Deep learning to automate phase identification from multi-phase XRD spectra



Nathan Szymanski, Ceder group at UC Berkeley

Preprint: <u>https://arxiv.org/abs/2103.16664</u> Github: <u>https://github.com/njszym/XRD-AutoAnalyzer</u> Email: nathan_szymanski@berkeley.edu



XRD: a cornerstone of inorganic materials research





Crystal structure \leftrightarrow X-ray diffraction pattern

XRD: a *chemical fingerprint* for phase ID



Characterization for exploratory syntheses:



What's in the sample? Solving this problem widely remains a manual task

Ludwig, npj Comp. Mat. (2019).

Grid of samples with distinct compositions

1) High-throughput and parallelized experiments

Robotic hardware is accelerating experimentation – can software keep up?





A growing need to reliably automate XRD analysis



Robotic hardware is accelerating experimentation – can software keep up?

2) Self-driving laboratories for closed-loop experiments

- Experiments are guided by previous results
- Goal is to remove human from the loop

Can we make inorganic materials synthesis a rapid and autonomous process?

Reliable phase identification is critical to learn from each synthesis attempt



Peak search-and-match Reference || FoM Measured ٠ Decomposition into *d-I* list Measured spectrum • Ο Ο Ο

2θ

Intensity

Approach:

- Extract peak positions and intensities
- Compare with known materials (e.g., from the ICSD or ICDD)

Limitations:

- Peak extraction is unreliable when:
 - Peaks overlap
 - Peaks "blend in" with noise
 - Impurities are present



Full-profile analysis

Measured spectrum & Simulated spectrum



Approach:

- Simulate spectra of known phases
- Quantify overlap between measured and simulated spectra

Limitations:

- Patterns become dissimilar when experimental artifacts are present, common artifacts include:
 - Strain, texture, small particle size, poor crystallinity

An improved approach: convolutional neural networks



Convolutional neural networks (CNNs) are widely used for image classification

Treat XRD pattern as a 1-D image



Recent work suggests that CNNs outperform traditional methods for:

- Symmetry classification
- Single-phase identification

Park et al., IUCrJ (2017). Oviedo et al., npj Comp. Mat. (2019). Lee et al., Nat. Commun. (2020). Maffettone et al., Nat. Comp. Sci. (2021).

1) Physics-informed data augmentation

2) Accounting for non-stoichiometry

3) A probabilistic treatment of multi-phase mixtures

Training the CNN using simulated XRD spectra

What data is used to train the CNN?

- Limited number of experimental XRD spectra
- Luckily, XRD spectra are easily simulated

Issues to consider:

- One spectrum per phase is not enough
- Ideal spectrum may not reflect experiment

Solution: perform *data augmentation*

Key question: How to augment simulated spectra?



Simulated XRD spectrum



Physics-informed data augmentation

Derive perturbations to spectra:

- 1) Shifts in peak positions \leftrightarrow strain in the unit cell:
- 2) Variations in peak intensities \leftrightarrow texture in the powder:
- 3) Broadening of peak widths \leftrightarrow small domain size:







Derive perturbations to spectra:

1) Shifts in peak positions \leftrightarrow strain in the unit cell:

2) Variations in peak intensities \leftrightarrow texture in the powder:

Intensity scaling $\propto \langle hkl \rangle_{\text{preferred}} \cdot [h'k'l']_{\text{peak}}$

Bounds: $\pm 50\%$ intensity

3) Broadening of peak widths \leftrightarrow small domain size:





Derive perturbations to spectra:

1) Shifts in peak positions \leftrightarrow strain in the unit cell:

2) Variations in peak intensities \leftrightarrow texture in the powder:

3) Broadening of peak widths \leftrightarrow small domain size:

Scherrer equation:
$$\tau = \frac{K\lambda}{\beta cos\theta}$$

FWHM (β) \propto^{-1} domain size (τ)

Bounds: 1 nm (broad) to 100 nm (narrow)





For a given reference phase:

50 spectra simulated from each artifact; changes sampled from normal distribution

ightarrow 150 augmented spectra (per phase) to train the CNN





Constant and

Li-Mn-Ti-O-F:

- Useful system for many battery materials
- Challenging test case for diffraction due to similarity in Mn/Ti and O/F scattering factors

Training data: 140 phases extracted from the ICSD \rightarrow 21,000 simulated spectra



Ensemble approach to yield probabilistic predictions

Test results on single-phase spectra

From the simulated spectra:

• 80/20 split for training/testing

Test results from the CNN:

• 94% of all test spectra are correctly classified

Accuracy is promising, but what is our baseline?

Test results on single-phase spectra

JADE from MDI:

• Peak search-match + full-profile comparison

Results from JADE:

• Only 78% of all spectra are correctly classified

Test results: our CNN outperforms traditional methods

Test results: probabilities provide insight into accuracy

Probabilities provide a measure of confidence:

- Correct/incorrect predictions are clearly distinguished by their probabilities
- High probability \rightarrow prediction is reliable
- Low probability \rightarrow use caution

1) Physics-informed data augmentation

2) Accounting for non-stoichiometry

3) A probabilistic treatment of multi-phase mixtures

Data augmentation accounts for minor perturbations from idealized reference patterns However, larger changes occur when composition deviates from expected stoichiometry

Perturbative treatment of artifacts

Beyond perturbations

Can we just use whatever solid solutions are available on the ICSD?

• <u>No</u>: the ICSD covers narrow regions of chemical space while leaving others sparse – bias toward highly-studied materials

FIZ Karlsruhe

Creating hypothetical solid solutions

- 1) Enumerate pairs of stoichiometric phases
- 2) For each pair, check whether the structures are isomorphic
- 3) Check whether equivalent sites contain ions comparable in size ($\leq 15\%$)

- 4) If both criteria are satisfied, assume solubility is possible
- 5) Interpolate a grid of hypothetical solid solutions by assuming **Vegard's law** holds for:
 - Lattice parameters
 - Atomic positions

$$= \text{ e.g., Mn} - \text{Mn}_{0.75}\text{Ti}_{0.25} - \text{Mn}_{0.5}\text{Ti}_{0.5} - \text{Mn}_{0.25}\text{Ti}_{0.75} - \text{Ti}_{0.75} - \text{Ti}_{0.75}$$

Site occupancies

Testing on non-stoichiometric materials: Li-Mn-Ti-O-F

Training set:

• Hypothetical solid solutions

Test set:

• Real non-stoichiometric phases from the ICSD

20 experimental structures with unique compositions

Classifications:

• **Classifications are discrete** – can only predict what we've trained on

Quantifying performance based on two metrics:

- **Structure:** is the predicted structure isomorphic to the true structure?
- **Composition:** mole fraction error between true and predicted compositions

Experimental compositions

Testing on non-stoichiometric materials: structure

Before including hypothetical solid solutions:

• **11/20 (55%)** of structures correctly identified

After introducing hypothetical solid solutions:

• **19/20 (95%)** of structures correctly identified

Including non-stoichiometry \rightarrow 40% improvement

Training set:

Without NS: 16,800 spectra
➤ Non-stoichiometric (hypothetical)
With NS: 13,800 spectra
✓ Non-stoichiometric (hypothetical)

Testing on non-stoichiometric materials: compositions

1.0Composition is more difficult to predict: (mole fraction) Lattice parameters often follow Vegard's Peak positions \rightarrow structure classification is reliable 8.0 The basis may not follow Vegard's law; Including NS Peak some ions can swap sites in the structure still reduces intensities \rightarrow limited accuracy for composition 0.6 errors composition Example: 0.4 Li/TM ion swap in spinel LiMn_{1-x}Ti_xO₄ 0.2 Error

0.0

Without NS

With NS

This algorithm is **not for refinement**

1) Physics-informed data augmentation

2) Accounting for non-stoichiometry

3) A probabilistic treatment of multi-phase mixtures

Possible methods for multi-phase identification

Approach #1 ^a:

- Choose phases with high probabilities
- Problem: model may confuse similar reference phases

Approach #2 ^b:

- Train on simulated multi-phase spectra
- Problem: linear combinations of spectra
 → combinatorial data explosion

a) Lee *et al.,* Nat. Commun. (2020). b) Maffettone *et al.,* Nat. Comp. Sci. (2021).

Possible methods for multi-phase identification

Approach #1:

- Choose phases with high probabilities
- Problem: model may confuse similar reference phases

Our solution:

Use probability to guide an iterative approach of phase identification and profile subtraction

Plug measured spectrum into CNN to identify the first phase

Fit the identified peaks to the measured spectrum

How to use probability as a guiding metric?

 \rightarrow First step is over-prioritized

A branching algorithm to maximize confidence

True phases: $Li_2TiO_3 + Mn_3O_4 + Li_2O$

Test results for multi-phase classification

Experimental procedure	Anticipated artifact	CNN	JADE	
Single-phase				
Pristine samples	None	10/10	9/10	
Kapton tape overlaid	Diffuse baseline	9/10	8/10	
Rapid XRD scan	Noisy baseline	10/10	7/10	
Thick samples	Shifts in 20	5/6	2/6	
Ball milled	Broadening	5/5	4/5	
Partially disordered	Intensity variation	5/6	4/6	
Solid solutions	Non-stoichiometry	4/4	3/4	
Multi-phase				
Two-phase mixtures	None	10/10	7/10	Good match with
Three-phase mixtures	None	13/15	9/15	simulated tests
	Overall accuracy:	71/76 (93.4%)	53/76 (71.4%)	

- Our approach can be **generalized to any arbitrary composition space**
- Code to automate data augmentation and model training:

https://github.com/njszym/XRD-AutoAnalyzer

• Only requires a set of reference phases

Preprint: https://arxiv.org/abs/2103.16664

Email: nathan_szymanski@berkeley.edu

Acknowledgements

Professor Gerbrand Ceder

Chris Bartel

Yan Zeng

Howard Tu

Preprint: <u>https://arxiv.org/abs/2103.16664</u> Github: <u>https://github.com/njszym/XRD-AutoAnalyzer</u> Email: nathan_szymanski@berkeley.edu Twitter: @NJSzymanski

