

REVISITING MATERIAL STRUCTURE IN THE TIME OF AI

Simon J. L. Billinge

¹Columbia University

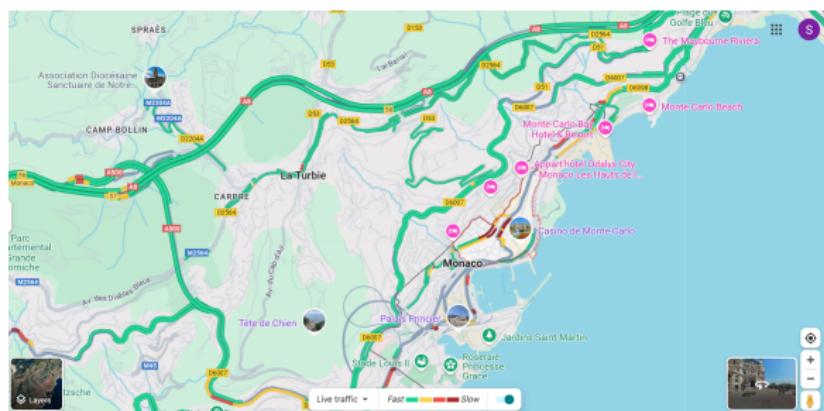
July 9, 2025

WHAT I LEARNED ABOUT MATERIALS FROM AI

What I learned about materials from AI

WHAT I LEARNED ABOUT MATERIALS FROM AI

Let me recap a bit of my journey

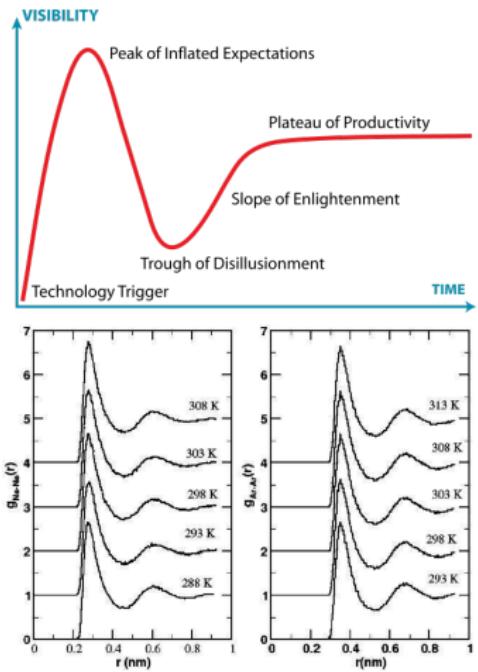


WHAT I LEARNED ABOUT MATERIALS FROM AI

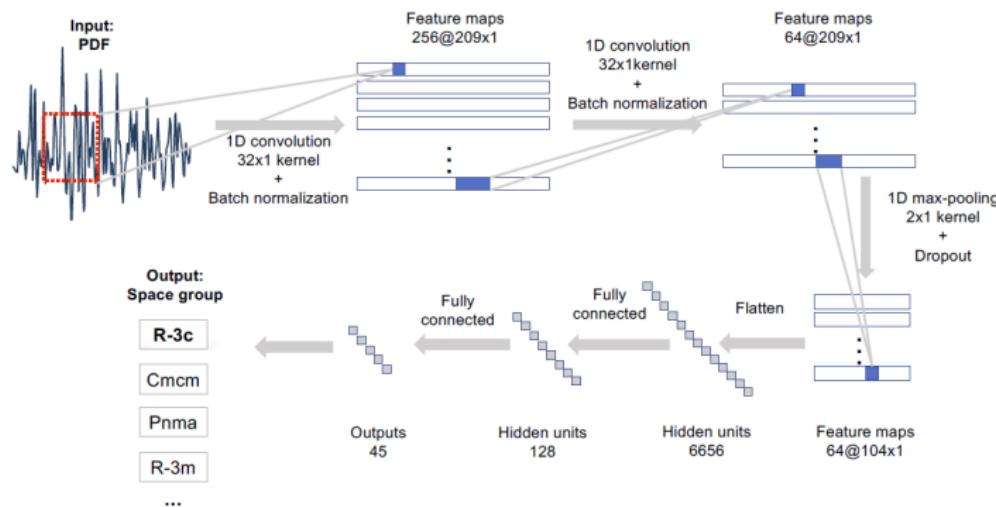
It started with a classification problem

- We can't get space-group from PDF
 - But PDF is just the FT of diffraction data
 - We can get SG from diffraction data

So the information is there, can a ML model figure out a surrogate forward problem?



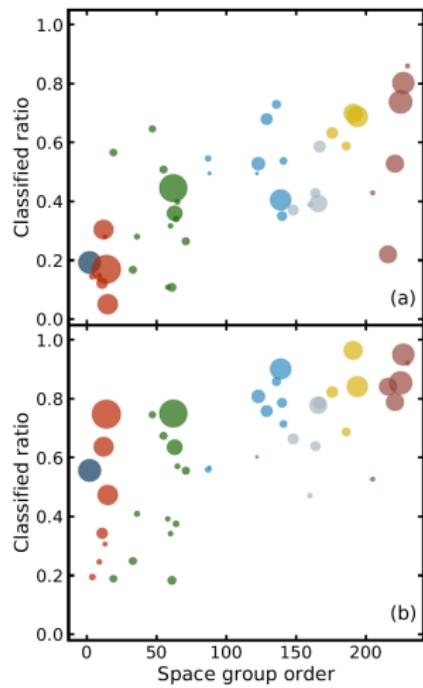
GIVEN A PDF, WHAT WAS THE SPACE-GROUP OF THE MATERIAL?



collaboration with group of Qiang Du (Columbia U)

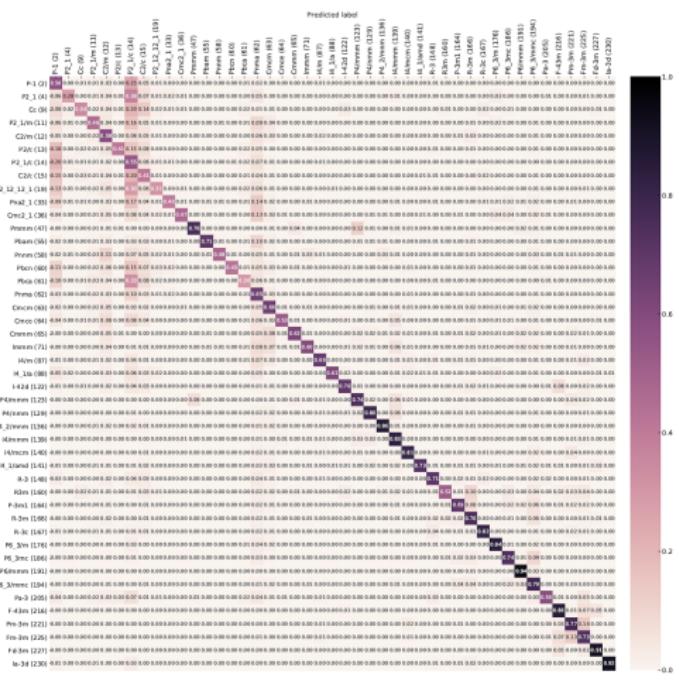
Chia-Hao Liu, SJLB et al., Acta Cryst. (2019), 10.1107/S2053273319005606

HOW WELL DID IT PERFORM?



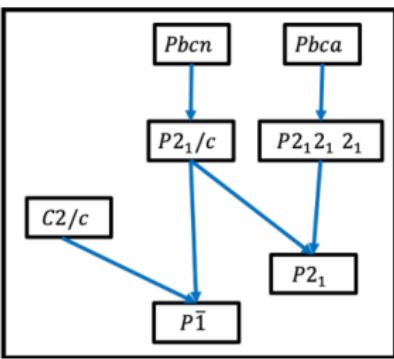
- Top six classification ratio
- (top) Logistic Regression baseline
- (bottom) CNN 91% overall
- Color indicates crystal system
- Size of spot indicates sample-size

WE LEARNED CRYSTALLOGRAPHY FROM THE MODEL



Matrix of confusion

- tear drops on space-groups
 $C2/1$, $Pnma$
 - connected by
group-sub-group
relationships



CAN WE DO STRUCTURE SOLUTION WITH GENERATIVE AI?

The goal

- Input: the PDF or a powder diffraction pattern
- Output: the structure of the material

That would be cool!

WHAT I LEARNED ABOUT AI BY APPLYING IT TO STRUCTURE SOLUTION

What I thought I knew

- Gen-AI is doing interpolation

WHAT I LEARNED ABOUT AI BY APPLYING IT TO STRUCTURE SOLUTION

What I thought I knew

- Gen-AI is doing interpolation
- But how do we interpolate structures into each other?

WHAT I LEARNED ABOUT AI BY APPLYING IT TO STRUCTURE SOLUTION

What I thought I knew

- Gen-AI is doing interpolation
- But how do we interpolate structures into each other?
- For close-packed structures, I can interpolate from fcc to bcc by introducing stacking-fault defects.

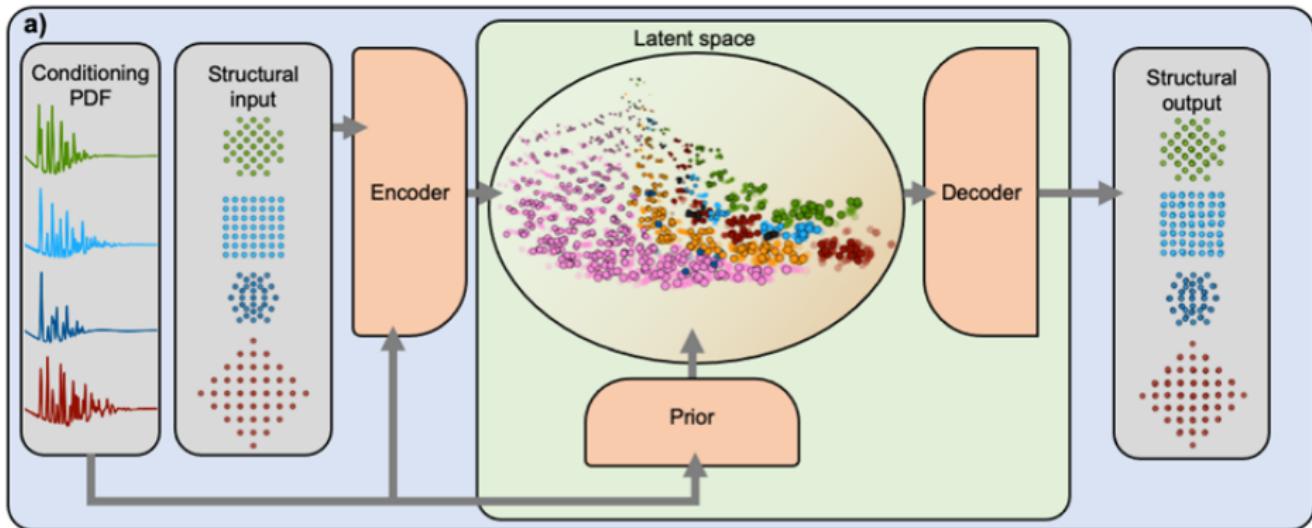
Let's do close-packed structures that interpolate

CASE 1: CLOSE PACKED STRUCTURES

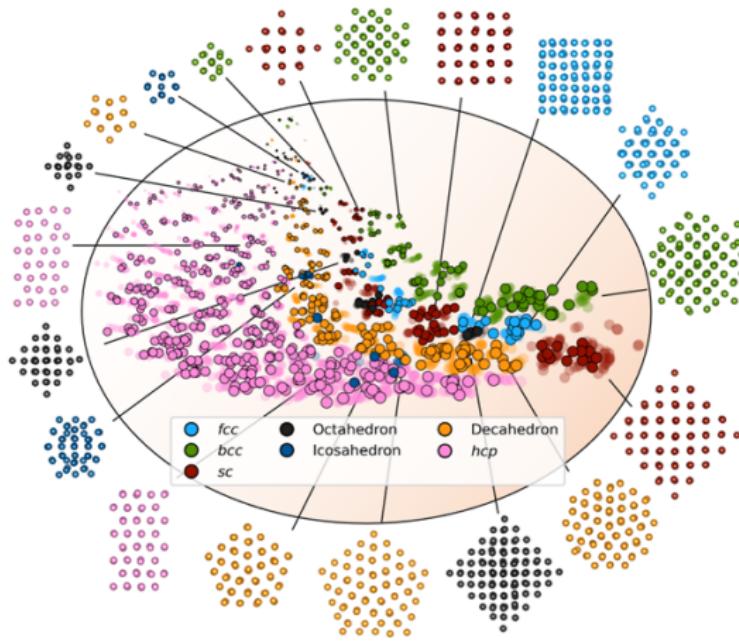
- Graph Convolutional Variational Autoencoder (CVAE)
- Collaboration with Kirsten Jensen, Raghav Selvan, U. Copenhagen
- Work of students Emil Kjaer and Andy Anker
- Kjaer, Anker, et al., Digital Discovery (2023),
[10.1039/D2DD00086E](https://doi.org/10.1039/D2DD00086E)

Structure solution using Neural Nets

DEEPSTRUCL: CLOSE PACKED NANOPARTICLES WITH A CVAE



DEEPSTRUCL: CLOSE PACKED NANOPARTICLES WITH A CVAE



WHAT I LEARNED ABOUT AI BY APPLYING IT TO STRUCTURE SOLUTION

What I thought I knew

- Gen-AI is doing interpolation
- But how do we interpolate structures into each other?
- For close-packed structures, I can interpolate from fcc to bcc by introducing stacking-fault defects.
- It worked!

WHAT I LEARNED ABOUT AI BY APPLYING IT TO STRUCTURE SOLUTION

What I thought I knew

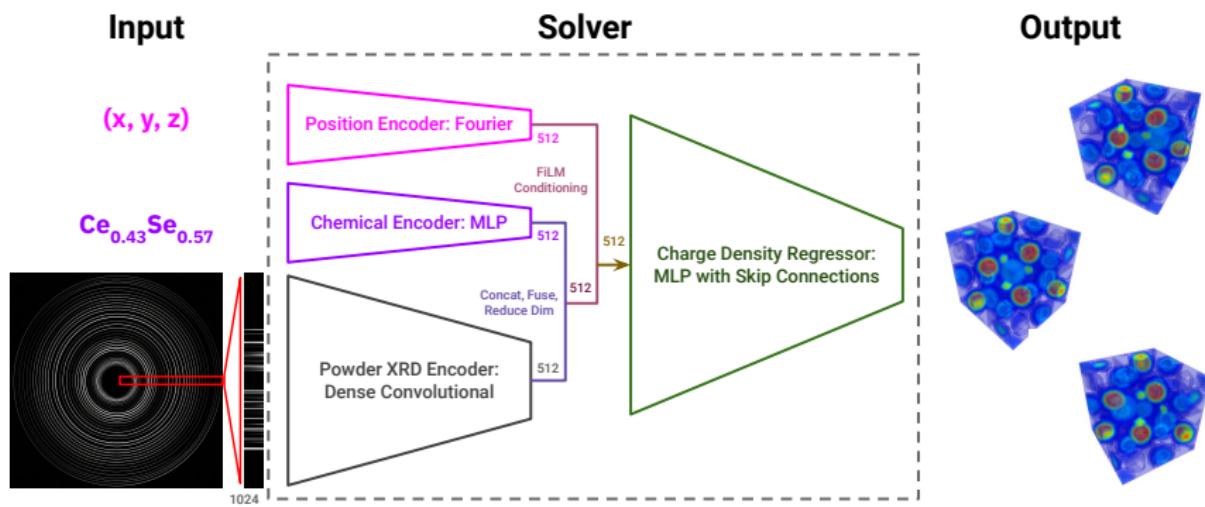
- Gen-AI is doing interpolation
- But how do we interpolate structures into each other?
- For close-packed structures, I can interpolate from fcc to bcc by introducing stacking-fault defects.
- It worked!
- But in general? If I draw a line from NaCl to nickel, I don't find CaCO_4 in between

Let's change the problem to something continuous....electron-density!

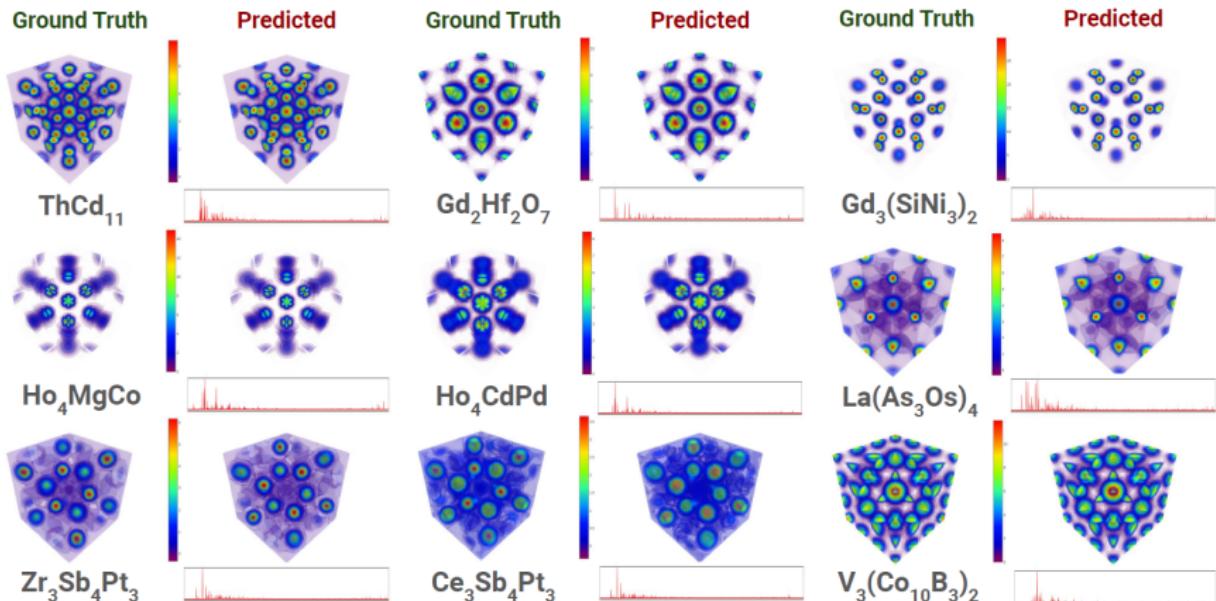
CASE 2: ELECTRON DENSITY DETERMINATION

- Collaboration with Hod Lipson (Columbia University)
- Work of students Gabriel Guo and Ling Lan
- G. Guo et al., npj-Computational Materials, to be published
- Graph Convolutional Variational Autoencoder (CVAE)
- input: powder diffraction pattern and composition
- output: electron density distribution
- Trained on cubic and trigonal (non-orthogonal but high symmetry) crystal systems

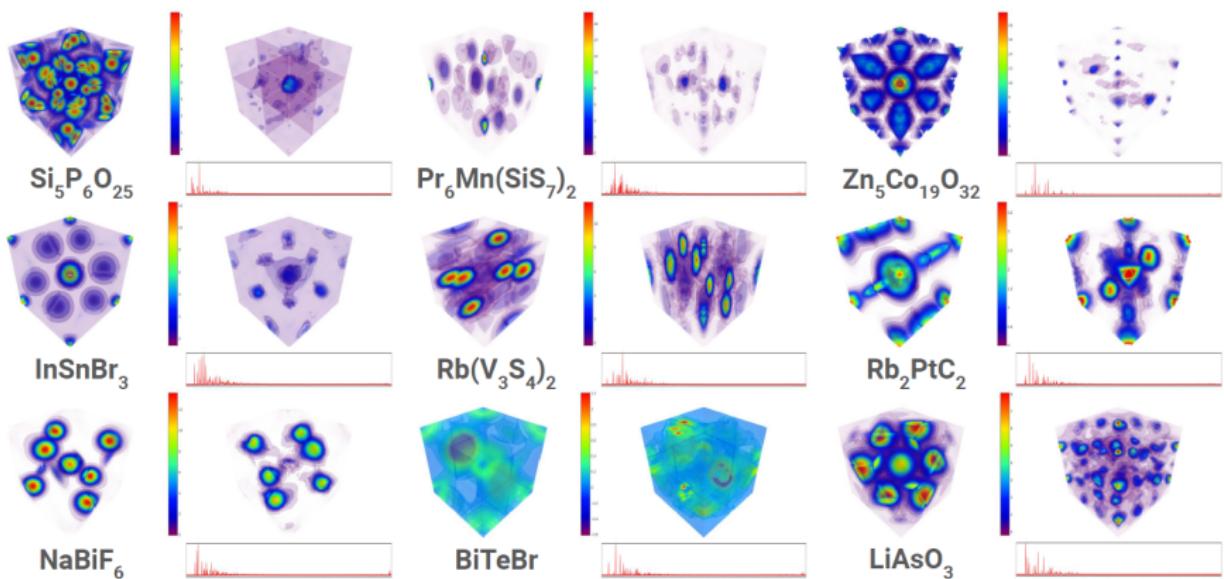
CRYSTALNET: ELECTRON DENSITY RECONSTRUCTION WITH CVAE



CRYSTALNET: ELECTRON DENSITY RECONSTRUCTION WITH CVAE



CRYSTALNET: ELECTRON DENSITY RECONSTRUCTION WITH CVAE



WHAT I LEARNED ABOUT AI BY APPLYING IT TO STRUCTURE SOLUTION

What I learned

- Gen-AI is doing interpolation
- But how do we interpolate structures into each other?
- In general? If I draw a line from NaCl to nickel, I don't find CaCO₄ in between

WHAT I LEARNED ABOUT AI BY APPLYING IT TO STRUCTURE SOLUTION

What I learned

- Gen-AI is doing interpolation
- But how do we interpolate structures into each other?
- In general? If I draw a line from NaCl to nickel, I don't find CaCO₄ in between
- But the interpolation is happening in the latent space of the auto-encoder, not in our familiar "structure-space"

WHAT I LEARNED ABOUT AI BY APPLYING IT TO STRUCTURE SOLUTION

What I learned

- Gen-AI is doing interpolation
- But how do we interpolate structures into each other?
- In general? If I draw a line from NaCl to nickel, I don't find CaCO₄ in between
- But the interpolation is happening in the latent space of the auto-encoder, not in our familiar "structure-space"

The magic of diffusion models

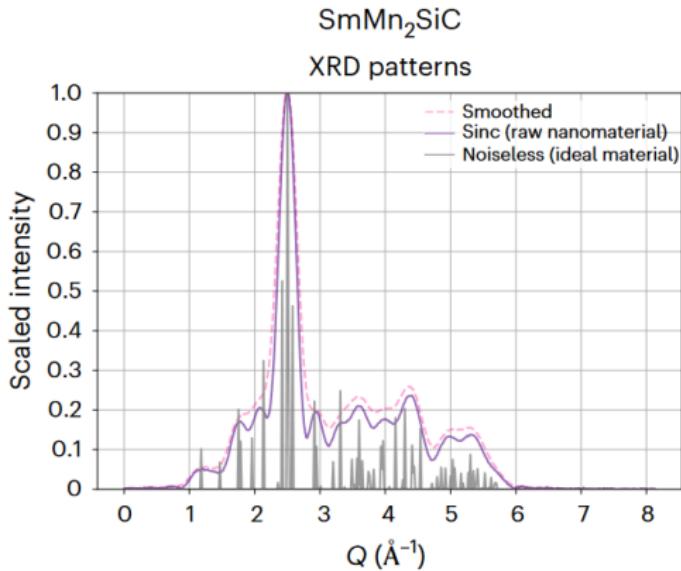
CASE 3: STRUCTURE SOLUTION OF NANOPARTICLES

- Collaboration with Hod Lipson (Columbia University)
- Work of students Gabriel Guo and Tristan Saidi
- G. Guo et al., Nature Materials, (2025)
<https://doi.org/10.1038/s41563-025-02220-y>
- input: nanostructure powder diffraction pattern and composition
- output: "structure": unit cell, coordinates and coloring

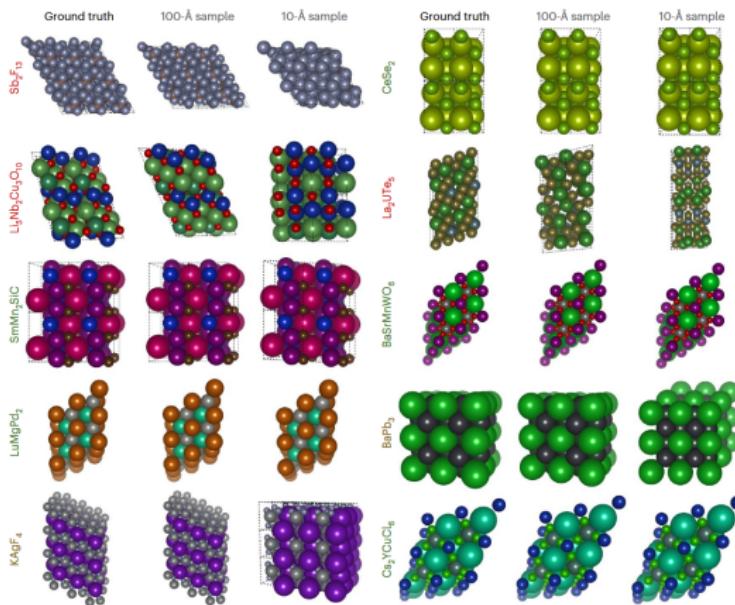
PXRDNet: STRUCTURE SOLUTION WITH DIFFUSION MODEL

PXRDNet

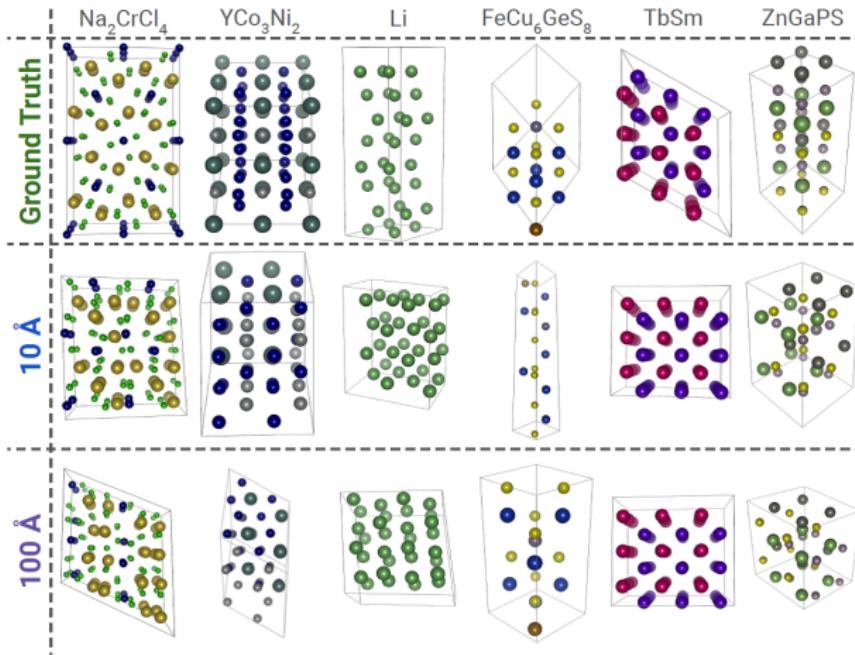
- 100 Å or 10 Å NP powder data.
Low information content!
- Based on CDVAE (T. Xie, Proc. International Conference on Learning Representations (2022)).
- SE(3) equivariant graph NN autoencoder for composition, lattice parameters and number of atoms.
- Denoising diffusion via noise-conditioned score networks for the atomic coordinates



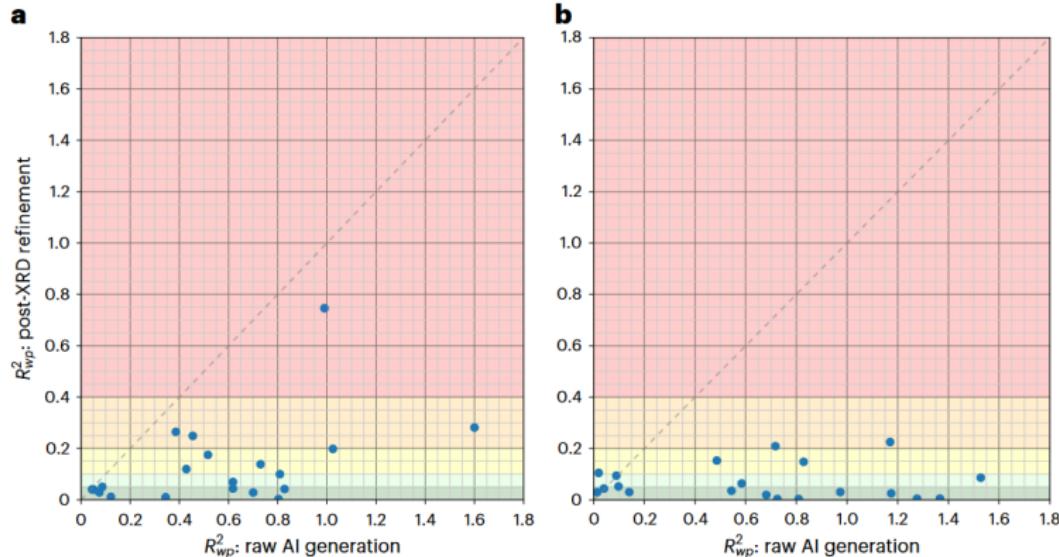
PXRDNET: STRUCTURE SOLUTION WITH DIFFUSION MODEL



PXRDNET: STRUCTURE SOLUTION WITH DIFFUSION MODEL



PXRDNET: STRUCTURE SOLUTION WITH DIFFUSION MODEL



DIFFUSION MODELS ARE THE S\$%T

Diffusion models are the s\$%t

- They build latent spaces/embeddings that can interpolate non-interpolatable things
- They do it by starting with a large population of “good” things
- They then randomize them in a well defined way
- Any point in this space can be denoised
- A point in the space can be picked by a prompt, which could be a desired target

DIFFUSION MODELS ARE THE S\$%T

Diffusion models are the s\$%t

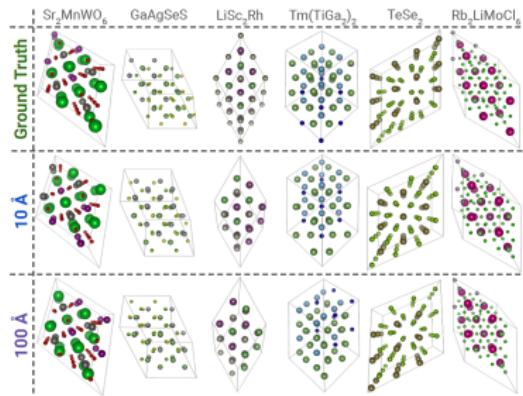
- They build latent spaces/embeddings that can interpolate non-interpolatable things
- They do it by starting with a large population of “good” things
- They then randomize them in a well defined way
- Any point in this space can be denoised
- A point in the space can be picked by a prompt, which could be a desired target

They work!

EXAMPLE: PXRDNET

It works!

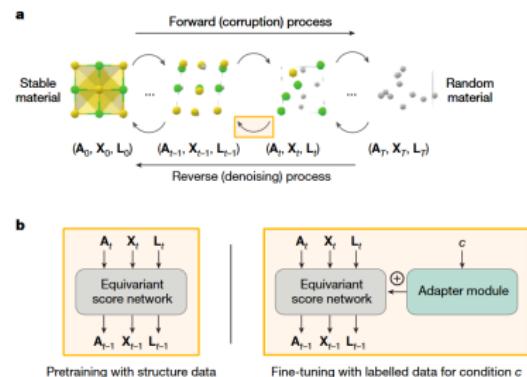
- After training the latent space and conditioning on (rather small amounts of) diffraction data, it makes a rather good guess of the structure solution
- Enforces $SO(3)$ equivariance so only interested in rigid bodies



EXAMPLE: MATTERGEN

Let's throw some serious Microsoft money at the same problem

- Trained on a dataset that is $> 10\times$ larger than PXRDnet using serious Microsoft compute
- In most respects, in common with PXRDnet
- Can be conditioned on many different labels. Has to be trained, but the foundation model means that it can be trained for the new task on a small amount of data
- At the time of writing of the MS could only be conditioned on scalars. I'm sure that has changed.



It works!

WHAT I LEARNED ABOUT AI - SUMMARY

- Diffusion models are the sh\$%t
- Emerging foundation models built on large, diverse, datasets are adaptable for many task with modest retraining

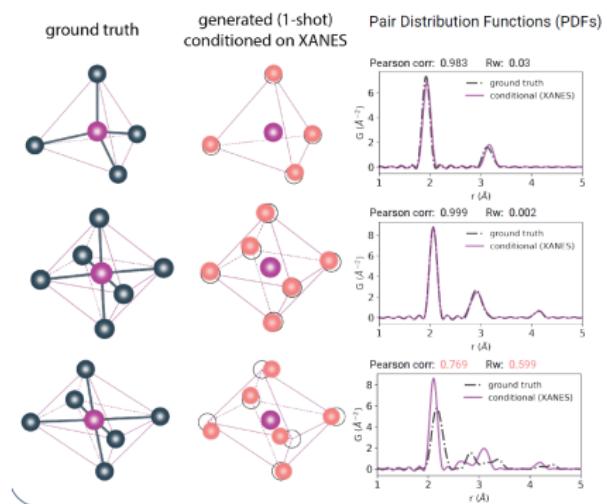
They work

LOCAL STRUCTURE FROM XANES WITH A DIFFUSION MODEL

Solving local structure from XANES (not XAFS!)

- Work of Aniv Ray and Tina Na Narong
- Collaboration with Steven Torrisi and Hod Lipson
- Use a diffusion model to extract local structure (with some surprises)

Please see Tina's poster



WHAT I LEARNED ABOUT MATERIALS BY APPLYING AI TO STRUCTURE SOLUTION

What I learned about materials by applying AI to structure solution

STRUCTURE SOLUTION

Structure of the problem

- Input is data, a diffraction pattern or PDF
- Output is a structure
- To train the model we need to tell it when it got the structure right vs wrong
- We know the structure, but we need a categorical success/failure measure
- How do we do that?

STRUCTURE SOLUTION

How do we assess if a structure solution is right?

- Need a measure of similarity between structures
- Need a threshold when structures are sufficiently similar to determine them as being “the same” and so “the structure is solved”

STRUCTURE SOLUTION

How do we assess if a structure solution is right?

- Need a measure of similarity between structures
- Need a threshold when structures are sufficiently similar to determine them as being “the same” and so “the structure is solved”
- btw, this problem affects Materials Discovery. MatterGen talks about SUN materials, “stable”, “unique”, “new” (i.e., not the same as a known structure).

WHEN ARE STRUCTURES THE SAME PROBLEM (WASTS PROBLEM)

But wait, there are many papers in the literature that do this without much discussion. How do they do it?

- They use a structure similarity metric in pymatgen “structureMatcher”
- But this appears not to have been heavily validated. We are taking on blind faith that is working as hoped
- Some chemists seem to disagree

WHEN ARE STRUCTURES THE SAME PROBLEM (WASTS PROBLEM)

But wait, there are many papers in the literature that do this without much discussion. How do they do it?

- They use a structure similarity metric in pymatgen “structureMatcher”
- But this appears not to have been heavily validated. We are taking on blind faith that is working as hoped
- Some chemists seem to disagree



Open Access

This article is licensed under CC-BY 4.0

Perspective

pubs.acs.org/cm

Artificial Intelligence Driving Materials Discovery? Perspective on the Article: Scaling Deep Learning for Materials Discovery

Anthony K. Cheetham* and Ram Seshadri*

THE STRUCTURE DEFINITION PROBLEM

What is a structure?

THE STRUCTURE DEFINITION PROBLEM

The screenshot shows a web browser window with the title "Online Dictionary of Crystallography". The main content area displays the article for "Crystal structure". The article includes a brief definition, multiple language equivalents, and a "See also" section with links to related topics like Point space, Noncrystallographic symmetry, and International Tables of Crystallography. A "Category" box indicates the topic is "Fundamental crystallography". At the bottom of the page, there is a note about the last edit date and a link to the Creative Commons Attribution license.

Online Dictionary of Crystallography

Navigation ▾

Search Online Dictionary of... Actions

Crystal structure

From Online Dictionary of Crystallography

Structure cristalline (Fr). Kristallstruktur (Ge). Struttura cristallina (It). 結晶構造 (Ja). Estructura cristalina (Sp).

A **crystal structure** is a **crystal pattern** consisting of atoms.

See also

- Point space
- Noncrystallographic symmetry
- International Tables of Crystallography, Volumes A and A1

Category: Fundamental crystallography

This page was last edited on 9 November 2017, at 17:29.
Content is available under [Creative Commons Attribution](#) unless otherwise noted.
[About Online Dictionary of Crystallography](#)

THE STRUCTURE DEFINITION PROBLEM

The screenshot shows a web page from the "Online Dictionary of Crystallography". The header includes the logo, the site name, a navigation menu, a search bar, and user account icons. The main content area has a title "Crystal pattern" with a subtitle "From Online Dictionary of Crystallography". Below the title is a definition: "Motif cristallin (Fr). Unendlicher Idealkristall (Ge). Motivo cristallino (It). 結晶模様 (Ja). Cristal ideal infinito (Sp)." A synonym is listed as "Infinite ideal crystal". A descriptive paragraph follows: "An object in the n -dimensional point space E^n is called an n -dimensional **crystallographic pattern** or, for short, **crystal pattern** if among its symmetry operations:

1. there are n translations, the translation vectors t_1, \dots, t_n of which are linearly independent;
2. all translation vectors, except the zero vector $\mathbf{0}$, have a length of at least $d > 0$.

When the crystal pattern consists of atoms, it takes the name of **crystal structure**. The crystal pattern is thus the generalization of a crystal structure to any pattern, concrete or abstract, in any dimension, which obeys the conditions of periodicity and discreteness expressed above."

See also

- Point space
- Noncrystallographic symmetry
- International Tables for Crystallography, Volumes A and A1

Category: Fundamental crystallography

This page was last edited on 16 November 2018, at 15:33.
Content is available under [Creative Commons Attribution](#) unless otherwise noted.
[About Online Dictionary of Crystallography](#)



THE STRUCTURE DEFINITION PROBLEM

What is a structure, chemist's definition?

- Definition based on space-group classification
- Structures that change without breaking symmetry are the same structure

THE STRUCTURE DEFINITION PROBLEM

What is a structure, chemist's definition?

- Definition based on space-group classification
- Structures that change without breaking symmetry are the same structure
- In this world it is easy to say that two structures are the same (continuous crystallography)
- much harder to say whether two unequal structures are the same in the chemist's sense

THE STRUCTURE DEFINITION PROBLEM

Case 1: thermal expansion

- Take material
- Warm it up
- The unit cell expands (without breaking (point) symmetry)
- The structure hasn't changed (according to Chemists)
- The structure has changed (according to continuous crystallography measures like PDD)
- Actually, it has changed, all the bonds are different lengths which affects properties

THE STRUCTURE DEFINITION PROBLEM

Case 2: minuscule phase transition

- Take material
- Move one atom off a special position by a tiny amount
- The structure has changed (according to Chemists), the space group has changed
- The structure has changed (according to continuous crystallography measures like TPP), but by a minuscule amount

THE STRUCTURE DEFINITION PROBLEM

Case 2: minuscule phase transition

- Take material
- Move one atom off a special position by a tiny amount
- The structure has changed (according to Chemists), the space group has changed
- The structure has changed (according to continuous crystallography measures like TPP), but by a minuscule amount

Case 2 the continuous crystallography thinks the materials are closer than Case 1, but chemists think otherwise

THE STRUCTURE DEFINITION PROBLEM

The screenshot shows the header of the Online Dictionary of Crystallography. It includes the logo, the site name "Online Dictionary of Crystallography", a "Navigation" dropdown, a search bar with placeholder "Search Online Dictionary of...", and user icons for "Actions".

Isostructural crystals

From Online Dictionary of Crystallography

Cristaux isotropiques (Fr). Isotypic Kristalle (Ge). Cristalli isostrutturali (It). 同形結晶 (Ja). Cristales isotípicos (Sp).

Definition

Two crystals are said to be *isostructural* if they have the same structure, but not necessarily the same cell dimensions nor the same chemical composition, and with a 'comparable' variability in the atomic coordinates to that of the cell dimensions and chemical composition. For instance, calcite CaCO_3 , sodium nitrate NaNO_3 and iron borate FeBO_3 are isostructural. One also speaks of *isostructural series*, or of *isostructural polymorphs* or *isostructural phase transitions*.

The term **isotypic** is synonymous with isostructural.

See also

- Chapter 3.3 of *International Tables for Crystallography, Volume D*

Category: [Crystal chemistry](#)

This page was last edited on 26 March 2019, at 13:56.

Content is available under [Creative Commons Attribution](#) unless otherwise noted.

[About Online Dictionary of Crystallography](#)

THE STRUCTURE DEFINITION PROBLEM

Summary

As a community we need more practical definitions of what it means for two structures to be the same or similar

WAYS TO COMPARE STRUCTURES: STRUCTURE METRICS

How do crystallographers compare structures?

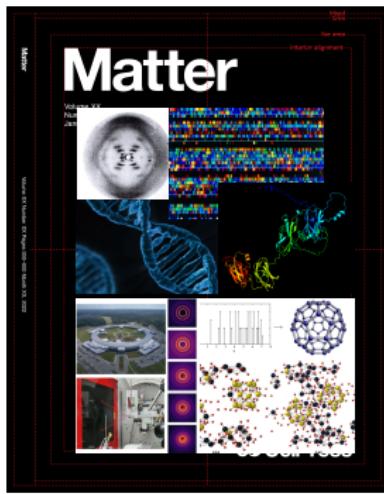
- Gold standard is R_w between measured and computed structure-factors or diffraction data coupled with crystallographic expertise
- CIFcheck (does some symmetry and chemistry common sense)

WAYS TO COMPARE STRUCTURES: CONTINUOUS STRUCTURE METRICS

continuous structure representations

- PDF as a “materials genome”, a 1D function that “codes for” 3D structure
- AMD (reduce the PDF to a scalar) - metric of Vitaliy Kurlin (U Liverpool)
- PDD (matrix of partial PDFs)
- smooth overlap of atomic positions (SOAP)
- atom-centered symmetry functions (ACSF)

AMD and PDD are not experimentally accessible but are valid metrics, and PDD is complete



COMPARING STRUCTUREMATCHER AND PDF

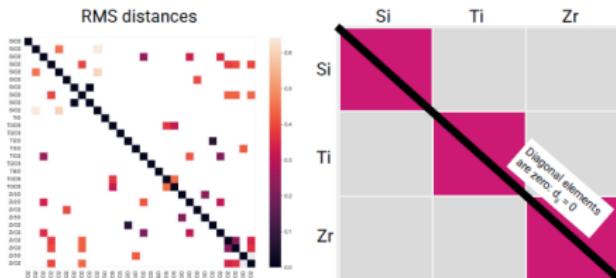
Baseline: pymatgen.StructureMatcher module

- Pymatgen: open-source Python library for materials analysis
- StructureMatcher module compares structures based off their atomic coordinates.

$$\text{distance}(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

Distance measure:

- Root-mean-square distance between the atomic positions of two structures
- Ignores atomic species, only accounts for atomic sites and lattice parameters
- Permutes atomic sites to minimize RMS, account for rotations & symmetries



Pair Distribution Function (PDF):

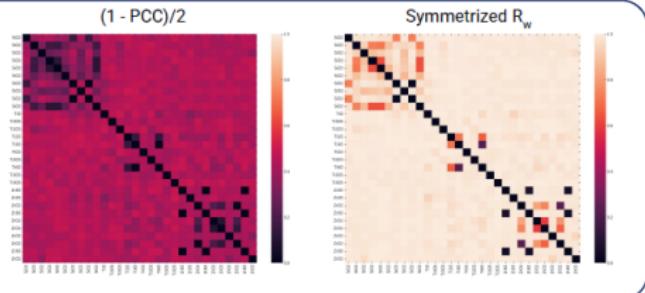
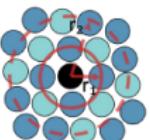
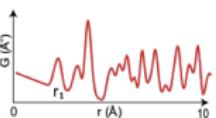
"histogram" of bond lengths in the structure, experimentally measurable from diffraction data.

Distance measures:

- Pearson correlation coefficient (PCC)
- Symmetrized R_w

$$R_w = \sum_i (x_i - \bar{x})^2$$

*Pairs of PDFs are "morphed" before computing PCC and R_w . This accounts for differences in scale between structures, which affects R_w , PCC greatly.



Work of Tina Na Narong, Zoe Zachko and Andrew Yang, Collaboration with Steven Torrisi

COMPARING PDF AND PDD/AMD

Pair Distribution Function (PDF):

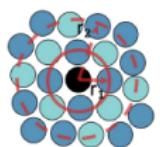
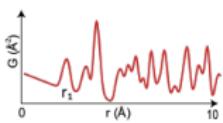
"histogram" of bond lengths in the structure, experimentally measurable from diffraction data.

Distance measures:

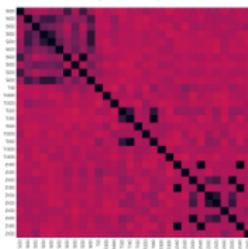
- Pearson correlation coefficient (PCC)
- Symmetrized R_w

$$R \sum_i (x_{i(\vec{r})} - \bar{x}_{\vec{r}})^2$$

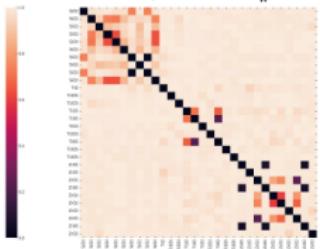
*Pairs of PDFs are "morphed" before computing PCC and R_w . This accounts for differences in scale between structures, which affects R_w , PCC greatly.



(1 - PCC)/2



Symmetrized R_w



Pointwise Distance Distribution (PDD)

- Matrix of interatomic distances in the structure (isometry invariant)
- Each row of a PDD matrix lists k shortest distances (ascending order) from an atom in the unit cell to other atoms

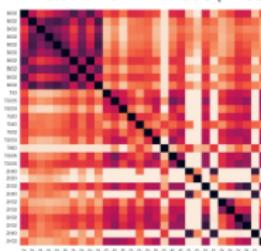
Distance measure:

Earth's Mover Distance (EMD)
→ treat PDD matrix elements as "piles" of Earth
→ solve a linear programming problem to quantify the minimum "mass transport work" needed to make the two PDDs identical

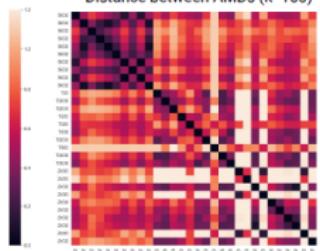
Atom in unit cell	Distances to neighbors				
	1	2	3	...	k
1	1	1	1.2	...	5

AMD = Average bond distances
average of PDD rows

Distance between PDDs ($k=100$)



Distance between AMDs ($k=100$)



Average Minimum Distance (AMD)

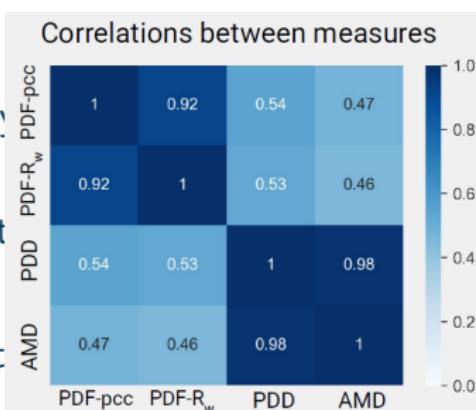
- AMD is an average of PDD: vector of k average shortest distances

Distance measure:

Chebyshev distance

COMPARING STRUCTURE METRICS

- different similarity metrics differ in which structures they think are similar
- Which is the rightest
- We don't know
- How similar of result do they provide?



Correlations with StructureMatcher
(only distances below threshold included)

PDF-pcc	0.72
PDF-R _w	0.83
PDD	0.63
AMD	0.55

SUMMARY

- Diffusion models are the s\$%t
- They can be used to make really great foundation models
- As with LLMs, we will likely all be using these in the future
- But:

SUMMARY

- Diffusion models are the s\$%t
- They can be used to make really great foundation models
- As with LLMs, we will likely all be using these in the future
- But: be careful what you wish for...the answers are only as good as the question that was asked
- Bias in the data, bias in the network/noising protocol, choice of success metric
- These things are still being figured out

With GenAI, the problem is the problem, not the solution

ACKNOWLEDGEMENTS



- My current and former students and post-docs
- Beamline and software teams
- Collaborators
- Funding (DOE-BES, NSF-DMR, TRI, Columbia-DSI)
- The Facilities!