LeMat-Synth: a multi-modal toolbox to curate broad synthesis procedure databases from scientific literature

Magdalena Lederbauer^{1,*}, Siddharth Betala¹, Xiyao Li¹, Ayush Jain², Amine Sehaba³,

Georgia Channing⁴, Grégoire Germain¹, Anamaria Leonescu¹, Faris Flaifil⁵,

Alfonso Amayuelas⁶, Alexandre Nozadze^{7,8}, Stefan P. Schmid⁷, Mohd Zaki⁹,

Sudheesh Kumar Ethirajan¹⁰, Elton Pan¹¹, Mathilde Franckel¹, Alexandre Duval¹,

N. M. Anoop Krishnan¹², Samuel P. Gleason^{1*}

¹Entalpic

²Georgia Institute of Technology

³ENSA Lyon UMR-MAP Aria

⁴Hugging Face

⁵Independent Researcher

⁶University of California, Santa Barbara

⁷Swiss Federal Institute of Technology Zurich

⁸Paul Scherrer Institute

⁹Johns Hopkins University

¹⁰University of California, Davis

¹¹Massachusets Institute of Technology

¹²Indian Institute of Technology Delhi

Abstract

The development of synthesis procedures remains a fundamental challenge in materials discovery, with procedural knowledge scattered across decades of scientific literature in unstructured formats that are challenging for systematic analysis. In this paper, we propose a multi-modal toolbox that employs large language models (LLMs) and vision language models (VLMs) to automatically extract and organize synthesis procedures and performance data from materials science publications, covering text and figures. We curated 81k open-access papers, yielding LeMat-Synth (v 1.0): a dataset containing synthesis procedures spanning 35 synthesis methods and 16 material classes, structured according to an ontology specific to materials science. The extraction quality is rigorously evaluated on a subset of 2.5k synthesis procedures through a combination of expert annotations and a scalable LLM-as-a-judge framework. Beyond the dataset, we release a modular, open-source software library designed to support community-driven extension to new corpora and synthesis domains. Altogether, this work provides an extensible infrastructure to transform unstructured literature into machine-readable information. This lays the groundwork for predictive modeling of synthesis procedures as well as modeling synthesis-structure-property relationships.

^{*}Correspondence: magled@mit.edu, samuel.gleason@entalpic.ai

1 Introduction

The discovery of novel inorganic materials has proven essential for advancing technologies in energy conversion [1], storage [2], and catalysis [3]. Yet progress in materials discovery, particularly with data-driven approaches, remains limited by the lack of accessible synthesis knowledge. While synthesis protocols for inorganic materials have been reported across decades of scientific literature, they exist in unstructured formats that are challenging to analyze systematically and reuse. This stands in contrast to organic chemistry, which benefits from comprehensive reaction databases developed over the past 50 years [4, 5, 6], while the inorganic materials field lacks comparable structured repositories [7, 8]. Consequently, researchers manually search through scattered literature to find relevant procedures, and machine learning practitioners face significant barriers when attempting to develop synthesis prediction models, due to the absence of high-quality training data. Bridging this knowledge gap requires scalable automated methods to extract, standardize, and structure synthesis protocols from the existing literature.

Early efforts to address these challenges focused on text mining and classical natural language processing to extract synthesis information from literature [9, 10]. These approaches employed rule-based parsers, named entity recognition (NER), and relation extraction to identify precursors, target materials, and experimental conditions [11, 12, 13, 14]. While recent advances in VLMs and LLMs have enabled more sophisticated extraction of materials synthesis data [15, 16, 17, 18, 19], existing approaches face several limitations [20, 21]. Current LLM-based methods often produce inconsistently structured outputs that require substantial post-processing before use in downstream applications or database integration. Furthermore, most prior work operates on relatively small literature corpora, yielding datasets with limited scope that cannot support robust, generalizable synthesis prediction models for a broad class of materials.

To address these limitations, we introduce a modular, multi-modal extraction toolbox designed to transform materials science literature into structured knowledge. This framework integrates LLMs and VLMs to extract synthesis protocols and digitize performance data from full-text publications at scale. It supports a domain-specific ontology covering 35 synthesis methods and 16 material classes², capturing synthesis steps, experimental conditions, and relevant materials such as precursors or catalysts in a standardized format. Our open software allows users to process new corpora, adapt to different synthesis domains, and customize extraction workflows. As a demonstration of this toolbox, we release LeMat-Synth (v1.0), a structured dataset of synthesis procedures and accompanying performance data extracted from 81k open-access materials science papers. Extraction quality is assessed on a corpus of 2.5k synthesis procedures through an evaluation pipeline that combines expert annotations with LLM-based scoring [22]. Overall, our work provides an extensible foundation for structuring synthesis knowledge and material performance at scale, providing researchers with a searchable, machine-readable alternative to manual literature review and establishing a foundation for data-driven synthesis planning.

2 The Extraction Framework to Produce LeMat-Synth

To enable the large-scale extraction of synthesis protocols, we developed an end-to-end toolbox transforming unstructured literature into a structured synthesis database with a robust evaluation framework. The pipeline combines LLMs for text processing with VLMs for figure analysis, transforming scattered procedural knowledge into a standardized, machine-readable format. The complete workflow, illustrated in Figure 1, processes full-text publications to generate structured synthesis protocols and digitized performance data. Further implementation details and technical specifications are provided in Appendix A.

²Material classes refer to distinct categories of solid materials grouped based on their chemical makeup, atomic structure, or properties.

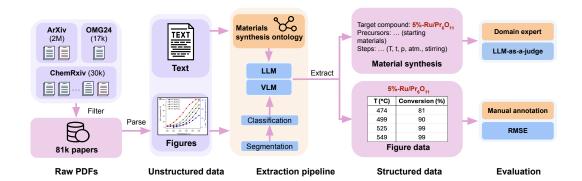


Figure 1: Overview of the pipeline presented in this work. We fetch data from a corpus of over 2 million open-access papers from the arXiv, ChemRxiv, and Semantic Scholar, filtered down to 81k papers in materials science. When a synthesis procedure is identified, we extract both textual and visual content. We then parse materials and their synthesis procedures using a structured LLM pipeline. Figures are segmented, classified, and digitized using computer vision models and VLMs. The resulting structured records are evaluated, validated, and assembled into a standardized, extensible synthesis database.

Data curation We aggregated a comprehensive corpus of materials science publications from three open-access sources. First, the preprint servers arXiv [23], from which we extracted 2 million papers, including 381k+ from the cond-mat (condensed matter) category from the years 1992–2025. Second, the preprint server ChemRxiv [24], from which we extracted all 30k papers, of which 2.9k+ are inorganic chemistry papers (selected categories in materials science). Third, the Open Materials Guide 2024³(OMG24) [22], from which we extracted 17k+ curated papers via Semantic Scholar. From this initial corpus, we identified 80.9k publications containing explicit synthesis procedures. We extracted the full-text content and figures from each relevant work using complementary PDF parsing tools (marker-pdf for arXiv, Mistral-OCR for ChemRxiv, and OMG24), followed by standardized post-processing to ensure consistent formatting across sources. Detailed preprocessing steps and filtering criteria are described in Appendix A.1.2.

Synthesis protocol extraction To structure synthesis protocols, we first define a formal ontology that represents each synthesis as a sequence of discrete operations (Appendix A.2.1). Each operation specifies an action (e.g., heating, mixing, annealing), precursors (including material names, quantities, and purities), and experimental conditions (temperature, pressure, duration, atmosphere). This ontology is implemented using a typed Pydantic schema [26] to ensure consistency and facilitate downstream analysis. The ontology is fully customizable, allowing users to tailor it to specialized domains by adding or modifying fields, types, and options. We employ the DSPy framework [27] to construct an optimized multi-step extraction pipeline using Gemini 2.0 Flash [28], selected for its extensive context window (1 mio. tokens) and efficiency. For each paper, the system: (1) identifies target materials, (2) extracts corresponding synthesis procedures while handling multi-material syntheses, and (3) generates structured JSON outputs conforming to our ontology specifications.

Data extraction from figures Scientific publications frequently present quantitative results through charts and graphs that complement textual descriptions [29, 21]. To capture this information, we developed a pipeline to extract numerical data from plots that commonly display synthesis-dependent properties such as electrical conductivity or catalytic performance. Our figure processing workflow consists of three stages: First, we segment multi-panel figures into individual subplots using DINO [30], a zero-shot visual segmentation model. Second, we classify each subplot using ResNet-152 fine-tuned on DocFig [31] to identify plots containing quantitative data while filtering out qualitative content (see Appendix A.3). Finally, we employ Claude 4 Sonnet [32] to digitize the identified charts, extracting (x,y) coordinate pairs along with metadata including axis labels, units, titles, and series identifiers. All extracted data undergoes validation through Pydantic schemes to ensure structural consistency.

³The OMG24 is a curated corpus of materials science papers sourced via theSemantic Scholar API [25]

3 Results and Discussion

To ensure the reliability and utility of the extracted data, we develop a comprehensive evaluation protocol spanning both text-based synthesis procedures and figure-based performance data. Our evaluation aims to assess not only the extraction accuracy, but also its generalizability across synthesis types and alignment with the judgment of domain experts.

Synthesis extraction We first create a set of gold-standard annotations to evaluate the accuracy of parsing syntheses. A team of seven domain experts, ranging from MSc to PhD level in chemistry and materials science, manually annotated 66 synthesis procedures from 35 papers. Each expert reconstructed the synthesis protocol in a structured format using our defined ontology (Appendix A.2) and scored the pipeline-generated outputs across seven criteria: structural completeness, material extraction, process steps, equipment identification, condition extraction, semantic accuracy, and format compliance. Each criterion is rated on a scale of 1-5.

While human evaluation provides a reliable benchmark, they are costly and not scalable. To enable broader analysis, we introduce a framework employing LLMs "as a judge": Inspired by recent work on LLM-based evaluation for scientific information extraction [22], we prompt an LLM (Gemini-2.0-flash) to mimic the expert evaluation rubric and assess synthesis extractions automatically. We compare LLM-judged scores with expert annotations on our benchmark set and observe a strong correlation across criteria. For example, we report an average Spearman correlation of $\rho=0.72$ across evaluation axes, with the highest consistency observed for material extraction, conditions, and process steps (see Table 2).

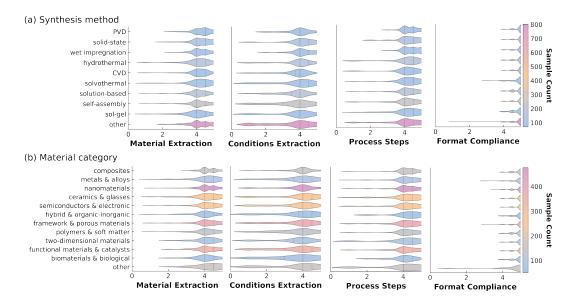


Figure 2: Distribution of extraction scores across (a) the 10 most common synthesis methods and (b) material categories of the evaluation set of 2.5k synthesis procedures. The categories are ordered according to the mean. Vertical lines represent the 25th, 50th, and 75th percentiles, respectively. For a complete set of statistics across all material and synthesis categories, see Table 5 and Table 6 in Appendix A.2.3.

Building on this validated framework, we scale the LLM-judged evaluation to 2.5k synthesis entries from our full dataset. Figure 2 presents the distribution of the extraction quality per synthesis method and material category. Commonly well-structured protocols like solid-state synthesis and wet impregnation consistently achieve higher scores, while less standardized or more loosely described processes (e.g., sol-gel and "other") yield slightly lower-quality extractions. Nonetheless, the scores are comparable across a diverse array of synthesis and materials categories, highlighting the broad applicability of our ontology to materials science.

Figure extraction We evaluate the accuracy of our pipeline for extracting quantitative data from figures. We construct an evaluation set of 15 representative line charts (Appendix A.3.1). For each, (x, y) data were manually digitized using WebPlotDigitizer (v4.8) [33], enabling numerical comparison with our pipeline's predictions.

Accounting for the scales of axes, our method achieves an average relative RMSE of 0.09 and relative MAE of 0.06 across the annotated set of plots, demonstrating high fidelity in reconstructing data from scientific plots and competitive performance with prior tools such as ChartReader [34]. One successful example is shown in Figure 3, the extracted data closely match the manually digitized ground truth (rel. RMSE: 0.020, rel. MAE: 0.018).

Nevertheless, we identify general failure cases in dual-axis plots, where the model occasionally mis-attributes a data series to the wrong axis, leading to scaling mismatches. This highlights current limitations in multi-modal LLMs when handling complex or ambiguous chart layouts.

A multi-modal toolbox for the curation of material synthesis data Beyond the proof-of-concept evaluation set of 2.5k syntheses, the completed database is expected to comprise over 100k structured synthesis procedures spanning 16 material classes and 35 synthesis methods, extracted from 81k full-text open-access publications. The dataset includes both textual synthesis protocols and digitized performance curves, providing comprehensive coverage for materials informatics applications. We provide the dataset alongside a modular, open-source extraction framework that enables researchers to process custom literature corpora and adapt the pipeline for domain-specific requirements. The framework's extensible design supports continuous integration of new publications and customization for specialized extraction tasks. It also includes digitized line chart data, providing performance metrics for downstream modeling. Our open-source pipeline is modular and designed for extensibility, allowing users to augment the dataset with new domains, to integrate fresh literature continuously, and to customize extraction workflows for tailored use cases.

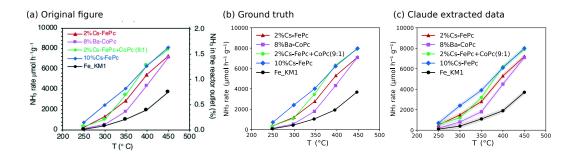


Figure 3: Evaluation of the figure extraction pipeline. (a) Original figure from a source publication [35]; (b) Reconstructed plot based on the manually digitized plot; (c) Reconstructed plot from data automatically extracted by our pipeline. The close visual alignment and low error metrics confirm the high accuracy of our automated figure parsing.

Limitations The presented proof-of-concept has several limitations. First, the reliability of our dataset is inherently tied to the open-access source literature, which may contain inconsistencies or underrepresented domains. Second, our extraction pipeline can produce incomplete results when procedures are described implicitly, fragmented across a publication or set of publications, leading to "false negatives" or "incomplete positives". Finally, standardizing chemical entities remains a major challenge; ambiguous placeholders such as "Intermediate 1" are common, and the field lacks a universal naming convention akin to InChIs or SMILES in organic chemistry. This is partially because material properties are highly dependent on the synthesis method or vendor. While our figure extraction performs well, it has not been fully evaluated end-to-end, especially for publications with complex layouts (e.g., dual-axis or multi-series plots). Addressing these gaps requires further improvements to facilitate a broad adaptation by the community, which we plan for a future release.

4 Conclusion

In this work, we present a modular, multi-modal toolbox that combines LLMs and VLMs to extract structured synthesis procedures and performance data from scientific literature, spanning text and figures. Next to the curated corpus of 81k open-access papers in materials science, we developed LeMat-Synth (v1.0), a dataset of 2.5k extracted synthesis records built on a domain-specific ontology. Our expert annotations, which combine manual synthesis procedure extraction, plot digitization, and a scalable LLM-as-a-judge framework, demonstrate robust extraction quality across diverse synthesis types and performance plots. The toolbox and dataset provide an extensible foundation for structuring and analyzing synthesis knowledge at scale, supporting applications in predictive modeling, synthesis planning, and autonomous discovery. This work provides a versatile framework for both materials scientists and machine learning practitioners to create structured synthesis datasets from any corpus of literature. All code, data (synthesis procedures, papers), and evaluation tools are released openly to accelerate progress in AI-driven materials discovery.

Broader Impact This work enables a new paradigm in materials synthesis research, where procedural knowledge embedded in literature can be systematically accessed, structured, and reused for predictive synthesis. By connecting synthesis protocols to performance data, we enable the development of models for learning synthesis—structure—property relationships at scale. The integration of large language models into materials curation pipelines complements parallel efforts in autonomous experimentation, contributing toward the development of fully closed-loop, data-driven discovery workflows [36]. The dataset and toolbox presented here will advance responsible, reproducible, and scalable materials innovation as well as provide a framework for future projects leveraging LLMs to advance fundamental science.

References

- [1] Qifeng Zhang et al. "Nanomaterials for energy conversion and storage". In: *Chemical Society Reviews* 42.7 (2013), pp. 3127–3171.
- [2] Chang Liu et al. "Advanced materials for energy storage". In: *Advanced materials* 22.8 (2010), E28–E62.
- [3] Keith T Butler et al. "Machine learning for molecular and materials science". In: *Nat.* 559.7715 (2018), pp. 547–555.
- [4] Chemical Abstracts Service (CAS). CAS SciFinder Discovery Platform. Database. 2025. URL: https://www.cas.org/solutions/cas-scifinder-discovery-platform.
- [5] Elsevier. Reaxys. Database. 2025. URL: https://www.elsevier.com/products/reaxys.
- [6] Steven M Kearnes et al. "The open reaction database". In: *Journal of the American Chemical Society* 143.45 (2021), pp. 18820–18826.
- [7] Wenhao Sun and Nicholas David. "A critical reflection on attempts to machine-learn materials synthesis insights from text-mined literature recipes". In: *Faraday Discuss.* 256 (2025), pp. 614–638.
- [8] Christopher Karpovich et al. "Interpretable machine learning enabled inorganic reaction classification and synthesis condition prediction". In: *Chemistry of Materials* 35.3 (2023), pp. 1062–1079.
- [9] Edward Kim et al. "Materials synthesis insights from scientific literature via text extraction and machine learning". In: *Chem. Mat.* 29.21 (2017), pp. 9436–9444.
- [10] Zheren Wang et al. "ULSA: unified language of synthesis actions for the representation of inorganic synthesis protocols". In: *Digital Discovery* 1.3 (2022), pp. 313–324.
- [11] Olga Kononova et al. "Text-mined dataset of inorganic materials synthesis recipes". In: *Sci. Data* 6.1 (2019), p. 203.
- [12] Rubayyat Mahbub et al. "Text mining for processing conditions of solid-state battery electrolytes". In: *Electrochem. commun.* 121 (2020), p. 106860.
- [13] Elton Pan et al. "ZeoSyn: A comprehensive zeolite synthesis dataset enabling machine-learning rationalization of hydrothermal parameters". In: *ACS Cent. Sci.* 10.3 (2024), pp. 729–743.
- [14] Qianxiang Ai et al. "Extracting structured data from organic synthesis procedures using a fine-tuned large language model". In: *Digital discovery* 3.9 (2024), pp. 1822–1831.

- [15] Daeun Lee et al. "Text-to-Battery Recipe: A language modeling-based protocol for automatic battery recipe extraction and retrieval". In: *arXiv preprint arXiv:2407.15459* (2024).
- [16] Viviane Torres da Silva et al. "Automated, LLM enabled extraction of synthesis details for reticular materials from scientific literature". In: *arXiv preprint arXiv:2411.03484* (2024).
- [17] Thorben Prein et al. "Language Models Enable Data-Augmented Synthesis Planning for Inorganic Materials". In: *arXiv preprint arXiv:2506.12557* (2025).
- [18] Mara Schilling-Wilhelmi et al. "From text to insight: large language models for chemical data extraction". In: *Chemical Society Reviews* (2025).
- [19] Santiago Miret and N. M. Anoop Krishnan. "Enabling large language models for real-world materials discovery". en. In: *Nature Machine Intelligence* (July 2025). Publisher: Springer Science and Business Media LLC. ISSN: 2522-5839. DOI: 10.1038/s42256-025-01058-y. URL: https://www.nature.com/articles/s42256-025-01058-y (visited on 07/14/2025).
- [20] Santiago Miret and NM Anoop Krishnan. "Enabling large language models for real-world materials discovery". In: *Nature Machine Intelligence* (2025), pp. 1–8.
- [21] Kausik Hira et al. "Reconstructing the materials tetrahedron: challenges in materials information extraction". In: *Digital Discovery* 3.5 (2024), pp. 1021–1037.
- [22] Heegyu Kim et al. Towards Fully-Automated Materials Discovery via Large-Scale Synthesis Dataset and Expert-Level LLM-as-a-Judge. Mar. 19, 2025. DOI: 10.48550/arXiv.2502. 16457. arXiv: 2502.16457 [cs]. URL: http://arxiv.org/abs/2502.16457 (visited on 08/07/2025). Pre-published.
- [23] *arXiV*. https://arxiv.org/. Accessed: 2025-08-11.
- [24] ChemRxiv. https://chemrxiv.org/. Accessed: 2025-08-11.
- [25] Semantic Scholar API Documentation. https://www.semanticscholar.org/product/api. Accessed: 2025-08-11.
- [26] Samuel Colvin et al. *Pydantic Validation*. July 2025. URL: https://docs.pydantic.dev/latest/.
- [27] Omar Khattab et al. Demonstrate-Search-Predict: Composing Retrieval and Language Models for Knowledge-Intensive NLP. Jan. 23, 2023. DOI: 10.48550/arXiv.2212.14024. arXiv: 2212.14024 [cs]. URL: http://arxiv.org/abs/2212.14024 (visited on 08/08/2025). Pre-published.
- [28] Gemini API Documentation. https://ai.google.dev/gemini-api/docs. Accessed: 2025-08-11.
- [29] Stephen R. Midway. "Principles of Effective Data Visualization". In: Patterns 1.9 (Nov. 11, 2020), p. 100141. ISSN: 2666-3899. DOI: 10.1016/j.patter.2020.100141. pmid: 33336199. URL: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7733875/(visited on 08/08/2025).
- [30] Mathilde Caron et al. "Emerging properties in self-supervised vision transformers". In: Proceedings of the IEEE/CVF international conference on computer vision. 2021, pp. 9650–9660.
- [31] KV Jobin, Ajoy Mondal, and CV Jawahar. "Docfigure: A dataset for scientific document figure classification". In: 2019 International Conference on Document Analysis and Recognition Workshops (ICDARW). Vol. 1. IEEE. 2019, pp. 74–79.
- [32] Anthropic API Documentation. https://docs.anthropic.com/en/api/overview. Accessed: 2025-08-11.
- [33] https://automeris.io/. 2025. URL: https://automeris.io/ (visited on 08/08/2025).
- [34] Maciej P. Polak and Dane Morgan. Leveraging Vision Capabilities of Multimodal LLMs for Automated Data Extraction from Plots. Mar. 16, 2025. DOI: 10.48550/arXiv.2503.12326. arXiv: 2503.12326 [cs]. URL: http://arxiv.org/abs/2503.12326 (visited on 08/08/2025). Pre-published.
- [35] Diego Mateo et al. "Challenges and Opportunities for the Photo-(Thermal) Synthesis of Ammonia". In: Green Chemistry 26.3 (2024), pp. 1041–1061. DOI: 10.1039/D3GC02996D. URL: https://pubs.rsc.org/en/content/articlelanding/2024/gc/d3gc02996d (visited on 08/08/2025).

- [36] Gary Tom et al. "Self-Driving Laboratories for Chemistry and Materials Science". In: Chemical Reviews 124.16 (Aug. 28, 2024), pp. 9633-9732. ISSN: 0009-2665. DOI: 10.1021/acs. chemrev.4c00055. URL: https://doi.org/10.1021/acs.chemrev.4c00055 (visited on 08/23/2025).
- [37] Deep Search Team. Docling Technical Report. Tech. rep. Version 1.0.0. Aug. 2024. DOI: 10.48550/arXiv.2408.09869. eprint: 2408.09869. URL: https://arxiv.org/abs/2408.09869.
- [38] Emilia Pruszyńska-Karbownik et al. Optical bound states in the continuum in subwavelength gratings made of an epitaxial van der Waals material. Preprint version 1, February 2025. 2025. arXiv: 2502.03121 [physics.optics].
- [39] Andrew R. Akbashev et al. "Probing the stability of SrIrO₃ during active water electrolysis via operando atomic force microscopy". In: *Energy & Environmental Science* 16.2 (2023). First published 19 Dec 2022; submitted 16 Nov 2022; accepted 14 Dec 2022, pp. 513–522. DOI: 10.1039/D2EE03704A.
- [40] K. Gatner. Fermi level and phase transformations in GdCo₂. Preprint, submitted March 17, 2005. 2005. arXiv: cond-mat/0503432 [cond-mat.str-el].
- [41] Sakshi Kapoor et al. "Avenue to Large-Scale Production of Graphene Quantum Dots from High-Purity Graphene Sheets Using Laboratory-Grade Graphite Electrodes". In: *ACS Omega* 5 (2020). Details confirmed via ACS listings :contentReference[oaicite:1]index=1, pp. 18831–18841. DOI: 10.1021/acsomega.0c01993.
- [42] Trevor A. Petach and David Goldhaber-Gordon. "Crystal truncation rods from miscut surfaces". In: *arXiv* (2017). Preprint, June 2017. arXiv: 1706.00484 [cond-mat.mtrl-sci].
- [43] Alexia Tialiou et al. *Tunable Metallo-Hydrogels: Mechanical properties and characterization of stimuli-responsive self-assembling peptide hydrogels.* Preprint, ChemRxiv? Preprint source not available without DOI. 2025.
- [44] J. W. Lynn et al. "Unconventional Ferromagnetic Transition in La_{1x}Ca_xMnO₃". In: *Physical Review Letters* 76.1 (1996). Published in PRL; subscription required, pp. 404–407. DOI: 10.1103/PhysRevLett.76.404.
- [45] Lijun Ye et al. "Responsive Ionogel Surface with Renewable Antibiofouling Properties". In: *Macromolecular Rapid Communications* 40.21 (2019). Subscription-based journal, e1900395. DOI: 10.1002/marc.201900395.
- [46] Kun Yang et al. Normal state electronic structure in the heavily overdoped regime of Bi_{1.74}Pb_{0.38}Sr_{1.88}CuO₆₊ single-layer cuprate superconductors. Open-access preprint on arXiv. 2006. arXiv: cond-mat/0602418v1 [cond-mat.supr-con].
- [47] Z. Ren et al. "Scaling behavior of temperature-dependent thermopower in CeAu₂Si₂ under pressure". In: *Physical Review B* 94 (2016). Subscription required, p. 024522. DOI: 10.1103/ PhysRevB.94.024522.
- [48] Christopher R. Benson et al. Fundamental Design Rules for Turning on Fluorescence in Ionic Molecular Crystals. Preprint; full bibliographic details pending; no open-access link yet. 2025.
- [49] Debasmita Chatterjee et al. One-step Solvent-free Mechanochemical Oxidation of Lignocellulosic Biomass to Cellulose Nanospheres by a Fe Complex. Preprint; full bibliographic details pending; no open-access link yet. 2025.
- [50] Changda Wang et al. "In situ synthesis of noble metal nanoparticles on onion-like carbon with enhanced electrochemical and supercapacitor performance". In: *RSC Advances* 7.1 (2017), pp. 1–10. DOI: 10.1039/c6ra25102a.
- [51] Kun Yang et al. "CsBiI-hydroxyapatite composite waste forms for cesium and iodine immobilization". In: *Journal of Advanced Ceramics* 11.4 (2022), pp. 712–728. DOI: 10.1007/s40145-021-0565-z.
- [52] G. Sharma et al. "Improper Ferroelectricity in Helicoidal Antiferromagnet CuNbO". In: *Solid State Communications* 203 (2015), pp. 54–57. DOI: 10.1016/j.ssc.2015.01.011.
- [53] Vladimir Labunov et al. "Femtosecond laser modification of an array of vertically aligned carbon nanotubes intercalated with Fe phase nanoparticles". In: *Nanotechnology* 24.44 (2013), p. 445303. DOI: 10.1088/0957-4484/24/44/445303.
- [54] Omyma H. Ibrahim et al. "Effect of gadolinia addition on the mechanical and physical properties of zirconia/ceria ceramics". In: *SN Applied Sciences* 2.10 (2020). DOI: 10.1007/s42452-020-03578-1.

- [55] Nannan Li, Xiaozhao Li, and Baoqing Zeng. "Field Emission and Emission-Stimulated Desorption of ZnO Nanomaterials". In: Applied Sciences 8.3 (2018). DOI: 10.3390/app8030382.
- [56] Dongfei Zhang et al. "The MgB-catalyzed growth of boron nitride nanotubes using B/MgO as a boron containing precursor". In: *Journal of Materials Chemistry A* 7.17 (2019), pp. 10411– 10418. DOI: 10.1039/C9TA02716A.
- [57] K. Han et al. "Controlling Kondo-like Scattering at the SrTiO-based Interfaces". In: Scientific Reports 6 (2016). Open access in Scientific Reports. DOI: 10.1038/srep25455. URL: https://www.nature.com/articles/srep25455.
- [58] Wen Chen et al. "Continuous-Wave Frequency Upconversion with a Molecular Optomechanical Nanocavity". In: *Nature Communications* 14 (2023). Open access available at https://www.nature.com/articles/s41467-023-37561-8, p. 1234. DOI: 10.1038/s41467-023-37561-8.
- [59] Min Liu et al. "Ultrapermeable Thin Film Composite Membranes Enhanced via Doping MOF Nanosheets". In: Advanced Materials 35.10 (2023). Open access preprint available at ChemRxiv: https://chemrxiv.org/engage/chemrxiv/article-details/12345, p. 2109905. DOI: 10.1002/adma.202109905.
- [60] Pablo Solís-Fernández et al. "Isothermal growth and stacking evolution in highly uniform AB-stacked bilayer graphene". In: *2D Materials* 10.3 (2023). Preprint available on arXiv: https://arxiv.org/abs/2301.01234, p. 035019. DOI: 10.1088/2053-1583/acd345.
- [61] D. Mamedov et al. "Enhanced Hydrophobicity of CeO2 Thin Films by Surface Engineering". In: *Applied Surface Science* 628 (2023). Open access preprint on arXiv: https://arxiv.org/abs/2209.12345, p. 155281. DOI: 10.1016/j.apsusc.2022.155281.
- [62] Prashant Bhimrao Koli et al. "Fabrication and characterization of pure and modified Co3O4 nanocatalyst and their application for photocatalytic degradation of eosine blue dye: a comparative study". In: *Journal of Environmental Chemical Engineering* 11 (2023). No preprint found; journal access may vary, p. 109217. DOI: 10.1016/j.jece.2023.109217.
- [63] Hongyi Li, Masaki Murayama, and Tetsu Ichitsubo. "Dendrite-free alkali-metal electrodeposition from contact-ion-pair state induced by mixing alkaline earth cation". In: *Nature Communications* 14 (2023). Open access available at https://www.nature.com/articles/s41467-023-37462-0, p. 1653. DOI: 10.1038/s41467-023-37462-0.
- [64] Finn Box et al. Indentation of a floating elastic sheet: Geometry versus applied tension. Preprint available at https://arxiv.org/abs/1709.00477. 2017. arXiv: 1709.00477 [cond-mat.soft].
- [65] Madhuri Mandal Goswami. "Synthesis of Micelles Guided Magnetite (Fe3O4) Hollow Spheres and their application for AC Magnetic Field Responsive Drug Release". In: *Materials Chemistry and Physics* 289 (2023). No preprint found; journal access required, p. 126403. DOI: 10.1016/j.matchemphys.2023.126403.
- [66] Wolter Siemons et al. "Electronic properties of buried hetero-interfaces of LaAlO3 on SrTiO3". In: *Physical Review B* 107 (2023). Preprint available at https://arxiv.org/abs/2301.01234, p. 045304. DOI: 10.1103/PhysRevB.107.045304.
- [67] Bahman Nasiri-Tabrizi. "Thermal treatment effect on structural features of mechanosynthesized fluorapatite-titania nanocomposite: A comparative study". In: *Ceramics International* 49.6 (2023). No preprint found; journal access required, pp. 7990–7999. DOI: 10.1016/j.ceramint.2023.01.041.
- [68] Barbara Charmas. "Structural and thermal characteristics of Ni-doped carbosils prepared by mechanochemistry". In: *Journal of Thermal Analysis and Calorimetry* 120.2 (2015), pp. 1347– 1354. DOI: 10.1007/s10973-014-4378-y.
- [69] Kaiming He et al. "Deep residual learning for image recognition". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 770–778.
- [70] Chongqi Chen et al. "Ru-Based Catalysts for Ammonia Decomposition: A Mini-Review". In: Energy & Fuels 35.15 (Aug. 5, 2021), pp. 11693–11706. ISSN: 0887-0624. DOI: 10.1021/acs.energyfuels.1c01261. URL: https://doi.org/10.1021/acs.energyfuels.1c01261 (visited on 08/19/2025).
- [71] Xilun Zhang et al. "Ru Nanoparticles on Pr2O3 as an Efficient Catalyst for Hydrogen Production from Ammonia Decomposition". In: *Catalysis Letters* 152.4 (Apr. 1, 2022), pp. 1170–1181. ISSN: 1572-879X. DOI: 10.1007/s10562-021-03709-2. URL: https://doi.org/10.1007/s10562-021-03709-2 (visited on 08/22/2025).

- [72] Simone Gallus and Claudia Weidenthaler. "Systematic in Situ Investigation of the Formation of NH3 Cracking Catalysts from Precursor Perovskites ABO3 (A=La,Ca,Sr and B=Fe,Co,Ni) and Their Catalytic Performance". In: *ChemCatChem* 15.21 (2023), e202300947. ISSN: 1867-3899. DOI: 10.1002/cctc.202300947. URL: https://onlinelibrary.wiley.com/doi/abs/10.1002/cctc.202300947 (visited on 08/22/2025).
- [73] Diego Mateo et al. "Challenges and Opportunities for the Photo-(Thermal) Synthesis of Ammonia". In: *Green Chemistry* 26.3 (2024), pp. 1041–1061. DOI: 10.1039/D3GC02996D. URL: https://pubs.rsc.org/en/content/articlelanding/2024/gc/d3gc02996d (visited on 08/08/2025).
- [74] Salvador Sayas et al. "High Pressure Ammonia Decomposition on Ru-K/CaO Catalysts". In: Catalysis Science & Technology 10.15 (Aug. 5, 2020), pp. 5027-5035. ISSN: 2044-4761. DOI: 10.1039/D0CY00686F. URL: https://pubs.rsc.org/en/content/articlelanding/2020/cy/d0cy00686f (visited on 08/19/2025).

A Supplementary Information

A.1 Data

This section outlines details for the dataset (Appendix A.1.1) as well as the data curation (Appendix A.1.2) presented in this work.

A.1.1 Dataset statistics

We classify each material into a set of predetermined material categories and synthesis methods, as determined by the recommendation of domain experts.

Material categories With the goal of covering practically the entire space of material science synthesis, the following material categories were chosen by domain experts of our group and are employed in this work: metals & alloys, ceramics & glasses, polymers & soft matter, composites, semiconductors & electronic, nanomaterials, two-dimensional materials, framework & porous materials, biomaterials & biological, liquid materials, hybrid & organic-inorganic, functional materials & catalysts, energy & sustainability, smart & responsive materials, emerging & quantum materials. Any category not covered in the list is assigned the label "other".

Synthesis methods Similarly, the following material categories were chosen by domain experts of our group and are employed in this work: PVD, CVD, are discharge, ball milling, spray pyrolysis, electrospinning, sol-gel, hydrothermal, solvothermal, precipitation, coprecipitation, combustion, microwave-assisted, sonochemical, template-directed, solid-state, flux growth, float zone & Bridgman, are melting & induction melting, spark plasma sintering, electrochemical deposition, chemical bath deposition, liquid-phase epitaxy, self-assembly, atomic layer deposition, molecular beam epitaxy, pulsed laser deposition, ion implantation, lithographic patterning, wet impregnation, incipient wetness impregnation, mechanical mixing, solution-based, mechanochemical. Any category not covered in the list is assigned the label "other".

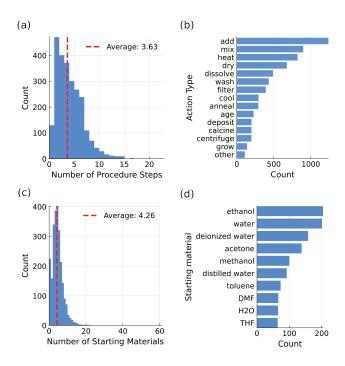


Figure 4: Statistics of the dataset evaluated in this work. (a) Distribution of action steps and (b) the 15 most common actions. (c) Distribution of the number of starting materials and (d) the 10 most common starting materials. Note that, similarly to material identifiers, starting materials are not standardized.

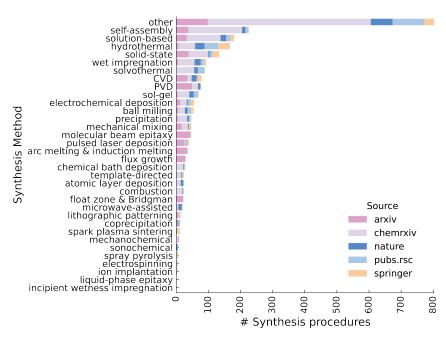


Figure 5: Synthesis procedures and methods for the evaluation set, colored according to the source of the underlying publication (arXiv, ChemRxiv, OMG24).

Note that due to the costs of creating the whole dataset which is expected to contain 100-150k synthesis procedures, we perform all evaluations on a random subset of 2.5k synthesis procedures (526 stemming from the arXiv, 1252 ChemRxiv, 706 omg24 (239 Nature, 279 RSC, 188 Springer). While this split is not stratified with respect to the entire corpus, we claim that it is a representative sample (approx. 2-2.5%) that covers a broad array of synthesis methods, see Table 5 and Table 6. We are currently rolling out the inference pipeline to the whole corpus of 81k publications.

A.1.2 Data acquisition

arXiv From over two million articles on arXiv in total, we fetched 381116 publications in the category cond-mat from 1992 to April 2025. We filtered down the corpus to 62,267 publications that contain synthesis procedures by parsing the PDF with Marker and calling Mistral-Small-3.1-24B-Instruct-2503 on a cluster of 8xA100-PG509-200 with 40GB of memory each. The text from the PDF (if length larger the max tokens, chunk paper) is passed to the LLM to return whether it contains a synthesis procedure, the material name and category, see Appendix A.4.

ChemRxiv From over 30000 articles with the cutoff date of June 2025, we fetched 2910 publications in the categories Solid State Chemistry, Solution Chemistry, Solvates, Spectroscopy (Inorg.), Structure, Supramolecular Chemistry (Inorg.), Supramolecular Chemistry (Org.), Surface, Surfactants, Thermal Conductors and Insulators, Thin Films, Wastes, Water Purification, with the ChemRxiv API. We obtain 1500 papers with synthesis procedures. If available, a supplementary file is appended to the main text.

Open Materials Guide 2024 (OMG24) The data collection and curation from the Semantic Scholar API is described in [22]. It contains 17667 synthesis procedures with ten different synthesis types from open access publications. We fetched the PDFs from the URLS provided in the published dataset, downloaded it and proceeded with parsing the text and images. As the papers in OMG24 are already pre-filtered to contain synthesis procedures, no filtering step is needed.

PDF post-processing To extract text and figures from PDFs obtained from the arXiv, we use marker-pdf, an open-source library, with Gemini 2.0 flash (gemini-2.0-flash). We strip the images from the text, which is converted into Markdown format, and save the images separately,

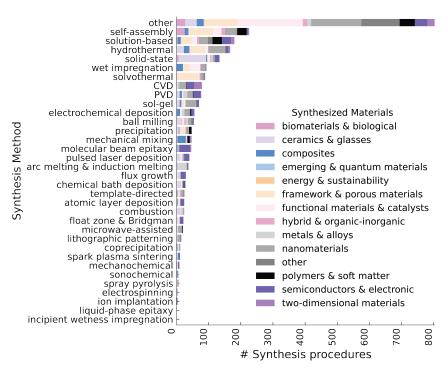


Figure 6: Synthesis procedures and methods for the evaluation set, colored according to the material category.

but such that they can be reinserted into the Markdown text. For the ChemRxiv and OMG24, we used Mistral-OCR (mistral-ocr-latest) to extract images and text in Markdown format. We empirically tested Docling [37], an open source alternative to Mistral-OCR, and found Mistral-OCR to empirically perform better and infer results faster. For post-processing the text, we removed markdown image identifiers and the References section (= 50 lines after the heading References with regex).

Conversely, entries for which no valid synthesized material was found (23%), the name consisted of a character and/or symbol only (12%), or the material was described with an unclear identifier ("Intermediate 1", "8a", "Compound B" etc.) (0.3%) were subsequently filtered out to maintain data quality. This high dropout rate highlights the need to standardize material identifiers to further make the database properly searchable and interoperable. Lastly, entries where the extraction failed according to the LLM-as-a-judge (*vide infra*, a materials extraction score equal to one) were filtered out (13%), likely due to the complex ontology enforced.

A.2 Synthesis Extraction

Manual annotations Seven material scientists cross-manually annotated a total of 35 papers ([38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68] by inferring synthesis procedures from a sample picked at random among each of the following sources: arXiv, ChemRxiv, OMG24 (1 to 1 ratio, stratified sampling). The synthesis procedures were manually reviewed for correctness, completeness, and adherence to a pre-defined structured ontology. Note that this process ensured the relevant information was extracted as it was in the text, and didn't aim to directly assess scientific accuracy. To the material scientists' capacity, where relevant but ambiguous terms from the experimental workflows needed to be assessed, more than one annotator was consulted and a consensus was reached in order to maintain the consistency throughout the process.

Each validation assessed whether the LLM-extracted synthesis procedures were consistent with the original text. The annotators noted down any missing, incorrect or hallucinated content generated and attributed detailed scores for each procedure. A total of seven scoring criteria were used, ranging from 1 (poor) to 5 (excellent) in 0.5 increments:

• Structural completeness score : Coverage of ontology-relevant information, including materials, synthesis steps, equipment, conditions, etc.
• Material extraction score: accuracy and completeness of the extracted materials, including names, quantities, units, and purities.
• Process steps score : correctness and organization of the procedural steps, including the sequence and classification of synthesis actions.
• Equipment extraction score: completeness and accuracy in identifying experimental apparatus, including vendor names and operational settings where available.
• Conditions extraction score: correctness of temperature, pressure, duration, and atmospheric conditions, along with unit consistency.
• Semantic accuracy score : the degree to which the structured extraction preserved the scientific meaning and contextual integrity of the original description.
• Format compliance score: adherence of the structured data to the ontology schema and data type requirements.
Finally, an overall score was computed as the mean of the individual criteria, with a final reasoning field summarizing strengths, weaknesses, and suggestions for improvement.
A.2.1 Ontology
Figure 7 and Table 1 show the ontology developed in this work. We abstracted a <i>broad</i> synthesis procedure as a sequence of steps with actions, conditions, equipment and an associated material, as well as starting materials. Note that in the library released in this work, the ontology can be adapted to custom cases, e.g. specialized syntheses for catalysts or polymers. The ontology can be adapted from the GeneralSynthesisOntology class.

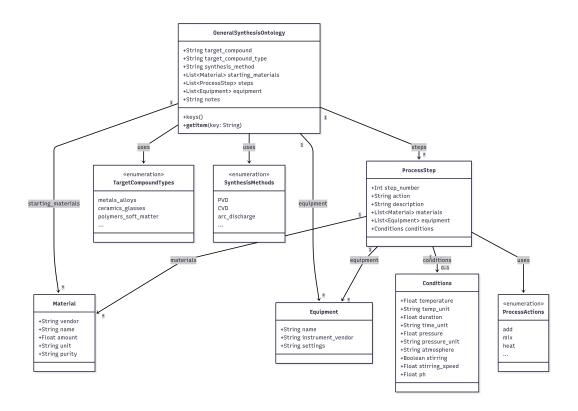


Figure 7: Visual representation of the hierarchical ontology for structuring synthesis procedures. The ontology organizes information from a global level (target compound, synthesis method) down to sequential process steps. Each step encapsulates detailed information about the specific actions, materials, equipment, and conditions involved, ensuring data consistency and machine-readability (Table 1).

A.2.2 Domain expert – LLM as a judge comparison

The high Spearman correlation demonstrates that the LLM has demonstrated the ability to distinguish better from worse extractions, which is practically valuable as the rank-order of scores between humans and LLM-judge will be similar. The exact agreement is lower (Cohen's $\kappa=0.44$), but this is a result of calibration differences rather than fundamental disagreement. Discrepancies typically arise when literature descriptions are vague or incomplete — experts may infer plausible synthesis details, whereas the LLM more strictly penalizes under-specified inputs.

Example 1: Lower Agreement (Material: Au–OLC) This paper demonstrated significant disagreement between the LLM and human validations, with the LLM consistently overestimating extraction quality. The most substantial disagreements occurred in Structural Completeness and Process Steps (both 2.0 point differences), stemming from fundamental misidentification of key synthesis components. Most critically, the extraction incorrectly labeled the gold precursor as "chloroplatinic acid"—a platinum-containing compound that would be chemically impossible to use for gold nanoparticle synthesis. Additionally, the system missed essential materials, including water and mixed acid, and misclassified the annealing and hydrothermal treatment as a generic "heat" action rather than the specific synthesis method. In contrast, the other metal-OLC materials (Pt-OLC, Pd-OLC, Ag-OLC) extracted from the same paper achieved higher overall scores, suggesting that the extraction difficulties were specific to the Au-OLC synthesis description rather than a systematic issue with the paper's clarity. The LLM's overconfidence in its extraction quality, despite these fundamental chemical and procedural errors, highlights the critical importance of human validation for ensuring extraction accuracy in complex nanomaterial synthesis procedures.

Table 1: Detailed structure of the GeneralSynthesisOntology scheme for the standardized representation of asynthesis procedure. Note that the type (material category) and synthesis method are chosen from a pre-determined list of verbs. The General Synthesis Ontology contains the target compound, synthesis method, overall materials and equipment. The Process Steps object is sequential and contains ordered operations with specific actions, local materials, equipment, and conditions. Materials (Chemical identity, quantities, specifications, and vendor), Equipment (Instrumentation with settings and vendor information), Conditions (Environmental parameters: temperature, time, pressure, atmosphere, pH) are set.

Component	Attributes	Description & Examples					
Target Compound	compound type	Chemical composition and description Material category: metals & alloys, ceramics, nano materials, polymers, semiconductors, etc.					
	synthesis method	Technique: sol-gel, hydrothermal, CVD, precipitation, electrodeposition, etc.					
	notes	Additional observations or variations					
	name	Chemical name (e.g., Nickel Nitrate, Deionized Water)					
Material	amount	Quantity used (numeric value)					
	unit	Mass (g, mg), Volume (mL, L), Molar (mol, mmol). Concentration (M, mM), etc.					
	vendor	Supplier information					
	purity	Grade specification (99%, ACS grade, etc.)					
Equipment	name	Instrument type (autoclave, tube furnace, magnetic stirrer)					
	vendor	Manufacturer (Thermo Fisher, Agilent, Bruker etc.)					
	settings	Operating parameters (500 rpm, heating rate 5°C/min)					
	temperature duration	Process temperature with units (°C, K, °F) Time period with units (h, min, s, days)					
	pressure	Applied pressure with units (atm, bar, Pa, torr)					
~	atmosphere	Gas environment (air, N ₂ , H ₂ , Ar, vacuum)					
Conditions	stirring	Boolean and speed (rpm)					
	pН	Solution acidity/basicity					
	step number	Sequential order in procedure					
	action	Primary operation: add, mix, heat, cool, reflux					
Process Step		age, filter, wash, dry, etc.					
1 Tocess Step	description	Detailed procedure text					
	materials	List of materials used in this step					
	equipment	List of equipment used in this step					
	conditions	Environmental parameters for this step					

Example 2: High Agreement (Material: Fluorapatite–Titania Nanocomposite) This example demonstrates excellent agreement between LLM and human evaluations, with perfect consensus across six of seven criteria and only a minor 0.5-point difference in Semantic Accuracy. The extraction successfully captured all key aspects of the mechano-chemical synthesis procedure, correctly identifying the starting materials (CaHPO₄, Ca(OH)₂, CaF₂, and TiO₂), process steps (mixing, ball milling, annealing), and reaction conditions. The LLM accurately extracted specific parameters such as the 20 wt% TiO₂ content, 600 rpm milling speed, and 700°C annealing temperature, while properly classifying the synthesis method as ball milling followed by thermal treatment.

Furthermore, the LLM only evaluates synthesis procedures that are extracted, and does not point out procedures that failed to extract.

Table 2: Comparing domain expert evaluations to LLM-as-a-judge. μ_{exp} , $\mu_{1/2,exp}$ and σ_{exp} refer to the mean, median and standard deviation for all six annotators and μ_{LLM} , $\mu_{1/2,LLM}$ and σ_{LLM} to the mean, median and standard deviation of the LLM (Gemini-2.0-flash), respectively.

Criterion	Spearman	p-value	Cohen	ICC(2,1)	ICC(3,1)	μ_{exp}	$\mu_{1/2,exp}$	σ_{exp}	μ_{LLM}	$\mu_{1/2,LLM}$	σ_{LLM}
Structural Completeness	0.4209	0.0004	0.2029	0.2286	0.2304	4.12	4.00	0.65	4.02	4.00	0.40
Material Extraction	0.7107	0.0002	0.5790	0.5996	0.5964	4.08	4.00	0.89	4.11	4.00	0.59
Process Steps	0.5547	0.0002	0.2867	0.2620	0.2626	4.15	4.25	0.82	4.27	4.25	0.55
Equipment Extraction	0.5842	0.0002	0.6287	0.6229	0.6325	4.05	4.50	1.19	3.80	4.00	1.18
Conditions Extraction	0.6201	0.0002	0.4747	0.4283	0.4565	4.27	4.00	0.70	4.01	4.00	0.68
Semantic Accuracy	0.5407	0.0002	0.3919	0.4170	0.4133	4.39	4.50	0.64	4.39	4.50	0.38
Format Compliance	0.2690	0.0350	0.1129	0.2141	0.2137	4.77	5.00	0.53	4.83	5.00	0.30
Overall	0.7195	0.0002	0.4407	0.5411	0.5399	4.25	4.30	0.52	4.20	4.25	0.42

Table 3: Evaluation scores for a low-agreement synthesis procedure extraction for Au-OLC from paper id 9a889c1a671fd3cae48285eaa95069d189d02fe3443.

Criterion	Human	LLM	Difference
Structural Completeness	2.0	4.0	2.0
Material Extraction	2.0	3.0	1.0
Process Steps	2.0	4.0	2.0
Equipment Extraction	5.0	4.0	1.0
Conditions Extraction	5.0	4.5	0.5
Semantic Accuracy	2.0	3.5	1.5
Format Compliance	4.0	5.0	1.0
Overall	3.1	4.0	0.9

A.2.3 Scaling LLM-as-a-judge across the dataset

Figure 8, Figure 9, Figure 11, Figure 10, Figure 12, Figure 13, Table 5 and Table 6 show the performance of LLM-as-a-judge across the dataset. For the sample on which we assess human–LLM agreement (n=66), we report Spearman rank correlations (ρ) between human and model scores, but compute their p-values using a permutation test (10,000 resamples, two-sided) rather than relying on the standard asymptotic approximation. This choice is motivated by the modest sample size and the bounded, quasi-ordinal nature of the scores, which induce many ties and can render asymptotic p-values anticonservative and unreliable. As the SciPy documentation recommends, "for small samples, consider performing a permutation test instead of relying on the asymptotic p-value," especially when ties and discrete data violate large-sample assumptions. The permutation procedure generates the exact finite-sample null distribution of ρ by permuting only one input (human scores) relative to the other while preserving marginal distributions. This approach provides valid inference under exchangeability, naturally handles ties, and ensures robust significance testing even with small, discrete datasets.

⁴https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.spearmanr.html

Table 4: Evaluation scores for a high-agreement synthesis procedure extraction for Fluorapatite—Titania Nanocomposite from paper id ccc7c5d70ae3ca3f9e975d0dc3b4d631586c1586.

Criterion	Human	LLM	Difference
Structural Completeness	4.0	4.0	0.0
Material Extraction	4.0	4.0	0.0
Process Steps	4.5	4.5	0.0
Equipment Extraction	4.0	4.0	0.0
Conditions Extraction	4.5	4.5	0.0
Semantic Accuracy	4.0	4.5	0.5
Format Compliance	5.0	5.0	0.0
Overall	4.4	4.3	0.1

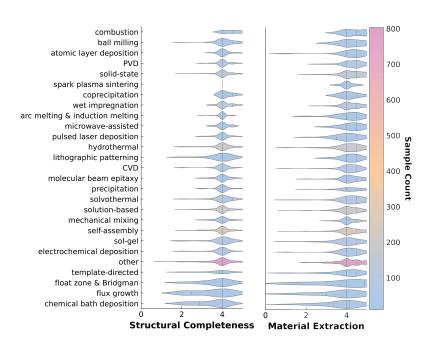


Figure 8: Distribution of LLM-judged overall extraction scores across different synthesis methods (structural completeness and material extraction score). See Table 5 for the full score overview. Each violin plot shows the probability density of the scores for a given synthesis type.

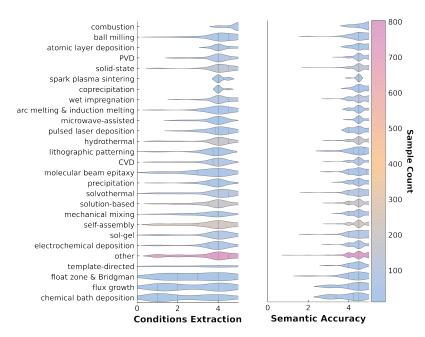


Figure 9: Distribution of LLM-judged overall extraction scores across different synthesis methods (condition extraction and semantic accuracy score). See Table 5 for the full score overview. Each violin plot shows the probability density of the scores for a given synthesis type.

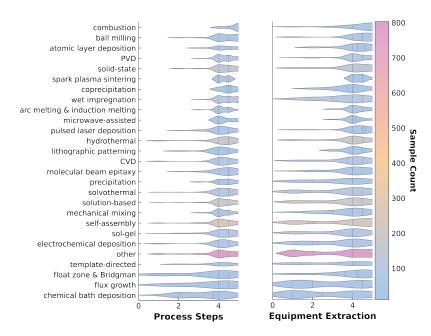


Figure 10: Distribution of LLM-judged overall extraction scores across different synthesis methods (process steps and equipment extraction score). See Table 5 for the full score overview. Each violin plot shows the probability density of the scores for a given synthesis type.

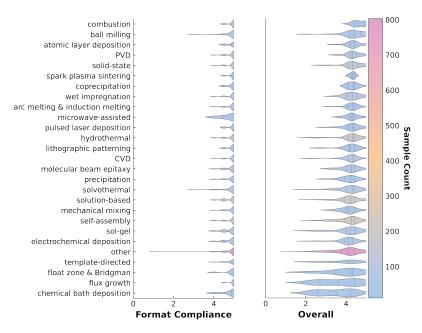


Figure 11: Distribution of LLM-judged overall extraction scores across different synthesis methods (format compliance and overall score). See Table 5 for the full score overview. Each violin plot shows the probability density of the scores for a given synthesis type.

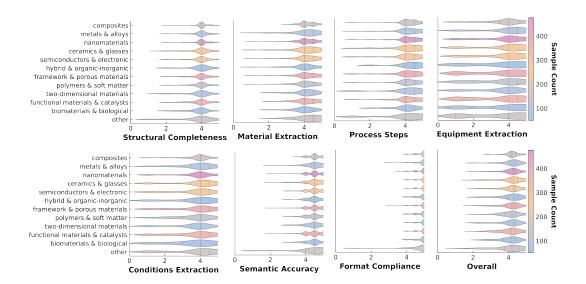


Figure 12: Distribution of LLM-judged overall extraction scores across different material classes. See Table 6 for a complete overview. Each violin plot shows the probability density of the scores for a given material category.

Table 5: Average LLM-judged extraction scores for the most frequent synthesis methods in the evaluated dataset subset (N=2483 procedures). Scores are reported as mean \pm standard deviation on a 1–5 scale. The Overall Score is the average of all seven evaluation criteria.

Synthesis	Structural	Material	Process	Equipment	Condition	Semantic	Format	Overall	Count
method	completeness	completeness	steps	extraction	extraction	accuracy	compliance	score	
other	3.85±0.73	4.14±0.65	4.00±0.99	3.22±1.55	3.47±1.31	4.42±0.57	4.83±0.43	3.99±0.70	803
self-assembly	3.94 ± 0.50	4.16 ± 0.58	4.13 ± 0.75	3.56 ± 1.53	3.56±1.16	4.49 ± 0.39	4.89 ± 0.26	4.10 ± 0.56	226
solution-based	4.01 ± 0.48	4.12 ± 0.57	4.23 ± 0.61	3.50 ± 1.32	3.71 ± 0.90	4.42 ± 0.40	4.84 ± 0.28	4.12 ± 0.51	180
hydrothermal	3.99 ± 0.57	4.09 ± 0.70	4.17 ± 0.86	3.88 ± 1.06	3.89 ± 0.93	4.47 ± 0.41	4.87 ± 0.27	4.20 ± 0.59	167
solid-state	4.09 ± 0.42	4.29 ± 0.56	4.29 ± 0.61	3.96 ± 1.10	4.13 ± 0.76	4.54 ± 0.41	4.92 ± 0.22	4.32 ± 0.44	134
wet impregnation	4.15 ± 0.37	4.23 ± 0.47	4.42 ± 0.46	3.49 ± 1.20	4.17 ± 0.60	4.53 ± 0.40	4.91 ± 0.23	4.28 ± 0.38	92
solvothermal	4.03 ± 0.52	4.21 ± 0.69	4.26 ± 0.62	3.47 ± 1.42	3.80 ± 1.11	4.47 ± 0.49	$4.84{\pm}0.37$	4.15 ± 0.57	89
CVD	3.96 ± 0.45	4.16 ± 0.55	4.18 ± 0.74	3.79 ± 1.11	3.71 ± 0.87	4.47 ± 0.34	4.92 ± 0.22	4.18 ± 0.46	79
PVD	4.06 ± 0.34	4.32 ± 0.49	4.34 ± 0.42	4.14 ± 0.78	3.92 ± 0.71	4.57 ± 0.34	4.90 ± 0.23	4.32 ± 0.33	77
sol-gel	3.94 ± 0.65	4.00 ± 0.71	4.20 ± 0.77	3.43 ± 1.30	3.76 ± 1.09	4.46 ± 0.53	4.81 ± 0.31	4.09 ± 0.64	70
electrochemical deposition	3.91 ± 0.49	4.14 ± 0.59	4.12 ± 0.82	3.26 ± 1.36	3.74 ± 0.92	4.41 ± 0.37	4.84 ± 0.29	4.06 ± 0.51	56
ball milling	4.17 ± 0.50	4.21 ± 0.56	4.38 ± 0.60	4.36 ± 0.93	4.07 ± 0.87	4.51 ± 0.52	4.88 ± 0.34	4.37 ± 0.52	54
precipitation	4.03 ± 0.38	4.20 ± 0.58	4.35 ± 0.47	3.36 ± 1.35	$3.82{\pm}0.68$	4.48 ± 0.38	4.89 ± 0.25	4.16 ± 0.44	47
mechanical mixing	4.01 ± 0.38	4.11 ± 0.36	4.21 ± 0.41	3.67 ± 1.07	3.46 ± 0.93	4.46 ± 0.37	4.90 ± 0.22	4.12 ± 0.39	47
molecular beam epitaxy	4.00 ± 0.37	4.21 ± 0.54	4.30 ± 0.63	3.96 ± 1.01	3.38 ± 1.10	4.46 ± 0.43	4.87 ± 0.27	4.17 ± 0.42	46
pulsed laser deposition	3.99 ± 0.35	4.01 ± 0.63	4.11 ± 0.61	4.26 ± 0.78	3.99 ± 0.67	4.46 ± 0.36	4.91 ± 0.19	4.25 ± 0.41	40
arc & induction melting	3.99 ± 0.22	4.07 ± 0.69	4.20 ± 0.36	4.22 ± 0.49	4.01 ± 0.75	4.47 ± 0.31	4.92 ± 0.22	4.27 ± 0.30	37
flux growth	3.45 ± 0.96	3.64 ± 1.14	3.59 ± 1.33	2.88 ± 1.63	2.74±1.47	4.24 ± 0.69	4.97 ± 0.13	3.64 ± 0.94	29
chemical bath deposition	3.39 ± 0.80	3.66 ± 0.95	3.29 ± 1.11	2.70 ± 1.65	2.41±1.47	4.02 ± 0.70	4.80 ± 0.37	3.47 ± 0.85	28
template-directed	3.75 ± 0.71	3.96 ± 0.85	3.94 ± 0.85	3.29 ± 1.55	3.27±1.41	4.33 ± 0.41	4.90 ± 0.25	3.92 ± 0.75	24
atomic layer deposition	4.10 ± 0.33	4.25 ± 0.77	4.27 ± 0.57	4.17 ± 0.83	4.17 ± 0.41	4.60 ± 0.36	4.88 ± 0.22	4.36 ± 0.34	24
combustion	4.46 ± 0.41	4.40 ± 0.51	4.75±0.39	4.42 ± 0.97	4.77±0.39	4.71 ± 0.39	4.94 ± 0.17	4.63 ± 0.28	24
float zone & Bridgman	3.73 ± 0.78	3.70 ± 1.32	3.52 ± 1.41	3.95 ± 1.25	2.91±1.41	4.32 ± 0.66	4.91 ± 0.29	3.87 ± 0.84	22
microwave-assisted	4.05 ± 0.28	$4.25{\pm}0.53$	4.32 ± 0.44	4.10 ± 0.45	3.80 ± 0.62	4.58 ± 0.24	4.72 ± 0.38	4.26 ± 0.29	20
lithographic patterning	3.87 ± 0.64	4.17 ± 0.49	4.03 ± 0.67	4.10 ± 0.57	3.67 ± 0.79	4.43 ± 0.53	4.93 ± 0.18	4.18 ± 0.45	15
coprecipitation	4.21 ± 0.33	4.12 ± 0.43	4.50 ± 0.37	3.83 ± 0.83	4.08 ± 0.19	4.62 ± 0.31	4.92 ± 0.19	4.32 ± 0.25	12
spark plasma sintering	4.00 ± 0.00	4.05 ± 0.27	4.23 ± 0.26	4.32 ± 0.34	4.14 ± 0.23	4.45 ± 0.15	4.95 ± 0.15	4.32 ± 0.12	11
mechanochemical	3.94 ± 0.88	3.94 ± 0.88	3.89 ± 1.34	4.11 ± 1.27	3.89 ± 1.11	4.44 ± 0.53	4.78 ± 0.36	4.14 ± 0.79	9
sonochemical	4.08 ± 0.20	4.25 ± 0.27	4.25 ± 0.27	4.00 ± 0.00	3.83 ± 0.26	4.50 ± 0.00	5.00 ± 0.00	4.28 ± 0.08	6
spray pyrolysis	4.33 ± 0.41	4.58 ± 0.38	4.50 ± 0.55	4.67 ± 0.41	4.25±0.76	4.75±0.27	5.00 ± 0.00	4.58 ± 0.33	6
electrospinning	4.00 ± 0.00	4.00 ± 0.00	4.38 ± 0.25	4.00 ± 0.00	4.12 ± 0.25	4.38 ± 0.25	5.00 ± 0.00	4.28 ± 0.13	4
ion implantation	3.83 ± 0.29	4.00 ± 0.00	4.50 ± 0.50	3.83 ± 1.61	3.67 ± 0.58	4.50 ± 0.50	5.00 ± 0.00	4.20 ± 0.53	3
liquid-phase epitaxy	4.00±nan	4.00±nan	$4.00\pm nan$	2.00±nan	4.00±nan	4.50±nan	5.00±nan	3.90±nan	1
incipient wetness impregnation	4.00±nan	4.00±nan	4.00±nan	4.00±nan	4.00±nan	4.50±nan	5.00±nan	4.20±nan	1
arc discharge	-	-	-	=	-	-	-	-	0

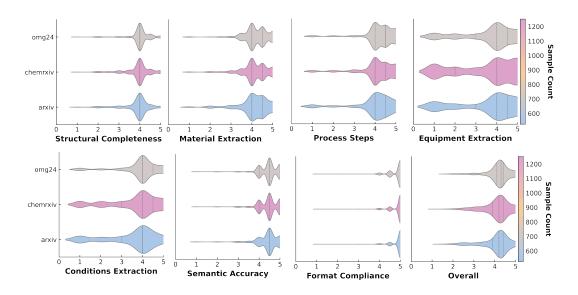


Figure 13: Distribution of LLM-judged overall extraction scores across different sources from LeMat-Synth. Each violin plot shows the probability density of the scores for a given synthesis type.

Table 6: Average LLM-judged extraction scores for the most frequent material types in the evaluated dataset subset (N=2483 procedures). Scores are reported as mean \pm standard deviation on a 1–5 scale. The Overall Score is the average of all seven evaluation criteria.

Material category	Structural completeness	Material completeness	Process steps	Equipment extraction	Condition extraction	Semantic accuracy	Format compliance	Overall score	Count
nanomaterials	4.01±0.47	4.14±0.57	4.21±0.68	3.65±1.24	3.76±0.97	4.48±0.41	4.85±0.29	4.16±0.51	476
framework & porous materials	$3.95{\pm}0.57$	4.15 ± 0.67	4.12 ± 0.84	$3.45{\pm}1.47$	3.63 ± 1.19	4.50 ± 0.43	4.88 ± 0.30	4.09 ± 0.61	385
functional materials & catalysts	3.93 ± 0.61	4.14 ± 0.63	4.12 ± 0.76	3.32 ± 1.51	3.52 ± 1.21	4.44 ± 0.45	4.88 ± 0.26	4.05 ± 0.61	351
ceramics & glasses	3.94 ± 0.65	4.10 ± 0.77	4.07 ± 0.95	3.80 ± 1.32	3.83 ± 1.15	4.43 ± 0.53	4.90 ± 0.26	4.15 ± 0.67	270
semiconductors & electronic	3.95 ± 0.57	4.16 ± 0.64	4.13 ± 0.84	3.64 ± 1.31	3.60 ± 1.16	4.48 ± 0.42	4.90 ± 0.23	4.13 ± 0.58	255
composites	4.06 ± 0.35	4.23 ± 0.41	4.27 ± 0.54	3.79 ± 0.97	3.90 ± 0.68	4.51 ± 0.34	4.86 ± 0.26	4.23 ± 0.35	154
other	3.75 ± 0.99	4.20 ± 0.69	$3.88{\pm}1.26$	3.26 ± 1.61	3.59 ± 1.36	4.33 ± 0.87	4.71 ± 0.76	3.96 ± 0.89	152
polymers & soft matter	3.96 ± 0.50	4.13 ± 0.61	4.20 ± 0.68	3.43 ± 1.38	3.62 ± 1.05	4.42 ± 0.42	4.84 ± 0.29	4.08 ± 0.54	132
metals & alloys	3.99 ± 0.45	4.11 ± 0.75	4.23 ± 0.66	3.87 ± 1.21	3.78 ± 1.01	4.48 ± 0.49	4.89 ± 0.31	4.19 ± 0.51	92
two-dimensional materials	3.88 ± 0.71	4.10 ± 0.63	4.05 ± 1.07	3.52 ± 1.30	3.56 ± 1.10	4.39 ± 0.49	4.90 ± 0.24	4.06 ± 0.66	89
biomaterials & biological	3.77 ± 0.60	4.01 ± 0.62	4.02 ± 0.69	3.48 ± 1.59	3.49 ± 1.25	4.40 ± 0.40	4.85 ± 0.30	4.00 ± 0.60	66
hybrid & organic-inorganic	3.93 ± 0.64	4.02 ± 0.70	4.25 ± 0.77	3.49 ± 1.50	3.71 ± 1.23	4.44 ± 0.38	4.86 ± 0.28	4.10 ± 0.65	51
energy & sustainability	4.31 ± 0.65	4.50 ± 0.46	4.50 ± 0.46	$4.12{\pm}1.33$	4.19 ± 0.65	4.69 ± 0.37	4.88 ± 0.23	$4.45{\pm}0.45$	8
emerging & quantum materials	4.50 ± 0.71	4.50 ± 0.71	4.75 ± 0.35	4.50 ± 0.71	4.50 ± 0.71	4.75 ± 0.35	4.75 ± 0.35	4.60 ± 0.57	2
liquid materials	-	-	-	-	-	-	-	-	0

A.3 Figure extraction

Segmenting large figures into sub-plots. To extract individual subplots from figures in research papers, we employ the DINO model [30] with zero-shot image segmentation. The prompt 'a plot' is used to guide the model in localizing subplot regions, with both text and box confidence thresholds set to 0.3. After initial detection, a post-processing step refines the bounding boxes to ensure complete coverage of each subplot, including axis labels and tick marks. To distinguish multi-panel figures from single-plot figures, we retain only bounding boxes that cover less than 50% of the total figure area; larger boxes are assumed to correspond to entire figures and are excluded. Empirical results indicate that this approach reliably identifies subplots across a variety of figure types.

Classifying plots with quantitative data. To classify segmented subplots and full-figure plots, we employ a ResNet-152 model [69], pretrained on ImageNet and fine-tuned on the DocFig dataset [31]. The dataset is split into 19,000 samples for training and 13,000 samples for testing. The model is trained with default hyperparameters for 20 epochs using the Adam optimizer with a learning rate of 1e-3. Our classification task focuses exclusively on the plot types "line chart", "bar plot" and "scatter plot" which are relevant for downstream information extraction; qualitative figures are excluded from further processing. The fine-tuned model achieves an F1-score of 88.03% on the test set, indicating strong performance in accurately identifying quantitative plots for subsequent analysis.

Extracting data with a vision LLM. To convert these numerical figures into a structured and interpretable format for further use, we explore the capabilities of advanced vision-language models to extract data from line plots, focusing on 2D coordinate retrieval. Inspired by [34], where multimodal models were used to extract and regenerate plots, we use Claude-Sonnet-4 (claude-sonnet-4-20250514) to extract 2D coordinates with their corresponding series names, as well as metadata fields like titles, axis labels, and units. The model is prompted to output a JSON object in a predefined schema, which is then parsed into a Pydantic object to ensure data consistency and structured integration into our data extraction pipeline.

A.3.1 Figure Extraction Evaluation

Manual annotations. For each series, the extracted coordinates are matched to the closest ground truth points using nearest-neighbor matching. This matching is performed in a normalized coordinate space, where both x and y axes are scaled to their respective ranges to ensure that errors are comparable across axes. The normalization scale is computed from the minimum and maximum values of the ground truth coordinates for each axis. We manually annotate 15 line charts from selected papers in catalysis [70, 71, 72, 73, 74]. For expanding the pipeline in the future, we plan to annotate larger samples from a more diverse array of plot types, e.g., scatter, bar, and box plots.

The evaluation is based on two error metrics:

- Root Mean Square Error (RMSE): which penalizes larger errors more heavily due to its quadratic nature.
- Mean Absolute Error (MAE): which treats all deviations linearly, providing a robust average error.

To compute the error metrics for a single series, we define the extracted points as:

$$\mathcal{P} = \{ (x_i, y_i) \mid i \in \{1, \dots, N\} \}$$
 (1)

and the ground truth points as:

$$\mathcal{G} = \{ (x_j^*, y_j^*) \mid j \in \{1, \dots, M\} \}$$
 (2)

Compute the normalization scales for each axis as:

$$S_x = \max_j x_j^* - \min_j x_j^*, \quad S_y = \max_j y_j^* - \min_j y_j^*$$
 (3)

For each extracted point (x_i, y_i) , we find the nearest ground truth point by computing the normalized Euclidean distance:

$$d_{i} = \min_{j} \sqrt{\left(\frac{x_{i} - x_{j}^{*}}{S_{x}}\right)^{2} + \left(\frac{y_{i} - y_{j}^{*}}{S_{y}}\right)^{2}} \tag{4}$$

The RMSE is then defined as:

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^{N} d_i^2}$$
 (5)

and the MAE as:

$$MAE = \frac{1}{N} \sum_{i=1}^{N} d_i$$
 (6)

A.4 Prompts

This section shows the system prompts employed and the full configurations used (incl. signatures and LLM configurations) to extract the data presented in this work.

Filtering papers

```
Prompt
Analyze the following text and answer the questions in JSON format:
{chunk}
Questions:
1. Does it contain a material synthesis recipe?
    (Answer with true or false)
2. If yes, what is the material name?
    (Answer with the material name or "N/A" if no recipe)
3. If yes, which category of materials does it belong to?
    (Answer with the specific material type or "N/A" if no recipe)
List of material categories:
Metals, Ceramics, Semiconductors, Superconductors, Composites,
Biomaterials, Nanomaterials, Polymers, Magnetic, Textiles, Chemicals, Other
Format your response as a JSON object with the following structure:
"contains_recipe": true/false,
"material name": "material name or N/A",
"material_category": "material category or N/A"
}}
```

Material extraction

```
Prompt
You are a helpful assistant that extracts ONLY the final synthesized materials
\hookrightarrow from scientific papers.
Your task is to identify ONLY the materials that are the final products of
\hookrightarrow synthesis procedures described in the paper.
IMPORTANT GUIDELINES:
- ONLY include materials that are the final synthesized products
- DO NOT include starting materials, precursors, supports, gases, solvents, or

→ other chemicals used in synthesis

- DO NOT include materials that are just mentioned or characterized but not
\hookrightarrow synthesized
- Focus on the main target materials that are actually synthesized
EXAMPLES OF WHAT TO INCLUDE:
- "Ni/Al203" (if Ni is deposited on Al203)
- "Ir/SiO2" (if Ir is supported on SiO2)
- "LiFePO4 nanoparticles" (if LiFePO4 is synthesized)
- "Co-doped LiFePO4" (if this specific material is synthesized)
EXAMPLES OF WHAT TO EXCLUDE:
- "Ni", "Ir", "Ru" (if these are just precursors)
- "H-ZSM-5", "Al203", "Si02" (if these are just supports)
- "Ammonia", "Argon", "Hydrogen" (gases)
- "Deionized water" (solvents)
- "Ammonium hydroxide" (reagents)
Return a simple comma-separated list of ONLY the final synthesized materials.
If no materials are synthesized in the paper, return "No materials
\hookrightarrow synthesized".
Keep the output simple and clean - just the final synthesized material names
\hookrightarrow separated by commas.
```

Configuration (YAML)

```
architecture:
  _target_: llm_synthesis.transformers.material_extraction.dspy_extraction.Dsp |
  \rightarrow yTextExtractor
  signature:
    _target_: llm_synthesis.transformers.material_extraction.dspy_extraction.m |
    \rightarrow ake_dspy_text_extractor_signature
    signature_name: "TextToMaterials"
    instructions: "Extract ONLY the final synthesized materials from the
    \hookrightarrow publication text."
    input_description: "The publication text to extract the final synthesized
    \hookrightarrow materials from."
    output_name: "materials"
    output_description: "The final synthesized materials as a comma-separated
    _target_: llm_synthesis.utils.dspy_utils.get_llm_from_name
    llm_name: "gemini-2.0-flash"
    model_kwargs:
     temperature: 0.0
    system_prompt:
```

```
_target_: llm_synthesis.utils.read_prompt_str_from_txt
prompt_path: "examples/system_prompts/material_extraction/default.txt"
```

Synthesis extraction

Prompt You are a helpful assistant that extracts the structured synthesis for a \rightarrow specific material from the paper text. Focus ONLY on the synthesis procedure for the specified material. Search \hookrightarrow through the entire paper text to find the synthesis procedure that describes how this specific material is made. IMPORTANT: You must output ONLY a valid JSON object with a → "structured_synthesis" field. Do not include any reasoning, explanations, \hookrightarrow or markdown formatting. If you cannot find a synthesis procedure for the specified material, return a \hookrightarrow minimal structure with the material name and an empty synthesis. The JSON output must follow this exact structure: "structured_synthesis": { "target_compound": "string (required) - should match the specified material "target_compound_type": "string (required) - choose from: 'metals & \hookrightarrow alloys', 'ceramics & glasses', 'polymers & soft matter', 'composites', \hookrightarrow 'semiconductors & electronic', 'nanomaterials', 'two-dimensional \hookrightarrow materials', 'framework & porous materials', 'biomaterials & \hookrightarrow biological', 'liquid materials', 'hybrid & organic-inorganic', $\,\hookrightarrow\,$ 'functional materials', 'energy & sustainability', 'smart & responsive \hookrightarrow materials', 'emerging & quantum materials', 'other'", "synthesis_method": "string (required) - choose from: 'PVD', 'CVD', 'arc discharge', 'ball milling', 'spray pyrolysis', 'electrospinning', 'sol-gel', 'hydrothermal', 'solvothermal', 'precipitation', \hookrightarrow coprecipitation', 'combustion', 'microwave-assisted', 'sonochemical', 'template-directed', 'solid-state', 'flux growth', 'float zone & \hookrightarrow Bridgman', 'arc melting & induction melting', 'spark plasma sintering', 'electrochemical deposition', 'chemical bath deposition', 'liquid-phase → epitaxy', 'self-assembly', 'atomic layer deposition', 'molecular beam → epitaxy', 'pulsed laser deposition', 'ion implantation', 'lithographic → patterning', 'wet impregnation', 'incipient wetness impregnation', 'mechanical mixing', 'other' "starting_materials": [{"name": "string", "amount": "number or null", → "unit": "string or null", "purity": "string or null", "vendor": "string \hookrightarrow or null"}], "steps": [{"step_number": "integer", "action": "string", "description": → "string or null", "materials": [{"name": "string", "amount": "number or → null", "unit": "string or null", "purity": "string or null", "vendor": → "string or null"}], "equipment": [{"name": "string", → "instrument_vendor": "string or null", "settings": "string or null"}], → "conditions": {"temperature": "number or null", "temp_unit": "string or → null", "duration": "number or null", "time_unit": "string or null", \hookrightarrow "pressure": "number or null", "pressure_unit": "string or null", \rightarrow "atmosphere": "string or null", "stirring": "boolean or null", $\ \, \rightarrow \ \, \text{"stirring_speed": "number or null", "ph": "number or null"}}],$ "equipment": [{"name": "string", "instrument_vendor": "string or null", → "settings": "string or null"}], "notes": "string or null"

```
}
}
Do not include any text before or after the JSON object. Output only the JSON.
```

```
Configuration (YAML)
architecture:
  _target_: llm_synthesis.transformers.synthesis_extraction.dspy_synthesis_ext |
  \hookrightarrow raction.DspySynthesisExtractor
  signature:
    _target_: llm_synthesis.transformers.synthesis_extraction.dspy_synthesis_e |

→ xtraction.make_dspy_synthesis_extractor_signature

    signature_name: "SynthesisSignature"
    instructions: "Extract the structured synthesis for a specific material
    \hookrightarrow from the paper text."
    paper_text_description: "The complete paper text to search for the

→ material's synthesis procedure."

    material_name_description: "The name of the specific material to extract
    \hookrightarrow synthesis for."
    output_name: "structured_synthesis"
    output_description: "The extracted structured synthesis for the specific

→ material."

 lm:
    _target_: llm_synthesis.utils.dspy_utils.get_llm_from_name
    llm_name: "gemini-2.0-flash"
    model_kwargs:
      temperature: 0.0
      max_tokens: 8000
      max_retries: 3
    system_prompt:
      _target_: llm_synthesis.utils.read_prompt_str_from_txt
      prompt_path: "examples/system_prompts/synthesis_extraction/default.txt"
```

Figure extraction

For figure extraction, we do not provide a separate DSPy configuration. Unlike material and synthesis extraction (which are wrapped with DSPy signatures and explicit input/output schemas), the figure extraction pipeline directly leverages the system prompt together with a Claude API client. In this setup, the model is invoked with the raw prompt and image data, and the parsing into structured objects (ExtractedLinePlotData) happens entirely within the custom transformer implementation. Because no DSPy signature or schema mediation is involved, there is no corresponding YAML configuration block to display. Instead, the logic is captured in the prompt (shown below) and the Python implementation excerpted below.

```
Prompt

LINE_CHART_PROMPT = """
You will be provided with a line chart. The chart may not be chunked very well, so you may need to read only the plot in the center of the image.
In the chart, there will be several lines representing different data series.

1. Identify the different lines by their colors and labels.
2. For each line, extract the coordinates of the points that make up the line.
   Do not include any points that are not part of the line.
3. If the chart has metadata such as a title, x-axis label, y-axis labels, or units, extract that information as well.
   Keep the scientific terms in Markdown format.
4. Output the data in the specified format:

Name_of_Line_1: [[x1, y1], [x2, y2], ...]
```

```
title:
  x_axis_label:
  x_axis_unit:
  y_left_axis_label:
  y_left_axis_unit:

Do not output any other text, just the data in the format above.
"""
```

```
Implementation excerpt (Python)
class ClaudeLinePlotDataExtractor(LinePlotDataExtractorInterface):
   def __init__(self, model_name: str,
                prompt: str = resources.LINE_CHART_PROMPT,
                max_tokens: int = 1024.
                temperature: float = 0.0):
       super().__init__()
       self.claude_client = ClaudeAPIClient(model_name)
       self.prompt = prompt
       self.max_tokens = max_tokens
       self.temperature = temperature
    def forward(self, input: FigureInfoWithPaper) -> ExtractedLinePlotData:
       figure_base64 = input.base64_data
       self.claude_client.reset_cost()
       claude_response_obj = self.claude_client.vision_model_api_call(
           figure_base64=figure_base64,
           prompt=self.prompt,
           max_tokens=self.max_tokens,
           temperature=self.temperature,
       return self._parse_into_pydantic(claude_response_obj)
    def _parse_into_pydantic(self, response: str) -> ExtractedLinePlotData:
         ""Parse text into Pydantic object with regex pattern matching""
```

Synthesis evaluation

In this case, the evaluation logic is fully captured within the DSPy configuration itself, so we do not provide a standalone prompt block. Both the task instructions and the system prompt are directly embedded inside the configuration file rather than stored separately. The complete configuration is shown below:

```
Configuration (YAML)
architecture:
  _target_: llm_synthesis.metrics.judge.general_synthesis_judge.DspyGeneralSyn |
  \hookrightarrow thesisJudge
 signature:
    _target_: llm_synthesis.metrics.judge.general_synthesis_judge.make_general_
    \ \hookrightarrow \ \texttt{\_synthesis\_judge\_signature}
    signature_name: "GeneralSynthesisJudgeSignature"
    instructions: >
      You are an expert materials scientist and data extraction specialist with
      \hookrightarrow extensive experience in:
        - Synthesis procedure analysis and documentation
        - Structured data extraction from scientific literature
        - Materials science ontology design and terminology standardization
        - Quality assessment of automated scientific information extraction
        \hookrightarrow systems
      Evaluate how well the GeneralSynthesisOntology extraction captures
      \hookrightarrow synthesis information from
      the provided source text.
```

```
IMPORTANT: Do NOT penalize the extraction system for failing to include
  \hookrightarrow information that is
  not present in the original paper. Missing elements should only be
  \hookrightarrow considered errors if they
  were clearly stated in the source but were not extracted. If an element
  \hookrightarrow is absent in both the
  source and the extraction, and is correctly left blank or omitted, this
  \hookrightarrow should be considered
  correct and scored highly.
  ASSESSMENT FOCUS:
    - Completeness: All synthesis components present in the source are
    \hookrightarrow captured
    - Accuracy: Correct values, units, and classifications based on the
    - Structure: Proper organization and logical sequencing of elements
    - Semantic Preservation: Scientific meaning and intent faithfully

→ maintained

    - Schema Compliance: Conforms to the expected ontology format and data
    \hookrightarrow types
  EVALUATION CRITERIA (Score 1-5 for each):
    1. Structural Completeness - Extraction of all relevant synthesis

→ components from the source (materials, steps, equipment,
    \hookrightarrow conditions)
    2. Material Extraction - Correct names, quantities, units, purities as
    \hookrightarrow specified in the paper
    3. Process Steps - Accurate step order and correct action
    \hookrightarrow \quad \textbf{classification} \quad
    4. Equipment Extraction - Proper identification of all equipment
    \hookrightarrow explicitly mentioned
    5. Conditions Extraction - Accurate recording of parameters such as
    \rightarrow temperature, time, atmosphere, pressure, etc.
    6. Semantic Accuracy - Faithful preservation of scientific meaning
    \hookrightarrow without misinterpretation
    7. Format Compliance - Adherence to ontology schema, data types, and
    For each criterion:
    - Assign a score between 1 and 5
    - Provide detailed technical reasoning for the assigned score
    - Offer specific, constructive recommendations for improvement, if
    \hookrightarrow applicable
_target_: llm_synthesis.utils.dspy_utils.get_llm_from_name
llm_name: "gemini-2.0-flash"
model_kwargs:
  temperature: 0.1
  max_tokens: 4096
system_prompt: >
  You are a senior materials scientist and data extraction expert with deep
  \hookrightarrow expertise in:
    - Inorganic and organic synthesis methodologies
    - Laboratory instrumentation and experimental workflows
    - Chemical nomenclature, stoichiometry, and unit conventions
    - Optimization of synthesis conditions and reaction parameters
    - Structured data modeling and materials science ontology design
    - Evaluation methodologies for automated information extraction systems
  Your assessments should reflect best practices in synthesis reporting and
```

28

 \hookrightarrow uphold the highest

standards of scientific accuracy, reproducibility, and structured data $\mbox{\ \hookrightarrow\ }$ quality.

When evaluating extracted synthesis data:

- Rely on your domain expertise to assess technical correctness,
- $\,\,\hookrightarrow\,\,$ semantic fidelity, and structural organization
- Emphasize clarity, precision, and alignment with real-world
- \hookrightarrow experimental protocols
- Consider the intended schema and use context to assess compliance and
- $\hookrightarrow \quad \texttt{completeness}$
- Do not penalize the extraction system for omitting elements that were $% \left(1\right) =\left(1\right) \left(1\right) \left($
- $\,\,\hookrightarrow\,\,$ not explicitly present in the source text

Your evaluation should be technically rigorous, yet fair, grounded in \hookrightarrow both materials science principles and data extraction best practices.

enable_reasoning_traces: true
confidence_threshold: 0.7