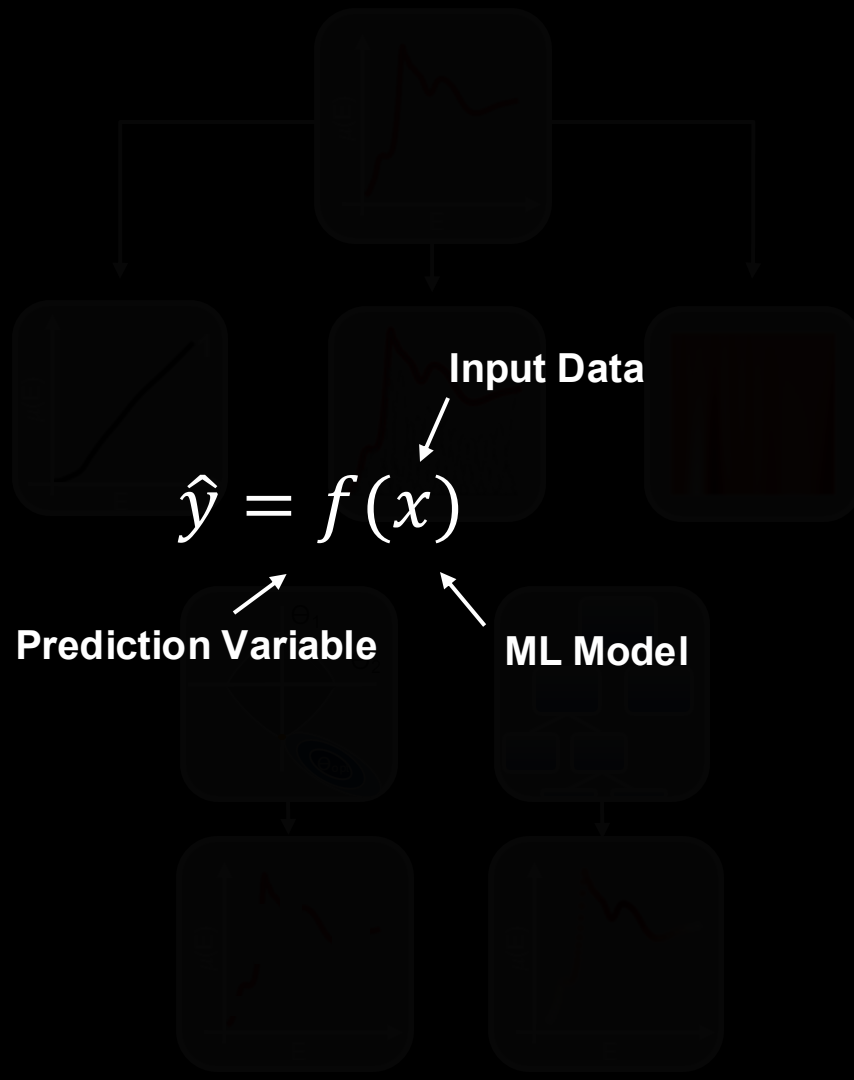
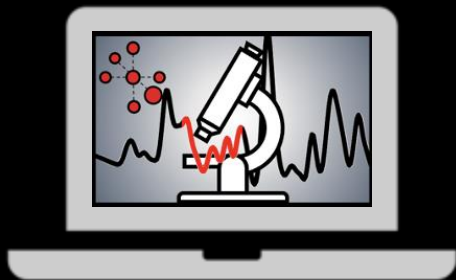
SpectraScope



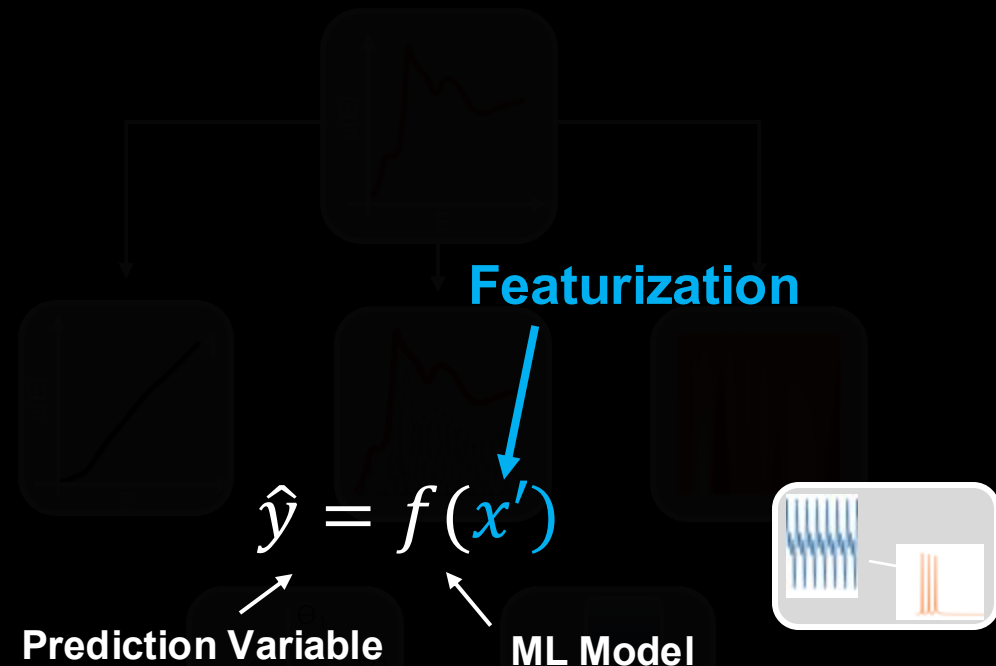
Hypotheses:

1. 'Raw' spectra may be suboptimal representations, especially when using linear models.
2. Simpler models are easier for humans to interpret.
3. Interpretable models are easier for humans to trust.

SpectraScope



SpectraScope



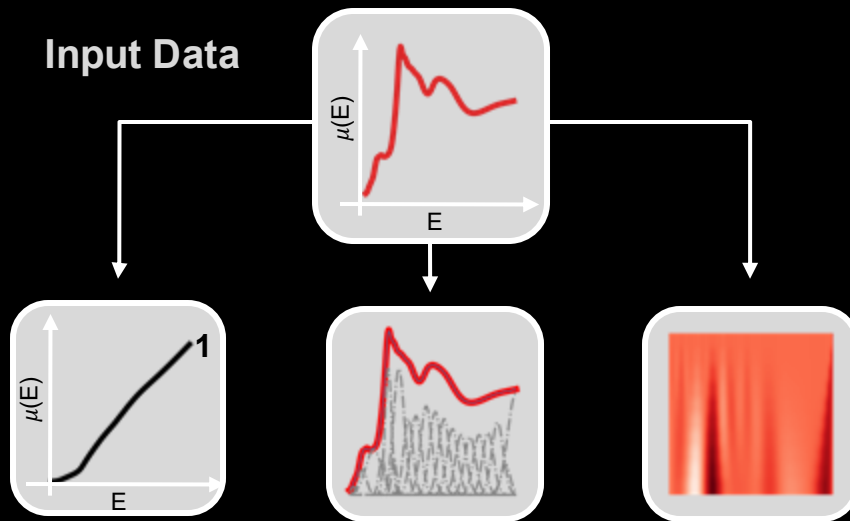
Why featurize?

- Improve performance
- Reduce overfitting

SpectraScope

Featurizers
 x'

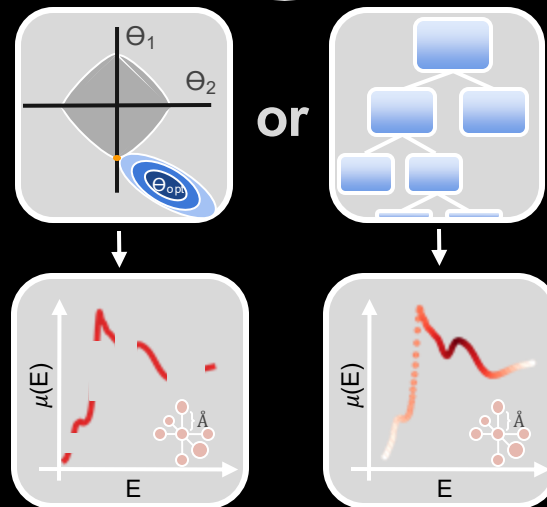
Input Data



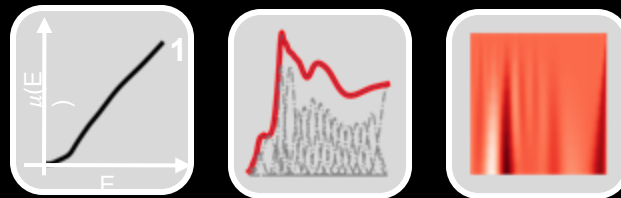
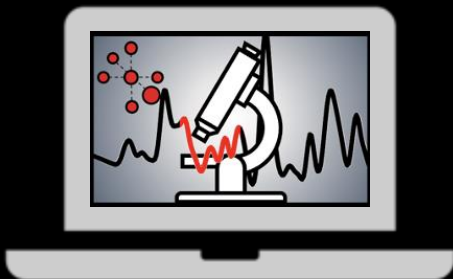
Models

or

Selected Features
& Prediction

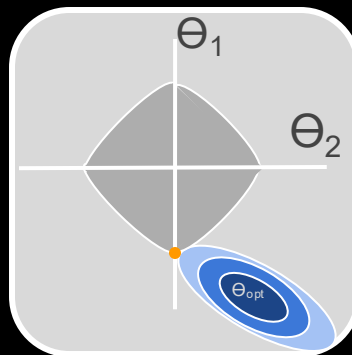


Interpretable & Parsimonious Models



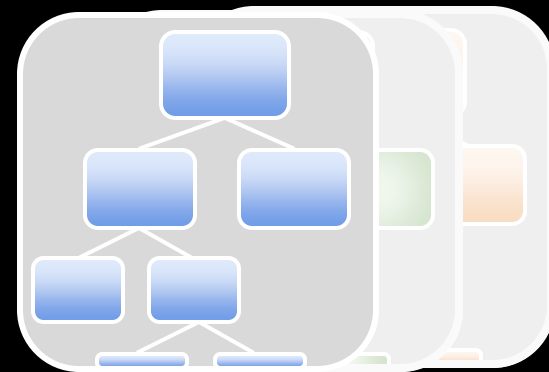
Lasso-Clip-Elastic Net

Models



or

Random Forests



1. **Lasso Regression**
2. **Clip: coefficients $< \text{cutoff} = 0$**
3. **Elastic Net Regression**
4. **Clip**

Seber, P. & Braatz, R. D. <https://arxiv.org/abs/2402.17120>
(2024).

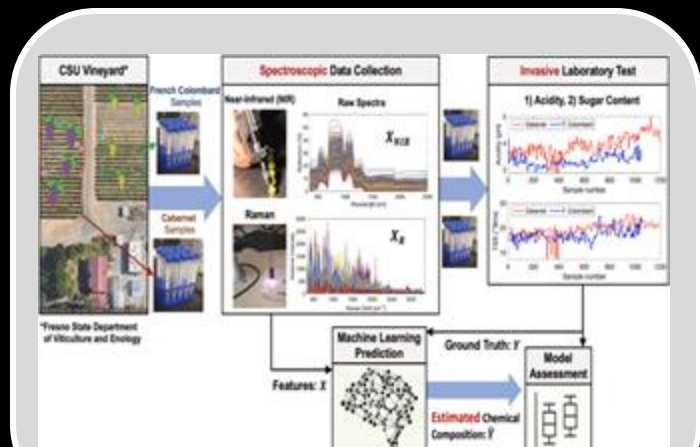
Predicting sugar content in wine grapes

Dataset

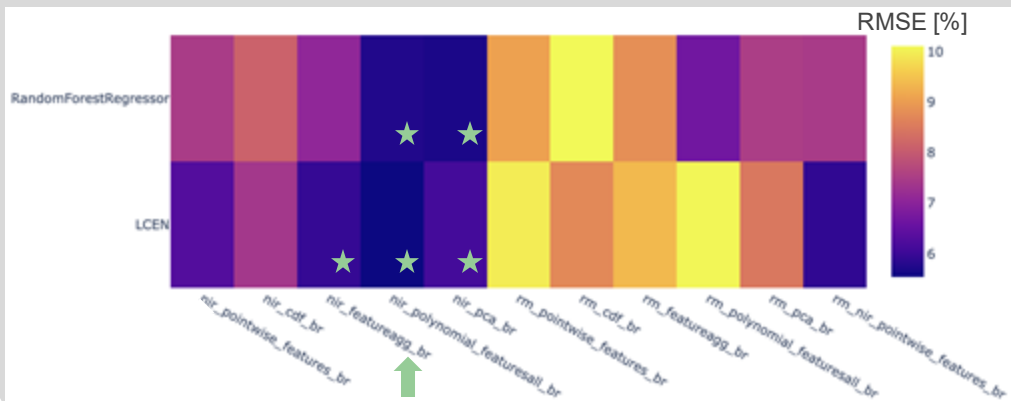
- X: Vis-NIR + Raman
- y: pH, TSS (°Brix)

SpectraScope

Comparing transformations & models



Predicting TSS with Vis-NIR Spectra:
RMSE 5~7% Reported



Well performing transformations

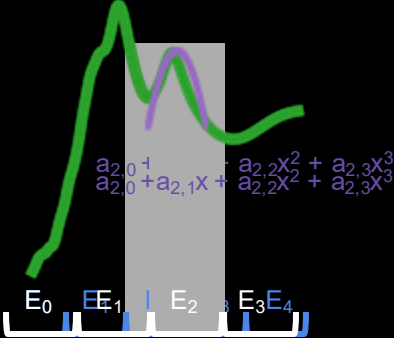
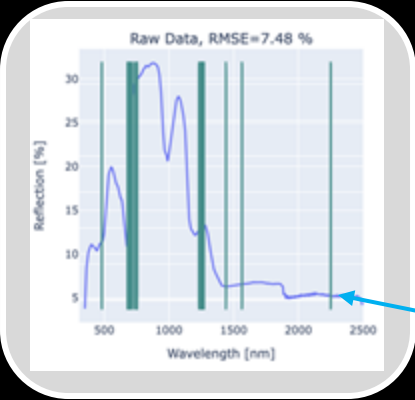
Comparable errors:
RMSE 5~7%

Predicting sugar content in wine grapes

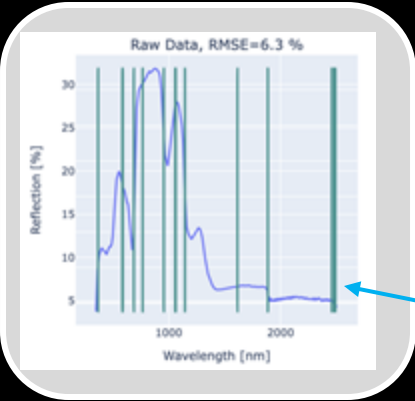
Raw Data

Polynomial Coefficients

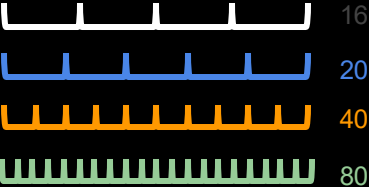
Random Forest



LCEN



$$\vec{a}_4 = [a_{0,0}, a_{0,1}, a_{0,2}, a_{0,3}, a_{0,4}, \dots, a_{20,0}, \dots, a_{33,0}]$$



New feature vector: all fit coefficients (156)

Predicting sugar content in wine grapes

Raw Data

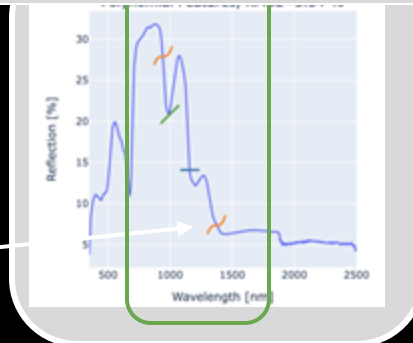
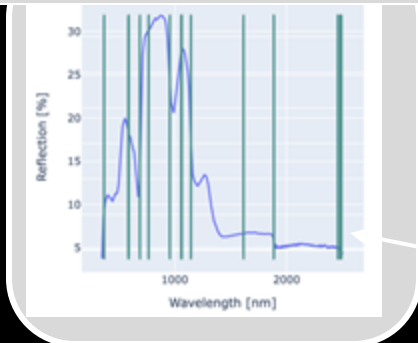
Polynomial Coefficients

Random Forest



1. Transformations can yield better prediction than the raw data
2. Feature selection can highlight individual important parts of spectra

LCEN



Top features

We are turning Spectra-Scope into an app


home_page

visualizing

Spectrascope Home

Deploy

SpectraScope



1. Load data

Done! (using st.cache_data)

☐ Show Full Dataframe

| | mp_id | structure | xanes |
|---|------------|--|---------|
| 0 | mp-1041683 | <code>["@class": "Structure", "@module": "pymatgen.core.structure", "charge": 0, "lattice": {"a":</code> | 0.16897 |
| 1 | mp-556341 | <code>["@class": "Structure", "@module": "pymatgen.core.structure", "charge": 0, "lattice": {"a":</code> | 0.12421 |
| 3 | mp-769525 | <code>["@class": "Structure", "@module": "pymatgen.core.structure", "charge": 0, "lattice": {"a":</code> | 0.09676 |
| 4 | mp-677246 | <code>["@class": "Structure", "@module": "pymatgen.core.structure", "charge": 0, "lattice": {"a":</code> | 0.33554 |
| 5 | mp-5020 | <code>["@class": "Structure", "@module": "pymatgen.core.structure", "charge": 0, "lattice": {"a":</code> | 0.03973 |

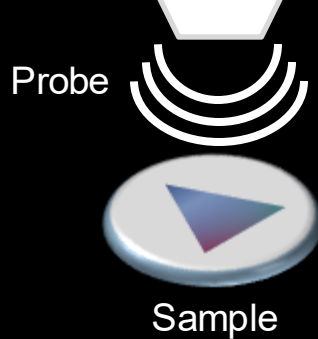
Discussion & Caveats

- These methods rely on data being available, or easily generated
 - Accelerated forward-modeling of complex spectra e.g. XANES may make it easier to assemble libraries of data
 - Same for other more demanding spectra
- Intended to be a guide for practitioners- not to replace them
 - Accelerate signal extraction from data
 - Highlight to practitioners relevant spectral components
- Possible applications may exist within manufacturing, process control

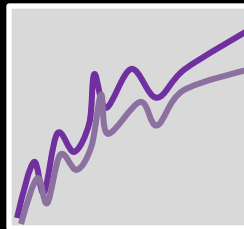
Other Related TRI Work



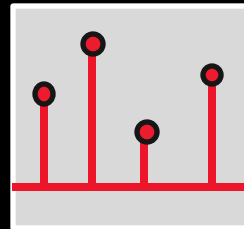
Experimental Data



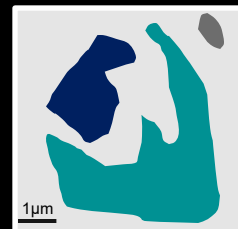
Light Absorption



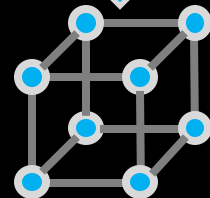
X-Ray Diffraction



Microscopy



Example Modalities



Atomic Structure

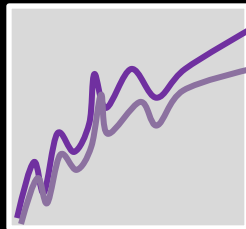
Congruent Idealizations



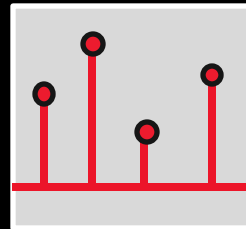
Experimental Data



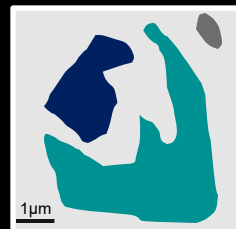
Light Absorption



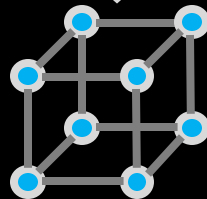
X-Ray Diffraction



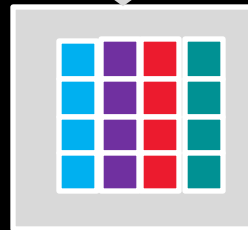
Microscopy



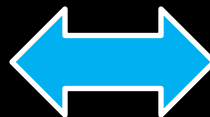
Example
Modalities



Atomic
Structure



Joint Embedding



Data
+
Structure

Congruent
Idealizations

\approx

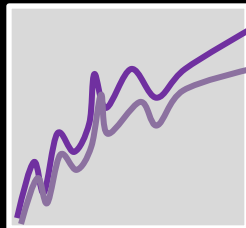
Aligned
Embeddings

Materials Multimodality

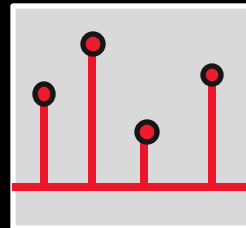
(4M)



Light Absorption



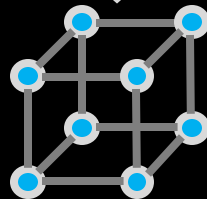
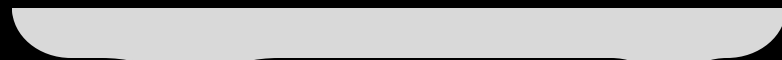
X-Ray Diffraction



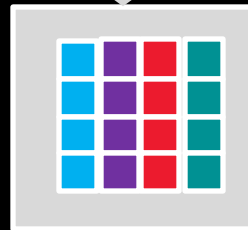
Microscopy



Example
Modalities



Atomic
Structure



Joint Embedding



Data
+
Structure

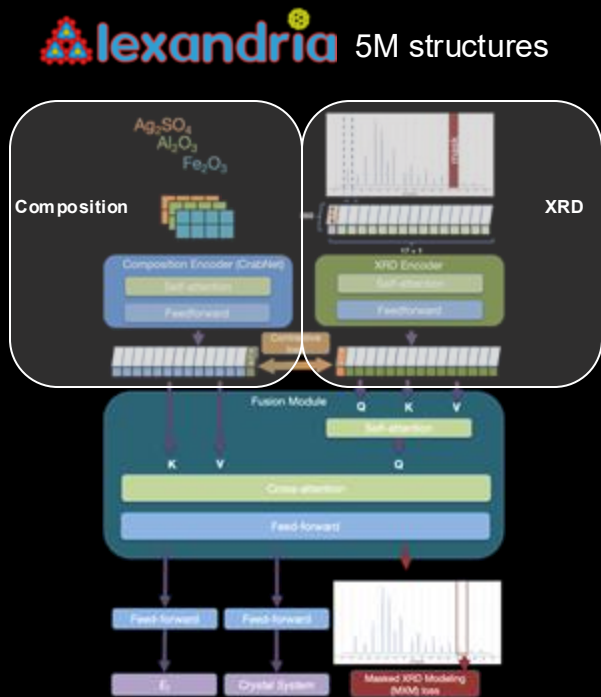
Seeing without crystal structure: structure agnostic multimodal learning for materials science

A strong bias in AI for materials: *atomic structure as the core input representation of materials*

- Cross-attention based multi-modal model based on composition + XRD
- Unsupervised pre-training: up to 4.2× speedup for task convergence
 - Contrastive loss
 - Masked XRD Modeling loss

Bimodal (XRD + Comp.) ~ Unimodal Structure > Unimodal (XRD/Comp.)

- Larger datasets favor multimodal models



XRD x-attention Composition
Transformer-NN (XxaCT-NN)

| | | E_f (MAE, meV/atom) | Crystal System (e.g. cubic, tetragonal, etc. acc. %) |
|----------|-----------------|----------------------------------|--|
| Unimodal | Composition | 132 | 67.8 |
| | XRD | 421 | 92.3 |
| | Structure | <u>16</u> [Schmidt et al., 2024] | - |
| Bimodal | Composition-XRD | <u>27.1</u> | <u>97.2</u> |

Read the preprint!
Subramanian, Hung,
Schweigert, Suram, Ye



Conclusions on Spectroscopy

i

A database approach can help us understand contributions of different spectral modes

ii

Feature generation and downselection can help reveal spectrum-property trends



iii

These might be combined for future optimal experiment planning & design

Talk Outline

1

TRI + AMDD Background

2

Challenges in AI-guided materials design

3

Characterization & bridging experiment & theory

Conclusions

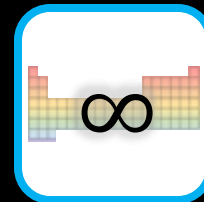
I

Conventional metrics for materials discovery may be limited in their applicability



II

We should focus on ways to improve synthesis success rates



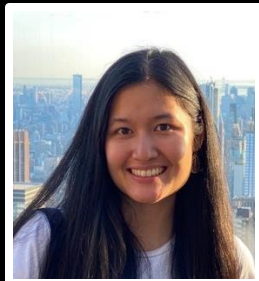
III

Interpretable ML may help practitioners make more trustworthy models





Zoe Zachko



Dr. Tina Na Narong



Amalya Johnson
(Stanford)



Weiike Ye



Leena
Sansguiri



Joey
Montoya

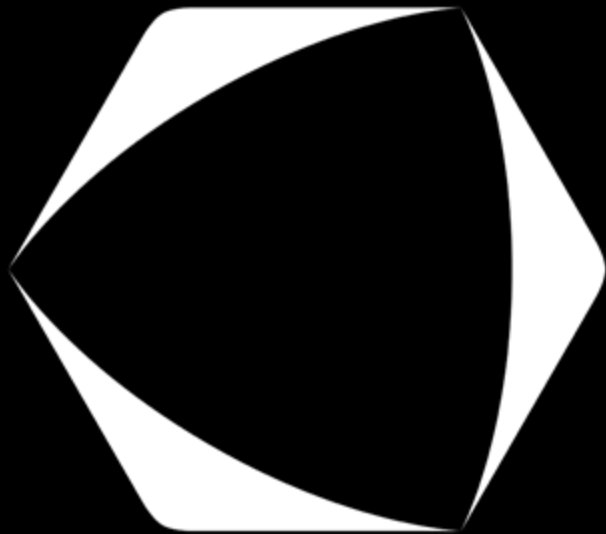


Prof. Simon Billinge

Thank you very
much to **the**
organizers, and to
you for your
attention!

Sabina Mohan,
NIST Conference Services

Daniel Wines
Kamal Choudhary
Francesca Tavazza
Brian DeCost
... and others!



Thank You

Questions?

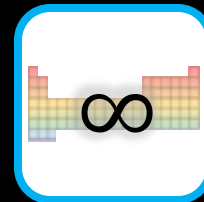
I

Conventional metrics for materials discovery are limited in their applicability



II

We should focus on ways to improve synthesis success rates



III

Interpretable ML can help practitioners make more trustworthy models



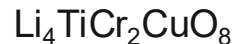
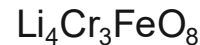
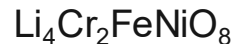
steven.torrisi@tri.global

What success rates might we see?

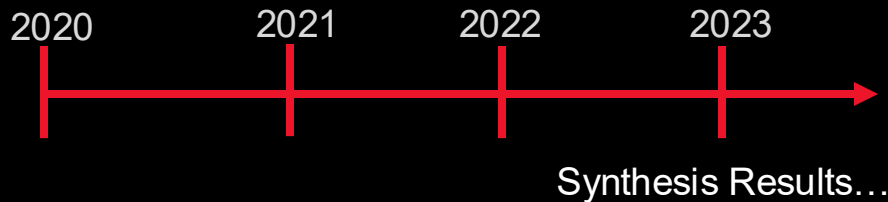
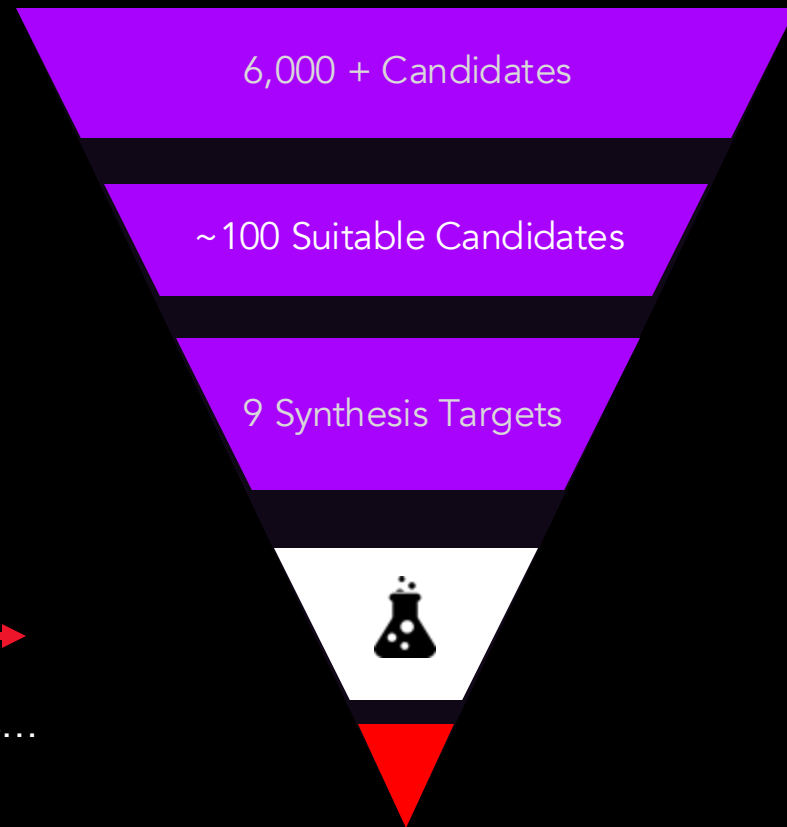
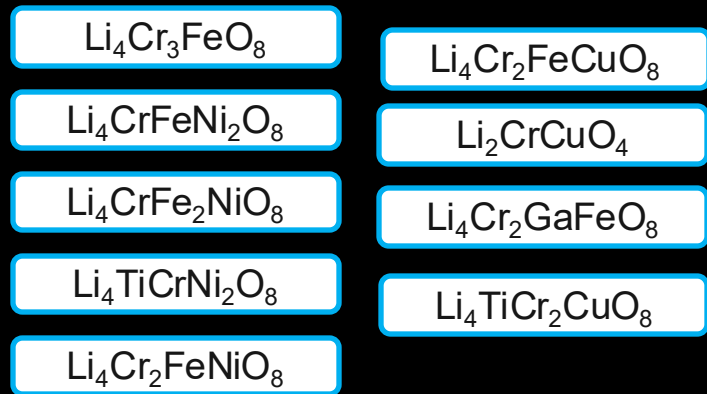
Attempted Syntheses based on Energy Metric

| | |
|---------|-------------|
| Ca-Ag-O | $CaAg_3O_4$ |
| | Ca_3AgO_4 |
| Li-Ag-O | $LiAg_3O_4$ |
| | $LiAgO_2$ |
| Ca-Pd-O | $CaPd_3O_6$ |
| | $CaPdO_3$ |
| Mg-Ru-O | $MgRu_2O_6$ |
| | $MgRuO_4$ |
| Ca-Ru-O | $CaRuO_4$ |
| | Ca_3RuO_6 |
| Li-Ru-O | Li_4RuO_5 |
| | Li_5RuO_5 |

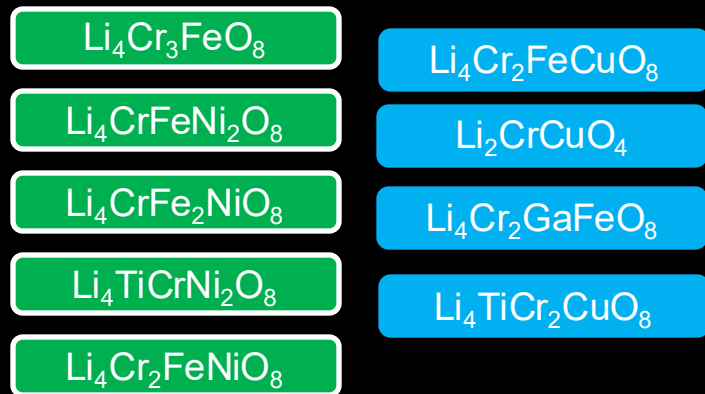
9 Target Candidates



Computer-to-Lab Pipeline



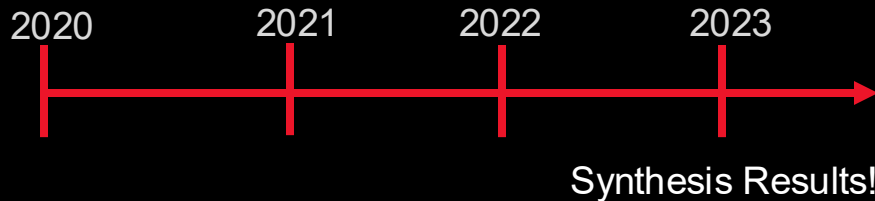
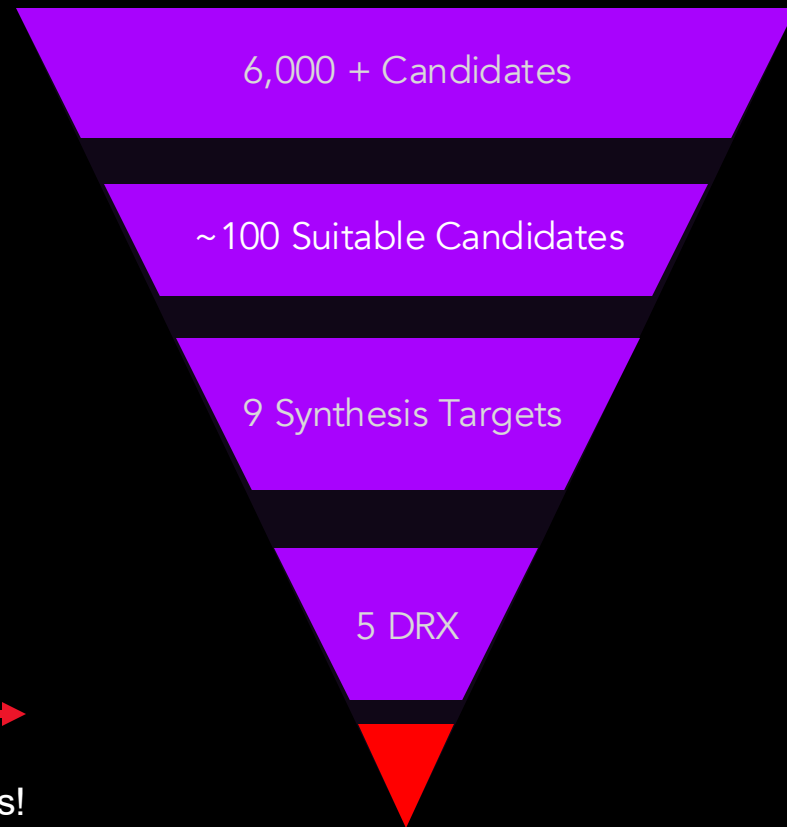
Computer-to-Lab Pipeline



DRX

Layered

5/9 in DRX!



Further problems for the field for discussion

Do we suffer from a
'computability bias' towards
ordered unit cells?



Streetlight effect

[Article](#) [Talk](#)

From Wikipedia, the free encyclopedia

Not to be confused with [Street light interference phenomenon](#).

The **streetlight effect**, or the **drunkard's search** principle, is a type of **observational bias** that occurs when people only search for something where it is easiest to look.^[1] Both names refer to a well-known joke:

Further problems for the field for discussion

Do many theorists suffer from a 'computability bias' towards ordered unit cells?



MP, OQMD, OCP etc are all wonderful resources and have come to shape a generation of scientific work-

But do they implicitly limit theorists' picture of what constitutes a material?

Speaking as someone with a DFT background: is it an issue that we focus so heavily on materials that 'fit through the keyhole' of DFT?

Further problems for the field for discussion



Things sometimes left out when working at the 10,000 ft level (and the combinatorics of these will be challenging to contend with!):

- Surfaces
- Interfaces
- Defects
- Finite temperature effects, dynamic disorder
- Long length & time scale events (e.g. reconstruction)
- Interaction with liquids, atmospheres
- Functionalization processes
- Micro-scale phenomena (e.g. grain boundaries)