

Automated Extraction and Database Creation for Ferroelectric Material Discovery

Ian Moog¹, Michael Wang², Brandon Schoener¹ (Advisor), Yuting Hu³, Dr. Hao Zeng¹, Dr. Jinjun Xiong^{3,4}

¹ Department of Physics, ² Williamsville East HS, ³ CSE Department, ⁴ IAD

Introduction

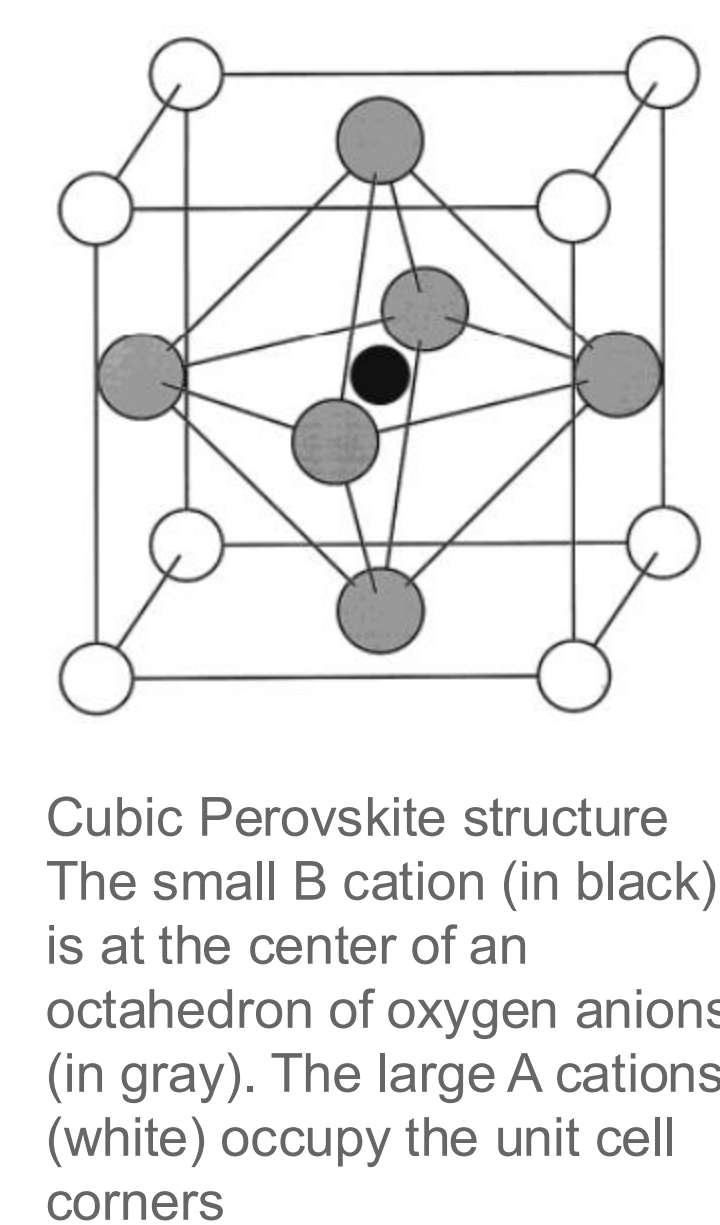
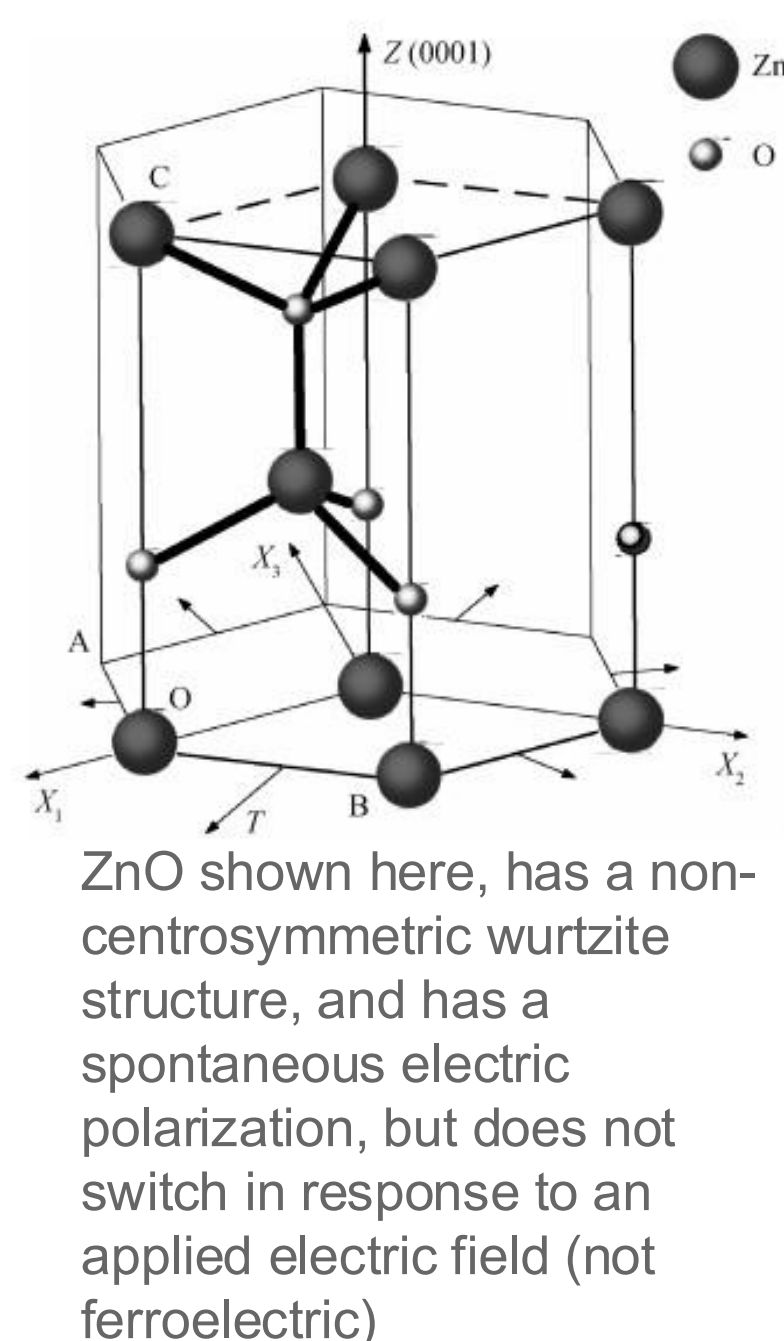
Ferroelectric materials are crystals that maintain an electric polarization even without an external field and can reverse polarity when a field is applied. Their coupled piezoelectric and pyroelectric properties enable applications in non-volatile memory, sensors, actuators, energy harvesting, tunable optics, or even neuromorphic computing, where their intrinsic non-volatility and polarization switching behavior can mimic synaptic activity for brain-like learning.

What Causes Ferroelectricity?

For ferroelectricity to arise in a material, it must have a non-centrosymmetric arrangement of ions and their corresponding electrons.

Many ferroelectric structures can be understood as derivatives of a *non-polar, centrosymmetric prototype phase*, the most widely studied being the perovskite structure.

Long-range Coulomb forces that favor polar distortion, compete with short-range repulsive forces, which favor symmetry. At lower temperatures, Coulomb forces tend to dominate, they can cause distortion in the center cation, pushing it off center. This is why ferroelectricity will only be present below a certain Curie temperature.



Methods

- Query OpenAlex + Unpaywall for URLs containing PDF links to open access papers.
- Direct download + Selenium download attempts.
- Conversion of PDFs to .md files via MinerU.

LLM Guided Property Extraction

- Utilize a long context model (Gemini-Flash-2.5) to analyze the .md files
- Return a structured response containing extracted properties

Example PDF Abstract + MinerU conversion

ABSTRACT
This study combines first-principles Density Functional Theory (DFT) calculations and experimental analysis to investigate the impact of Sr-doping on the structural, electronic, and piezoelectric properties of $\text{Ba}_{1-x}\text{Sr}_x\text{TiO}_3$ (BST) perovskites at different Sr concentrations ($x = 0, 0.125, 0.25$ theoretically; $x = 0, 0.1, 0.125, 0.15$, experimentally).

ABSTRACT
This study combines first- principles Density Functional Theory (DFT) calculations and experimental analysis to investigate the impact of Sr-doping on the structural, electronic, and piezoelectric properties of $\text{Ba}_{1-x}\text{Sr}_x\text{TiO}_3$ (BST) perovskites at different Sr concentrations ($x = 0, 0.125, 0.15, 0.2, 0.25$ experimentally).

Snippet from a .md file which can be passed to an LLM for extraction

Why use an ML Model?

The space of possible crystal structures is essentially infinite, and first-principles calculations like DFT are too slow to explore it fully. A machine learning model allows property prediction of these candidate structures.

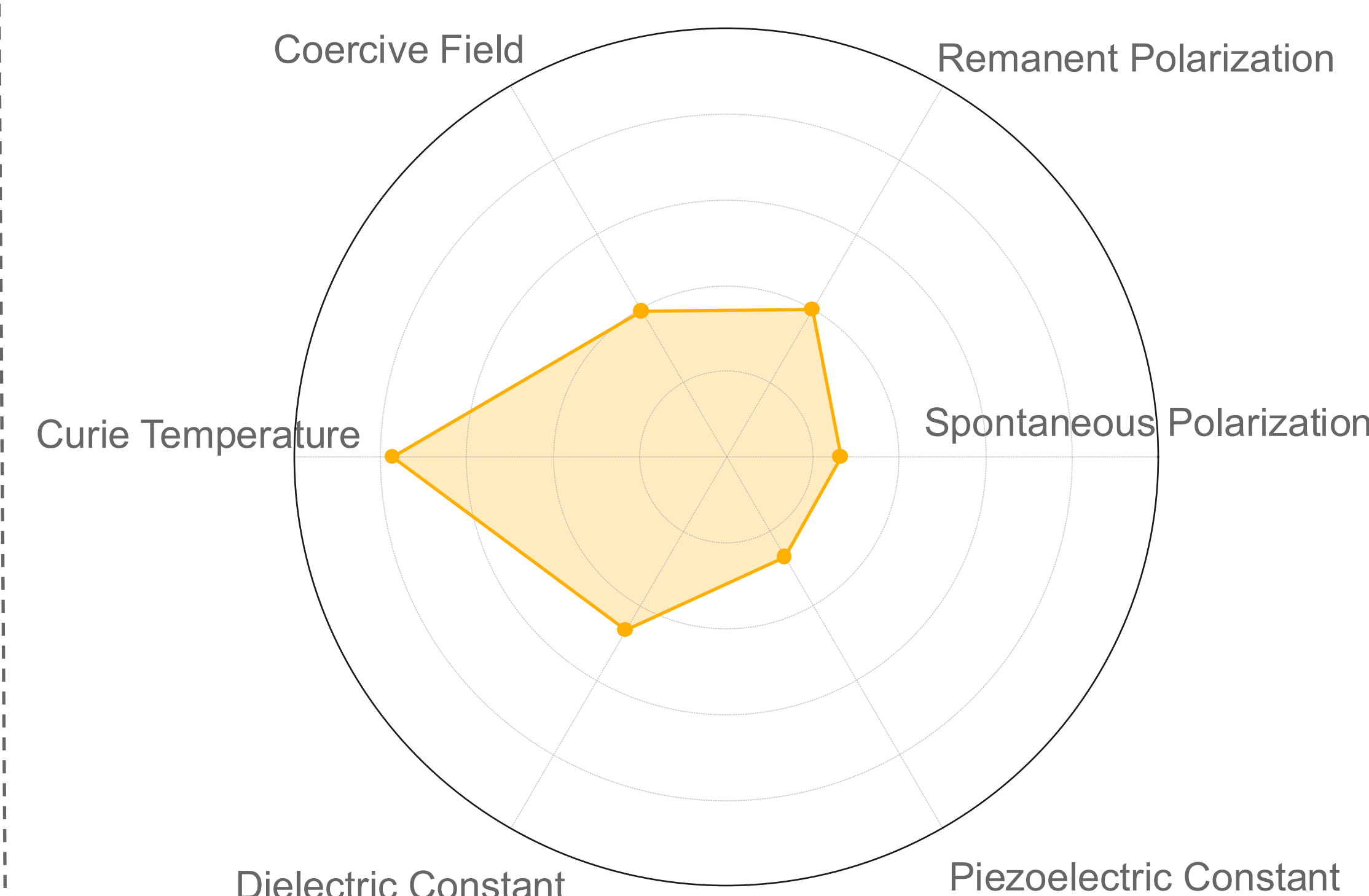
This enables a potential workflow:

- Generate** new candidate structures (e.g., with a crystal GAN).
- Screen** them rapidly with the ML model for desirable properties.
- Validate** only the most promising ones through detailed simulations or experiments

Extraction Results

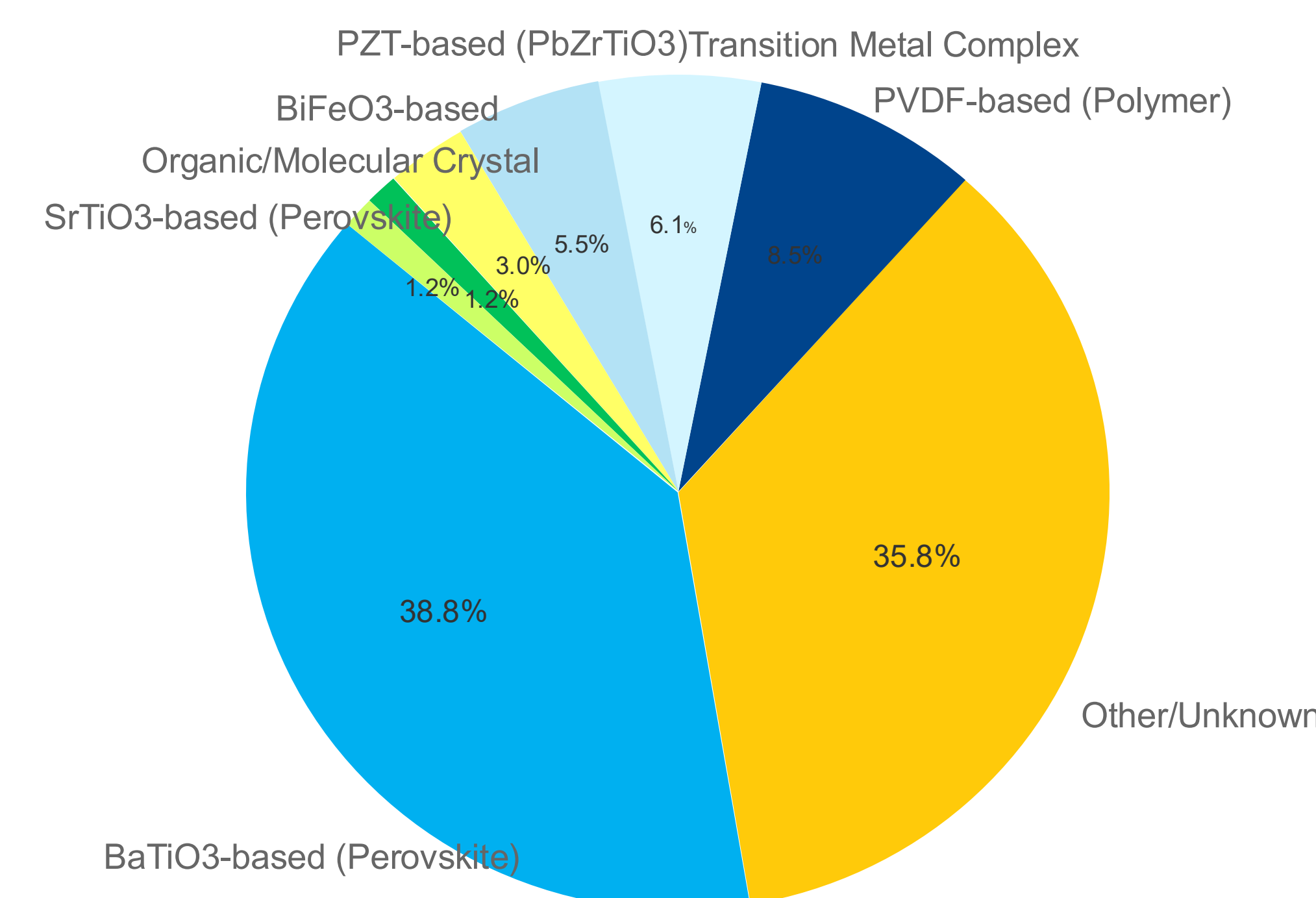
- From testing on 100 papers, we identified 165 unique ferroelectric materials
- Curie temperature was the most reported (77%)

Extraction Coverage by Property



This figure shows the extraction success rate by property, i.e., how often each property was extracted.

Distribution of Material Types in Extracted Dataset



This figure shows the distribution of material types that we extracted. We can see that BaTiO3 based Perovskites were the most recorded type of ferroelectric, which is expected. This is beneficial because the physics is well understood, and future discoveries will likely be of this type.

Prompting Techniques

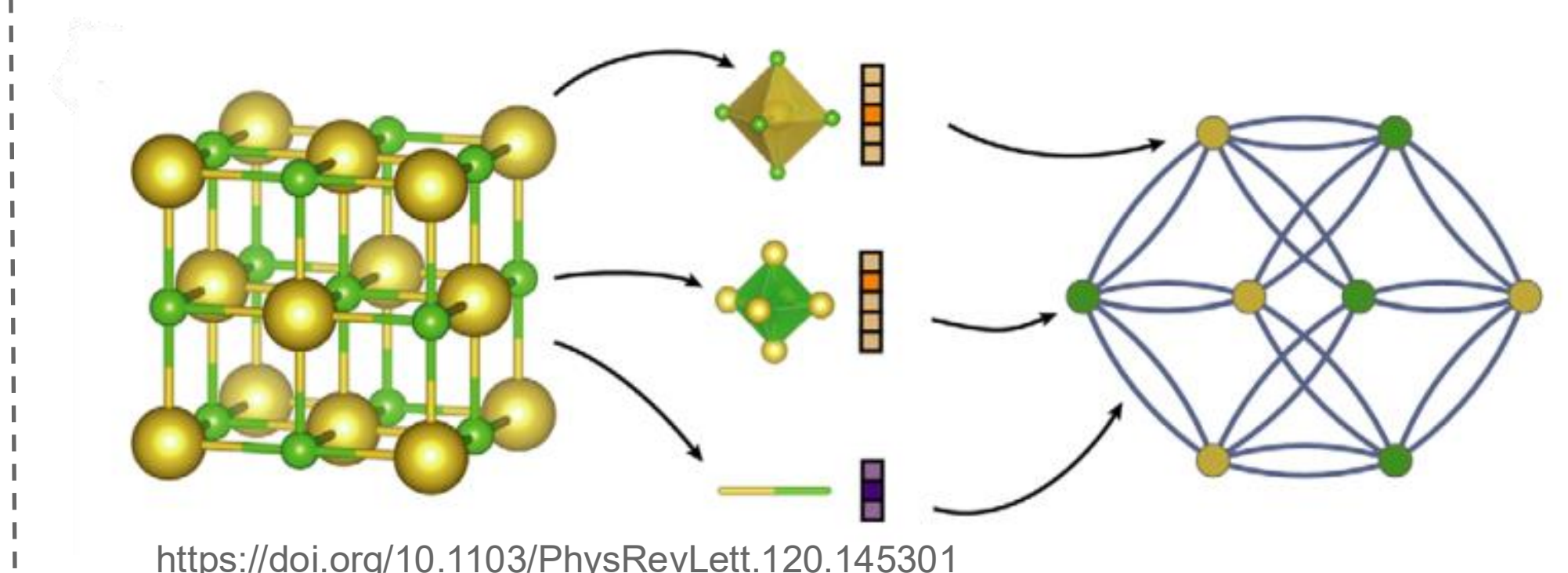
Due to the long context window (1 million tokens) of Gemini-Flash-2.5 we can safely pass our files. Prompting is straightforward, the API allows a schema to be defined, so it will only output one csv line. We simply describe the properties, their units, probable locations, etc. And the model can figure out the rest. Here is an example output:

```
BaTiO3,,14 μC/cm²,,2.2 kV/cm,,130 °C,,220 pC/N,DFT-based-and-experimental-study-on-Sr-doped-BaTiO3
Ba0.95Sr0.1TiO3,,16 μC/cm²,,2.4 kV/cm,,100 °C,,180 pC/N,DFT-based-and-experimental-study-on-Sr-doped-BaTiO3
Ba0.875Sr0.125TiO3,,17.5 μC/cm²,,2.1 kV/cm,,90 °C,,170 pC/N,DFT-based-and-experimental-study-on-Sr-doped-BaTiO3
```

As you can see, we successfully extracted properties for three materials from one paper, with 4/6 properties for each. We can flag extractions with low confidence, and send them to a higher cost model later, which will perform better.

Property Prediction

To predict properties, we can use a Crystal Graph Convolutional Neural Network (CGCNN). A graph is a data type made of nodes connected by edges. Crystal structures can be represented this way, where atoms are the nodes, and the distance between them can be represented by the edges. So, our CGCNN can take a crystal structure as an input and give property values as output.



Pipeline for Material Discovery

- Match crystal structures (CIF files) to our extracted properties
- Use this as input to train a CGCNN to predict properties from structure alone
- Once enough structure-property relationships are learned we can generate and feed it entirely new crystal structures.
- Use our model to rank these structures, and find potential candidates to test experimentally

References

- T. Xie & J.C. Grossman, Phys. Rev. Lett. 120, 145301 (2018)
- K. Rabe, Ch. H. Ahn & J.-M. Triscone, Springer (2007).
- A. Jain et al., APL Mater. 1, 011002 (2013).
- B. Schoener, UB, Comprehensive Exam