

Materials Synthesis Insights from Scientific Literature via Text Extraction and Machine Learning

Isha Kumari



- Background
- Methodology
- Results
- Conclusion
- Future Outlook



Background



Why This Matters: The Need for Data-Driven Synthesis

- The pace of discovering new materials is often slowed by the lack of efficient synthesis strategies.
- Traditionally, synthesis relies on heuristics, which do not scale with the rapid growth of computational material predictions
- This project leverages Natural Language Processing (NLP) and Machine Learning (ML) to extract and analyze synthesis parameters from thousands of published scientific articles.

By building a large, structured database of synthesis conditions, we can:

- Understand key parameters influencing material properties
- Improve reproducibility and guide new syntheses
- Enable data-driven automation in materials design





Framework for Literature-Based Synthesis Extraction













Article Retrieval

Plain-Text Conversion & Classification



Neural Net Word **Dependency Parsing** Parameter Extraction Material Detection & (SpaCy + Labeling (86%) via Tree Traversal Validation ChemDataExtractor) Accuracy) Noun phrases (e.g., A trained model Sentences are LiOH, ethanol) are From the verb, the classifies each word converted into parse tree is traversed to as Material, matched against trees. Key synthesis PubChem, an n-Operation, Condition, extract synthesis verbs (e.g., *sinter*, conditions (e.g., or Amount using gram classifier dissolve) are temperature, time, (82% accuracy), and embeddings from identified as root 5000+ annotated ChemDataExtractor stirring speed). nodes. examples.

Parsing and Extraction

9





Verification of Annotated Data Data Mining and Machine Learning







Calcination Temperature Distribution Across Metal Oxide Systems (Bulk vs Nano)

Bulk vs Nano Trend:

Bulk materials are typically calcined at **higher temperatures** than their nanostructured counterparts.

Complexity Increases Temperature:

Calcination temperature increases with **elemental complexity**, from binaries to pentanaries, especially in **nano** syntheses.

Binary Oxide Range:

Most binary oxides fall between 450–550° C

- Alumina: found more in bulk
- ZnO: more in nano form

Ternary oxide Range:

- **Bismuth Ferrite:** sharp peak at 750° c(bulk) and around 600° C for nano.
- Barium Titanate: wider range, 900-1100 °C

•Hydrothermal reactions occur at low temperatures (typically 150– 200 °C) and long durations (12–24 hours).

Reaction temperatures are limited by the critical point of common solvents (e.g., water, ethanol).
Calcination steps happen at much higher temperatures and usually shorter times.

•As material complexity increases, calcination temperatures and durations both rise.

•Simple systems concentrate around 400–500 °C, <5 h, while complex ones extend to 800–900 °C, >5 h.











Decision Tree for Titania Nanotube Formation

Key Features Identified:

NaOH concentration (most important, root node)

is commonly clustered around **1 M and 10 M**

Hydrothermal temperature also plays a significant role Some notable peaks at 150° C and 180° C

Annealing time was found to be less predictive, as its values didn't clearly distinguish between outcomes.

The model achieves 82% classification accuracy, showing that NaOH concentration and hydrothermal temperature are strong indicators of nanotube formation.





The diagram shows higher probability of nanotube formation (darker green) at:

- High NaOH concentrations
- Lower hydrothermal temperatures

This reduced parameter space allows for easy visualization of synthesis trends using only experimentally accessible variables.

Machine-Learned Phase Diagram for Titania Nanotube Formation





Machine-learned classifiers and predictions across materials systems

All four subfigures compare three

- Support Vector Machine (SVM)
- Heuristic (simple rule-based)

SVM consistently outperforms

learning models across materials and synthesis outcomes.



Conclusion

Conclusion



This work bridges the gap between published knowledge and predictive synthesis planning.

Built an automated framework to extract synthesis data using text mining and ML.

Compiled a **large**scale database of synthesis conditions from scientific literature. Identified **key parameters** (e.g., temp, time, solvent) that influence material outcomes.

Enabled accurate predictions without prior domain knowledge or manual input.



Future Outlook



Extend coverage to more synthesis methods

Thin film deposition, catalytic systems, non-crystalline materials

Refine data accuracy

Weighting highly cited papers Incorporating newly published research



Thank You

20 June 2025



Any Questions?