# Representation of Inorganic Synthesis Reactions and Prediction: Graphical Framework and Datasets

**Samuel Andrello**
Columbia University
New York, NY 10027

**Daniel Alabi**
University of Illinois at Urbana-Champaign
Urbana, IL 61801

**Simon J. L. Billinge**
Columbia University
New York, NY 10027

## Abstract

While machine learning has enabled the rapid prediction of inorganic materials with novel properties, the challenge of determining how to synthesize these materials remains largely unsolved. Previous work has largely focused on predicting precursors or reaction conditions, but only rarely on full synthesis pathways. We introduce the ActionGraph, a directed acyclic graph framework that encodes both the chemical and procedural structure, in terms of synthesis operations, of inorganic synthesis reactions. Using 13,017 text-mined solid-state synthesis reactions from the Materials Project, we show that incorporating PCA-reduced ActionGraph adjacency matrices into a $k$-nearest neighbors retrieval model significantly improves synthesis pathway prediction. While the ActionGraph framework only results in a 1.34% and 2.76% increase in precursor and operation F1 scores (average over varying numbers of PCA components) respectively, the operation length matching accuracy rises 3.4 times (from 15.8% to 53.3%). We observe an interesting trade-off where precursor prediction performance peaks at 10-11 PCA components while operation prediction continues improving up to 30 components. This suggests composition information dominates precursor selection while structural information is critical for operation sequencing. Overall, the ActionGraph framework demonstrates strong potential, and with further adoption, its full range of benefits can be effectively realized.

## 1   Introduction

Inorganic solid-state synthesis is notoriously challenging to predict and design due to the lack of a general theory describing how phases form and transform under given conditions. In particular, there is no comprehensive theoretical framework dictating how precursor phases evolve into products during heating and other processing steps [ESHH+24]. As a result, synthesis routes are often developed by trial-and-error or past experience. Moreover, real-world data on multi-step syntheses with intermediate phases is scarce[AW24]. Most reported synthesis recipes focus only on the starting materials and final products, omitting the transient *intermediate* phases that may appear along the reaction pathway. Even with recent efforts to compile text-mined datasets of inorganic synthesis recipes [KHH+19], the available data remain limited in capturing the full complexity of synthesis processes (especially intermediate steps). This combination of a knowledge gap and limited data motivates new computational approaches to guide inorganic materials synthesis.

In this work, we address the problem of identifying which *precursor* materials and intermediate steps can lead to the desired *products* in inorganic synthesis, given only the target products as input. We propose a novel graph-based framework, termed the ActionGraph, to represent an entire inorganic synthesis pathway. Using this representation, we employ $k$-nearest neighbors models [CH67] with dimensionality-reduced graph adjacency matrices to predict the reaction steps-specifically, to predict likely precursors and intermediates that yield the target product Our ultimate goal is to develop a predictive model that can suggest how to synthesize new materials by learning from known examples. Such a model could significantly reduce the trial-and-error in materials discovery and lay

the groundwork for automated synthesis planning. Indeed, data-driven methods are beginning to play a role in rationalizing and predicting synthesis outcomes [HBH$^+$22], and an accurate precursor-prediction model could eventually be integrated into robotic synthesis laboratories for closed-loop experimentation.



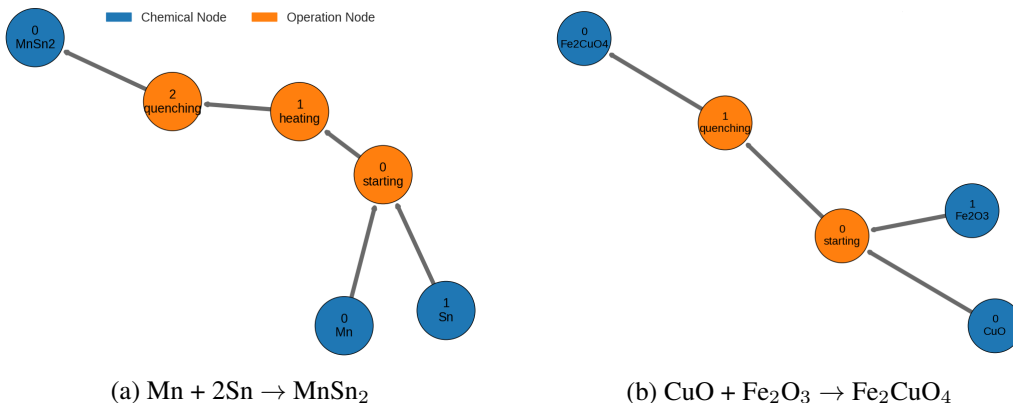(a) Mn + 2Sn $\rightarrow$ MnSn$_2$      (b) CuO + Fe$_2$O$_3$ $\rightarrow$ Fe$_2$CuO$_4$

Figure 1: Two synthesis reactions represented in ActionGraph form. Nodes are labeled with their corresponding indices. Chemical nodes placed prior to operations have a separate set of indices than chemical nodes placed after operations.

## 1.1 Inorganic Synthesis: Why is the problem difficult?

Inorganic synthesis presents unique challenges compared to organic synthesis, making it particularly difficult to predict and automate. While organic synthesis benefits from established retrosynthesis frameworks [WSQ$^+$25] that allow for systematic planning of reaction pathways, inorganic synthesis lacks analogous theoretical foundations. This absence of a comprehensive theory means that predicting how precursors transform into products during solid-state reactions remains largely empirical and experience-driven.

Several factors contribute to this complexity. First, inorganic syntheses often involve high-temperature solid-state reactions where atomic diffusion and phase transformations occur through mechanisms that are difficult to observe directly. Second, reaction pathways frequently involve transient intermediate phases that may not be reported in the literature or captured in databases. Third, the outcome of an inorganic synthesis depends heavily on processing conditions such as temperature profiles, atmospheres, and mechanical treatments, creating a vast parameter space that is challenging to navigate systematically.

This complexity has significant implications for materials innovation. For instance, the development of novel battery materials [SZH$^+$21] often involves extensive trial-and-error experimentation to discover viable synthesis routes, substantially slowing the transition from computational material design to physical realization. The field has consequently seen growing interest in "closed-loop" synthesis approaches [SZH$^+$21], where automated experimentation platforms could iteratively refine synthesis protocols with minimal human intervention.

Unlike organic chemistry, where retrosynthesis provides a structured method to solve the synthesis inverse problem (finding precursors and conditions that yield a target molecule), inorganic materials science lacks such systematic frameworks. This gap has motivated recent efforts to develop data-driven approaches that can learn from existing synthesis recipes to predict viable pathways for novel materials [KNG$^+$24, MDP21]. The ActionGraph framework we present aims to address this challenge by providing a structured representation of inorganic synthesis that captures both the chemical constituents and the procedural aspects of materials preparation.

Our contributions are as follows:

- The ActionGraph framework which in this paper is being used to represent inorganic synthesis reactions, but could be extended to other processes or any set of actions with a process-like structure.

- A dataset of 13,017 ActionGraph graphs transformed from open-source Materials Project data, each representing a single solid-state inorganic synthesis reaction.

Next, we discuss some related work.

## 2 Related Work

### 2.1 Machine Learning for Sciences

Our work on ActionGraph sits in the domain of research to further enable scientific discovery via the use of machine learning. Recently, the application of machine learning (ML) to scientific domains has grown rapidly. Such applications have led to new and more efficient simulations and has facilitated the discovery of patterns in complex, high-dimensional data. For example, in the physics community, ML has been employed to uncover hidden structures in quantum systems [CT17, MST21], model phase transitions [vNLH17], and approximate solutions to differential equations [RPK19]. In chemistry and materials science, deep learning techniques have achieved impressive performance in molecular property prediction [GSR$^+$17], drug discovery [ZIA$^+$19], and the generation of novel molecules [GBWD$^+$18]. However, none of these works develop frameworks nor datasets to address the important task of synthesis of materials. Our work aims to bridge this gap.

In terms of ML architectures, Graph neural networks (GNNs) or graph-based architectures have emerged as a key tool for modeling structured scientific data. They have been particularly effective in learning interatomic potentials [SSK$^+$18], modeling physical interactions [BHB$^+$18, BP24], and predicting protein structures [JEP$^+$21]. Meanwhile, generative models, including variational autoencoders and generative adversarial networks, have facilitated the inverse design of materials [SLAG18] and molecular synthesis [PIT18]. Recent work also explores how ML can augment simulation, such as learning surrogate models for climate prediction [RDK$^+$19] or accelerating Monte Carlo methods in statistical physics [LW17]. Furthermore, Reinforcement learning has been applied to optimize experimental design and automate laboratory processes [MPM$^+$20].

Despite these advances, challenges remain in developing ML methods that incorporate physical constraints, provide uncertainty estimates, and generalize beyond observed regimes [BL07, JGD$^+$08, NNZTB25]. Our work on ActionGraph presents a novel graph-based framework to capture the problem of inorganic synthesis and to incorporate physical constraints. We also provide datasets to evaluate the efficacy of predicting synthesis pathways.

### 2.2 Synthesis Prediction

Synthesis prediction has evolved along different trajectories in organic and inorganic chemistry, reflecting the distinct challenges in each domain. In organic chemistry, retrosynthetic analysis-working backward from target molecules to starting materials through known reaction types-has been formalized for decades and recently enhanced through machine learning approaches [WSQ$^+$25]. These methods benefit from large, structured datasets of reactions and well-established reaction mechanisms.

In contrast, inorganic synthesis prediction has emerged more recently and faces unique challenges. McDermott et al. [MDP21] introduced a graph-based framework to predict solid-state reaction pathways, focusing on thermodynamic driving forces rather than synthesis procedures. Their approach, while pioneering, primarily addresses reaction energetics rather than practical synthesis protocols.

He et al. [HHB$^+$23] developed a machine learning approach specifically for precursor recommendation, achieving significant accuracy in predicting starting materials for inorganic syntheses. However, their work focuses primarily on precursor selection rather than complete synthesis pathways including operations and conditions.

More recently, Kim et al. [KNG$^+$24] proposed an elementwise template formulation for predicting inorganic synthesis recipes, while another study by Kim et al. [KJS24] explored using large language models for inorganic synthesis predictions. Recently, Noh et al. [NLNP25] predicted reaction components using an attention mechanism. These approaches make important strides but focus primarily on either precursors or reaction conditions rather than the complete synthesis workflow.

Karpovich et al. [KJVO21] addressed reaction condition prediction using generative machine learning, emphasizing the parameter space of synthesis conditions rather than the structural aspects of synthesis pathways.

What distinguishes our work is the ActionGraph framework's holistic representation of inorganic synthesis as a directed process flow capturing both materials and operations. Unlike previous approaches that focus on either precursors or conditions in isolation, our methodology encodes the entire synthesis pathway as a graph structure, enabling prediction of complete synthesis routes rather than individual components.

Now, we provide preliminary definitions for our results.

## 3 Preliminaries

**ActionGraph definition.** We model an inorganic synthesis process as an ActionGraph, which is a directed, acyclic graph capturing the relationships between starting materials, intermediate compounds, and final products through the sequence of synthesis actions. Formally, an ActionGraph $G = (V, E)$ consists of nodes $V$ representing chemical compounds (phases) and operations (synthesis steps). The directed edges $E$ represent "action flow," or the flow of phases through the synthesis process. The graph is acyclic, reflecting the temporal order of synthesis steps.

**Precursors and products.** In an ActionGraph, we distinguish two roles for compound nodes. *Precursors* are the initial reactant materials provided at the start of the synthesis (these correspond to source nodes with no incoming edges in the graph). *Products* are the final target materials of the synthesis, which only have incoming edges. For example, if a recipe starts with powders of $A$ and $B$ (precursors), mixes and heats them to form an intermediate compound $C$, and upon further heating yields final product $D$, we would represent this as precursor nodes $A$, $B$ feeding into intermediate node $C$ (via a "heating" edge), which in turn connects to product node $D$ via another action edge.

**Operations.** In order to capture the nature of synthesis, an ActionGraph also contains operation nodes which represent the actual synthesis steps that transform precursors into products. They are topologically ordered after precursor nodes, but before product nodes. These nodes can represent operations such as heating, mixing, and milling that reflect a high-level view of a syntehsis process (i.e. what is actually being done to the precursors and reaction intermediates).

*Formal definition.* Define an ActionGraph as $G = (V, E)$ with

$V = C \cup O$ where $C$ is the set of all chemical nodes and $O$ is the set of all operation nodes,

$E \subset V \times V$ is the set of directed edges representing the flow of the synthesis process.

Define a chemical node $C_i = \{\text{element}_i : \text{subscript}_i\}_{i=1}^m$ where $\text{subscript}_i$ is the stoichiometric coefficient for $\text{element}_i$.

Define $C_{in}$ and $C_{out}$, the sets of all precursors and products, respectively. Then $C \equiv C_{in} \cup C_{out}$.

Define an operation node as

$O_i = \left(\text{type}_i, \{\text{condition}_j : \text{value}_j\}_{j=1}^m, \text{metadata}_j\right)$

where $\text{type}_i \in \{\text{StartingSynthesis}, \text{MixingOperation}, \text{ShapingOperation},$

$\text{DryingOperation}, \text{HeatingOperation}, \text{QuenchingOperation}\}$.

*Structural constraints.* We also consider some constraints on the structure of an ActionGraph that are required for its validity:

1. $\forall c \in C_{in}, \text{indegree}(c) = 0 \iff$ precursor nodes have no incoming edges.
2. $\forall c \in C_{out}, \text{outdegree}(c) = 0 \iff$ product nodes have no outgoing edges.
3. $\forall o \in O, \text{indegree}(o) \geq 1$ and $\text{outdegree}(o) \geq 1 \iff$ operation nodes have at least one input and one output.
4. $G$ is acyclic.

**Synthesis workflow and modeling challenges.** A typical solid-state synthesis workflow involves several sequential actions: precursors are first combined (mixed, ground together, sometimes pressed

or shaped into pellets), then subjected to thermal treatments such as drying (to remove solvents or binders) and one or more heating steps (calcination, sintering) at high temperatures, and possibly cooling or quenching steps at the end. These processing actions can lead to the formation of new crystalline phases as intermediates before the final product phase is obtained. Modeling this process computationally is challenging because the outcome of each step can depend non-linearly on the combination of precursors, the specific conditions (temperature, time), and the sequence of actions. The ActionGraph framework provides a structured way to encode all these aspects (materials and actions) into a single representation. However, learning to predict the sequence from precursors to product requires capturing potentially complex relationships in the graph (e.g., how multiple precursor compounds might combine to form a particular intermediate). The lack of abundant, richly annotated data detailing intermediate information such as crystal structure and reaction conditions further complicates direct modeling. In addition, it is desirable to capture the high-level synthesis process from this data, as the goal is to predict synthesis pathways that are directly usable in a laboratory setting. These challenges necessitate an adaptable representation for synthesis (hence our graph-based approach) along with benchmarks on the available data.

## 4 Methodology

### 4.1 Initial Dataset Construction

**Dataset and preprocessing.** We leverage a recently published dataset of inorganic materials synthesis recipes text-mined from the literature [KHH$^+$19]. This dataset (from the Materials Project repository) contains thousands of solid-state reaction examples, each including a list of precursors, details of synthesis operations (e.g. "heat at 800°C for 10 hours"), and the resulting product (and occasionally intermediate phases if reported by the authors). The initial, unfiltered set numbered 30,850 JSON files, each representing a single solid-state synthesis reaction. To focus on well-defined reactions, we filtered out any entries with variable or non-stoichiometric compositions (e.g. recipes that involve formula variables or unspecified proportions), as these are difficult to interpret in a standardized way. Furthermore, erroneous reactions, or those that contain invalid or empty chemical formula, were also removed. The final dataset numbered 13,017 reactions.

After cleaning, a second dataset was constructed by transforming each synthesis reaction into our ActionGraph representation. In this conversion, the precursors are set as the initial nodes, any operations as nodes following this, and the products as the final nodes. Edges are added to connect these nodes in the order of the described steps. For instance, if a paper describes that mixing compounds A and B followed by heating yields intermediate C, which upon further heating yields product D, we create edges $A \rightarrow C$, $B \rightarrow C$ (for the mixing+heating step forming $C$) and $C \rightarrow D$ (for the subsequent step yielding $D$). Each edge is annotated with the type of action (mixing, heating, etc.), though in our current model we focus primarily on the graph connectivity and not the exact conditions due to the lack of sufficient data. This second dataset was of the same length as the untransformed, filtered dataset: 13,017 serialized ActionGraph objects.

Data sharing in materials science is not encouraged due to copyright or legal reasons [AGM$^+$25]. As a result, we can rely on publicly-available datasets instead.

**Approach overview.**

In order to demonstrate the utility of the ActionGraph framework, we evaluate its performance through comparative experiments on two distinct datasets. The first dataset is derived from the filtered Materials Project (MP) data, which provides material compositions and associated properties in a conventional tabular format. Using this dataset, we construct a baseline $k$-nearest-neighbors ($k$-NN) model with $k = 1$. This model serves as a simple, interpretable benchmark that relies purely on standard feature vectors extracted directly from the material property data, without any graph-based augmentation.

Next, to assess the added value introduced by the ActionGraph framework, we train a second $k$-NN model using data processed through ActionGraph. This second dataset encodes domain knowledge and inter-entity dependencies by representing materials and reactions as nodes and edges in a graph. To incorporate this graph structure into the learning pipeline, we transform each data instance into an enriched feature vector that includes not only conventional material attributes, but also information derived from the graph topology. Specifically, we integrate structural features obtained from the

adjacency matrix of the ActionGraph, which encodes the connectivity patterns among nodes (e.g., reactions and precursors). This embedding captures relational and contextual cues, such as shared precursors, co-occurring reactions, or common synthesis pathways, that are otherwise absent in flat representations. By augmenting the feature space in this manner, we enable the second $k$-NN model to exploit both local material properties and higher-order structural correlations, thus providing a more holistic assessment of material similarities and model performance within the ActionGraph paradigm.

**Featurization.** Each synthesis reaction is featurized to serve as input to a knn. In our preliminary implementation, we use a simple composition-based feature vector: essentially a representation of the chemical formula of the compound. This can be, for example, an element fraction vector indicating the proportion of each element in the compound's formula, or a one-hot encoding of the compound identity if it is recognized from a known list. The intuition is to provide the model with information about the chemical makeup of each node (since knowing the elements present in a precursor vs in the product might help the model learn which precursors could lead to that product). this composition vector is then concatenated with another feature vector containing the averages of each of atomic mass, atomic radius, melting point, Pauling electronegativity, electron affinity, and ionization energy for each chemical. However, this straightforward featurization may not capture more subtle aspects of synthesis and is limited by the lack of structural characterization information of chemicals in the dataset. Exploring more sophisticated featurization (incorporating, e.g., known chemical descriptors or features learned from materials graphs [ESHH⁺24]) is an area of ongoing work.
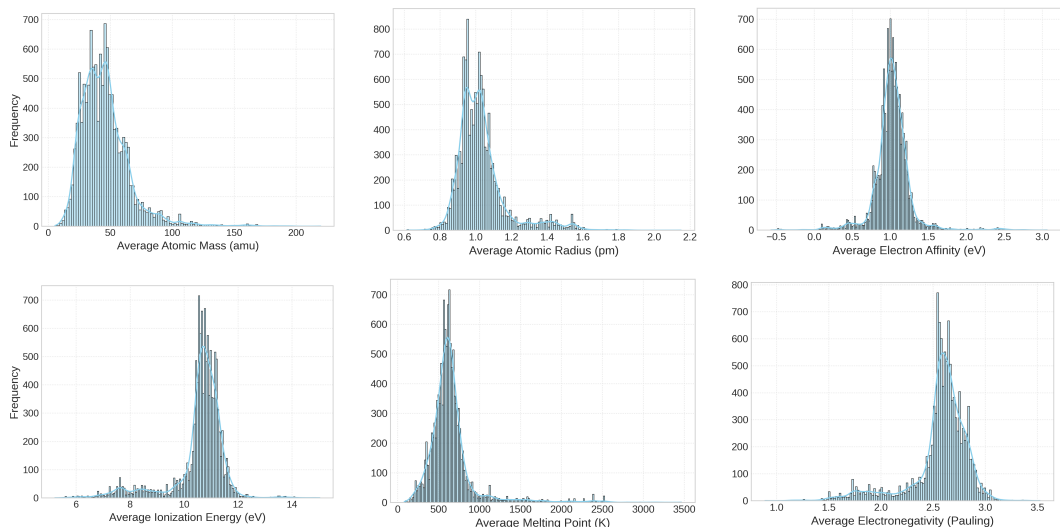


Figure 2: Property distributions of the featurized dataset. Frequency represents the number of total chemicals that contain a property with a value within that bin. Each property value is the average of the elemental properties for the composition of a specific chemical.

**Utilizing the ActionGraph.** In order to assess the impact of the topology of the synthesis reaction on predicting precursors and operations, we instead train a $k$-NN model on the second dataset. The ActionGraph structure is incorporated in featurization via a multi-step process:

1. Find the maximum number of nodes present in all of the ActionGraph graphs. This was determined to be 31.

2. During featurization of a target ActionGraph, the adjacency matrix is extracted then padded to a 31-by-31 matrix, flattened, then scaled using STANDARDSCALER.

3. Principal component analysis (PCA) is applied to reduce the dimensionality of the adjacency matrix such that the "adjacency matrix vector" for each ActionGraph is the same length. This is then concatenated with the existing feature vector, which was determined the same way as before.

While this method could suffer from being too simple given the complexity of inorganic synthesis, it provides a straightforward approach for examining the effectiveness of using the ActionGraph without requiring complex neural network architectures.

**Model architecture.** For both datasets, a $k$-NN model with $k = 1$ is used, acting as a nearest-analog retrieval system. Specifically, the NEARESTNEIGHBORS model from SCIKIT-LEARN is employed with cosine distance as a similarity metric [PVG$^+$11]. This approach treats synthesis prediction as an information retrieval problem: when queried with a target material, the model identifies the most similar synthesis reaction in the training set set and adopts its complete recipe as the prediction for the new material.

The feature vectors for both models undergo scaling using SCIKIT-LEARN's QUANTILETRANS-FORMER with output_distribution="normal" and n_quantiles=min(1000, X_TRAIN). This nonlinear transformation maps the original feature distributions to follow a normal distribution, which helps improve model performance by reducing the impact of outliers and ensuring consistent feature scales. This is notably important for the elemental properties as different units cause significant differences in their feature scales.

A 75-25 train-test split was used for all experiments to ensure consistent evaluation across the different model configurations. All experiments were performed in a Python 3.12.9 environment with scikit-learn 1.61.

## 5 Experimental Validation

**Evaluation metrics.** We adopt several metrics from synthesis prediction literature to evaluate our models. F1 score measures the harmonic mean of precision and recall for precursor and operation prediction, providing a balanced assessment of prediction accuracy. Exact Match (EM) [KJS24] is a strict metric that counts a prediction as correct only when the entire set of predicted precursors exactly matches the reference set, with no missing or extra compounds. Operation Length Matching [KJVO21] measures the percentage of predictions where the number of synthesis operations matches the ground truth, reflecting the model's ability to capture the procedural complexity of synthesis.

To ensure robust evaluation, we employed a 75-25 train-test split of our dataset, resulting in 9,763 samples for training and 3,254 for testing. All reported metrics are computed on the test set. Feature scaling was performed using scikit-learn's QuantileTransformer to ensure robustness to outliers and non-Gaussian distributions in the chemical property features.

### 5.1 Principal Component Analysis and Feature Space Visualization



(a) PCA explained variance ratio
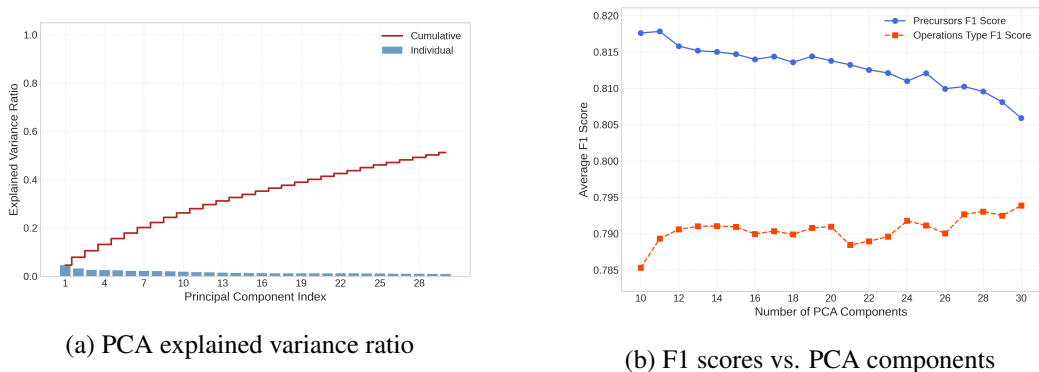
(b) F1 scores vs. PCA components

Figure 3: PCA analysis of ActionGraph adjacency matrices: (a) Individual and cumulative explained variance showing how information is distributed across components; (b) Trade-off between precursor and operation F1 scores as more structural information is incorporated.

The PCA variance analysis (Figure 3a) shows that structural information in ActionGraph adjacency matrices is distributed across many components, with the cumulative explained variance reaching

approximately 50% at 30 components. This indicates that synthesis graph structures are complex and no single component dominates the representation. The individual contribution of each component (blue bars) decreases gradually, suggesting that even later components contain meaningful information about synthesis pathways.

Figure 3b reveals a clear trade-off as we increase the number of PCA components: precursor F1 score peaks at 10-11 components (0.818), after which it gradually declines, while operation F1 score improves with additional components, rising from 0.785 at 10 components to 0.794 at 30 components. This divergence suggests that precursor prediction benefits from a more composition-dominated representation, while operation prediction leverages richer structural encoding.



(a) Baseline features

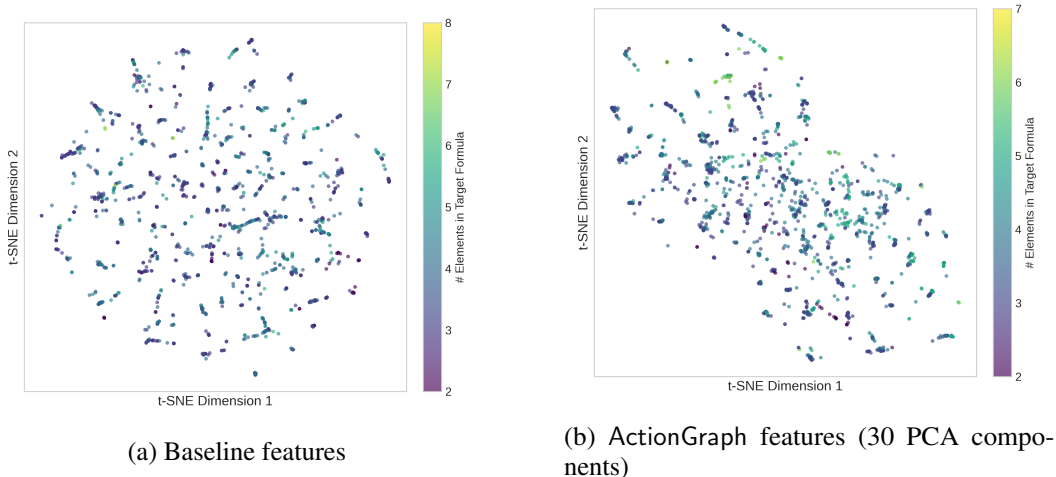(b) ActionGraph features (30 PCA components)

Figure 4: t-SNE visualizations of feature spaces colored by number of elements in target formula: (a) Baseline features showing less defined clustering; (b) ActionGraph features showing more pronounced structure, indicating better encoding of synthesis similarities.

To qualitatively assess the effect of ActionGraph structural features, we visualized the feature spaces using t-SNE (Figure 4). The ActionGraph-based features (Figure 4b) produce more pronounced and structured clusters compared to the baseline representation (Figure 4a). This suggests that the inclusion of synthesis topology enables the model to group reactions with similar procedural flows, not just similar compositions.

## 5.2 Quantitative Results and Discussion

Table 1 summarizes the key quantitative results for the baseline and ActionGraph-enhanced models at representative PCA component counts. The ActionGraph model with 10 PCA components (AG-PCA10) achieves the highest precursor F1 (0.818), while the 20-component model (AG-PCA20) yields the best operation F1 (0.791). Notably, operation length matching improves dramatically from 15.8% (baseline) to 53.3% (AG-PCA20), a $3.4\times$ increase, indicating that structural information is critical for capturing the procedural fidelity of synthesis.

Table 1: Performance comparison across models (test set $n = 2,968$)

| Metric | Baseline | AG-PCA10 | AG-PCA20 |
|---|---|---|---|
| **Precursors** | | | |
| F1 | 0.7995 | 0.8176 | 0.8138 |
| Exact Match | 0.509 | 0.567 | 0.550 |
| **Operations** | | | |
| F1 | 0.763 | 0.785 | 0.791 |
| Length Match (%) | 15.8 | 47.5 | 53.3 |

8

The observed divergence in F1 trends (Figure 3b) highlights the importance of tuning the number of PCA components to balance precursor and operation prediction objectives. The optimal regime appears to be 10–15 components for precursor-centric tasks, and 20–30 for operation-centric tasks.

*Note*: using different seeds for the train-test split did not significantly alter results.

**Structural analysis.** The PCA variance plot (Figure 3a) confirms that no single component dominates the structural encoding; rather, meaningful information is distributed across many components, reflecting the complexity of inorganic synthesis pathways. The t-SNE visualizations further support the conclusion that ActionGraph features capture procedural similarities not accessible to composition-only models.

**Interpretability and limitations.** While ActionGraph encoding substantially improves procedural fidelity (operation sequence and length), exact operation matches remain below 9%, reflecting the inherent ambiguity and sparsity of metadata in text-mined synthesis descriptions. Future work may address this gap by incorporating richer operation metadata or by learning hierarchical graph representations.

Overall, these results validate the core hypothesis: encoding synthesis structure via the ActionGraph framework enhances prediction of synthesis operations and improves procedural accuracy, representing a step toward robust inorganic synthesis pathway prediction from product formulas alone.

**Computational resources.** Experiments were conducted on an ASUS ROG Zephyrus G14 laptop with AMD Ryzen 9 4900HS processor and NVIDIA GeForce RTX 2060 Max-Q GPU. Each experiment completed in under one minute, with runtime increasing proportionally with the number of PCA components used.

# 6    Conclusion and Outlook

In this work, we introduced the ActionGraph framework and demonstrated its effectiveness for predicting inorganic synthesis pathways given only the target product. Our approach encodes synthesis reactions as directed acyclic graphs and utilizes $k$-nearest neighbors with PCA-reduced adjacency matrices to capture structural information. The results show that incorporating ActionGraph structure improves both precursor and operation prediction compared to composition-only baselines.

A key insight from our experiments is the trade-off between precursor and operation prediction performance as a function of structural encoding depth. Precursor F1 scores peak at 10-11 PCA components before declining, while operation F1 scores continue improving with additional components. This suggests that precursor selection relies primarily on composition information, while operation prediction benefits from richer structural encoding. The dramatic improvement in operation length matching accuracy (from 15.8% to 53.3%) demonstrates that ActionGraph structural information is particularly valuable for capturing procedural aspects of synthesis.

Looking forward, several promising directions emerge. First, more sophisticated featurization methods could enhance both the chemical and structural representations in our model. While our current element fraction and property-based features provide a reasonable baseline, incorporating crystal structure information or learned embeddings could improve prediction accuracy [GST+25]. Second, the ActionGraph framework could be extended to include more detailed operation conditions and intermediate phases, provided sufficient data becomes available. Third, alternative dimensionality reduction techniques or graph embedding methods might preserve more structural information than our current PCA approach.

The significant improvement in operation length matching suggests that the ActionGraph framework effectively capture synthesis procedure complexity, but the relatively modest gains in F1 scores indicate room for further refinement. Future work could explore dynamic weighting between compositional and structural features to optimize for specific prediction tasks.

We envision this framework becoming part of a closed-loop system where synthesis predictions guide automated experimentation, generating new data that further refines the model. Such a system could dramatically accelerate inorganic materials discovery by reducing reliance on trial-and-error approaches. While substantial challenges remain in data quality and model sophistication, the ActionGraph approach represents an important step toward data-driven synthesis planning for novel inorganic materials.

# References

[AGM⁺25] Daniel Alabi, Sainyam Galhotra, Shagufta Mehnaz, Zeyu Song, and Eugene Wu. Privacy and security in distributed data markets. In *Companion of the 2025 International Conference on Management of Data, SIGMOD/PODS 2025, Berlin, Germany, June 22-27, 2025*. ACM, 2025.

[AW24] Daniel Alabi and Eugene Wu. Empiredb: Data system to accelerate computational sciences, 2024.

[BHB⁺18] Peter Battaglia, Jessica Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018.

[BL07] Simon J. L. Billinge and Igor Levin. The problem with determining atomic structure at the nanoscale. *Science*, 316(5824):561–565, 2007.

[BP24] Simon J. L. Billinge and Thomas Proffen. Machine learning in crystallography and structural science. *Acta Crystallographica Section A*, 80(2):139–145, 2024.

[CH67] T. Cover and P. Hart. Nearest neighbor pattern classification. *IEEE Transactions on Information Theory*, 13(1):21–27, 1967.

[CT17] Giuseppe Carleo and Matthias Troyer. Solving the quantum many-body problem with artificial neural networks. *Science*, 355(6325):602–606, 2017.

[ESHH⁺24] Amer Marwan El-Samman, Incé Amina Husain, Mai Huynh, Stefano De Castro, Brooke Morton, and Stijn De Baerdemacker. Global geometry of chemical graph neural network representations in terms of chemical moieties. *Digital Discovery*, 3:544–557, 2024.

[GBWD⁺18] Rafael Gómez-Bombarelli, Jennifer N Wei, David Duvenaud, José Miguel Hernández-Lobato, Benjamin Sánchez-Lengeling, Dennis Sheberla, Jorge Aguilera-Iparraguirre, Timothy D Hirzel, Ryan P Adams, and Alán Aspuru-Guzik. Automatic chemical design using a data-driven continuous representation of molecules. *ACS Central Science*, 4(2):268–276, 2018.

[GSR⁺17] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, 2017.

[GST⁺25] Gabe Guo, Tristan Luca Saidi, Maxwell W. Terban, Michele Valsecchi, Simon J. L. Billinge, and Hod Lipson. Ab initio structure solutions from nanocrystalline powder diffraction data via diffusion models. *Nature Materials*, 2025.

[HBH⁺22] Haoyan Huo, Christopher J. Bartel, Tanjin He, Amalie Trewartha, Alexander Dunn, Bin Ouyang, Anubhav Jain, and Gerbrand Ceder. Machine-learning rationalization and prediction of solid-state synthesis conditions. *Chemistry of Materials*, 34(16):7323–7336, Aug 2022.

[HHB⁺23] Tanjin He, Haoyan Huo, Christopher J. Bartel, Zheren Wang, Kevin Cruse, and Gerbrand Ceder. Precursor recommendation for inorganic synthesis by machine learning materials similarity from scientific literature. *Science Advances*, 9(23):eadg8180, June 2023. arXiv:2302.02303 [cond-mat].

[JEP⁺21] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathrin Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.

[JGD⁺08] P. Juhas, L. Granlund, P. M. Duxbury, W. F. Punch, and S. J. L. Billinge. The Liga algorithm for ab initio determination of nanostructure. *Acta Crystallographica Section A*, 64:631–640, November 2008.

[KHH+19] Olga Kononova, Haoyan Huo, Tanjin He, Ziqin Rong, Tiago Botari, Wenhao Sun, Vahe Tshitoyan, and Gerbrand Ceder. Text-mined dataset of inorganic materials synthesis recipes. *Scientific Data*, 6(1):203, Oct 2019.

[KJS24] Seongmin Kim, Yousung Jung, and Joshua Schrier. Large Language Models for Inorganic Synthesis Predictions. *Journal of the American Chemical Society*, 146(29):19654–19659, July 2024.

[KJVO21] Christopher Karpovich, Zach Jensen, Vineeth Venugopal, and Elsa Olivetti. Inorganic Synthesis Reaction Condition Prediction with Generative Machine Learning, December 2021. arXiv:2112.09612 [cond-mat].

[KNG+24] Seongmin Kim, Juhwan Noh, Geun Ho Gu, Shuan Chen, and Yousung Jung. Predicting synthesis recipes of inorganic crystal materials using elementwise template formulation. *Chemical Science*, 15(3):1039–1045, 2024.

[LW17] Jun Liu and Lei Wang. Accelerated monte carlo simulations using restricted boltzmann machines. *Physical Review B*, 96(14):144426, 2017.

[MDP21] Matthew J. McDermott, Shyam S. Dwaraknath, and Kristin A. Persson. A graph-based network for predicting chemical reaction pathways in solid-state materials synthesis. *Nature Communications*, 12(1):3097, May 2021.

[MPM+20] Benjamin P MacLeod, Fraser G L Parlane, Thomas D Morrissey, Florian Häse, Loïc M Roch, Kevan E Dettelbach, Rafael Moreira, Peter J Yunker, Alán Aspuru-Guzik, Jason E Hein, et al. Self-driving laboratory for accelerated discovery of thin-film materials. *Science Advances*, 6(20):eaaz8867, 2020.

[MST21] Kyle Mills, Michael Spanner, and Isaac Tamblyn. The role of machine learning in scientific simulations. *Nature Reviews Physics*, 3(6):447–460, 2021.

[NLNP25] Heewoong Noh, Namkyeong Lee, Gyoung S. Na, and Chanyoung Park. Retrieval-retro: Retrieval-based inorganic retrosynthesis with expert knowledge, 2025.

[NNZTB25] Tanaporn Na Narong, Zoe N. Zachko, Steven B. Torrisi, and Simon J. L. Billinge. Interpretable multimodal machine learning analysis of x-ray absorption near-edge spectra and pair distribution functions. *npj Computational Materials*, 11(1):98, 2025.

[PIT18] Mariya Popova, Olexandr Isayev, and Alexander Tropsha. Deep reinforcement learning for de novo drug design. *Science Advances*, 4(7):eaap7885, 2018.

[PVG+11] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

[RDK+19] David Rolnick, Priya L Donti, Lynn H Kaack, Kelly Kochanski, Alexandre Lacoste, Krishna Sankaran, Andrew S Ross, Nicolas Milojevic-Dupont, Natasha Jaques, Anna Waldman-Brown, et al. Tackling climate change with machine learning. *arXiv preprint arXiv:1906.05433*, 2019.

[RPK19] Maziar Raissi, Paris Perdikaris, and George E Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686–707, 2019.

[SLAG18] Benjamin Sanchez-Lengeling and Alán Aspuru-Guzik. Inverse molecular design using machine learning: Generative models for matter engineering. *Science*, 361(6400):360–365, 2018.

[SSK+18] Kristof T Schütt, Huziel E Sauceda, Pieter-Jan Kindermans, Alexandre Tkatchenko, and Klaus-Robert Müller. Schnet: A continuous-filter convolutional neural network for modeling quantum interactions. *NeurIPS*, 2018.

[SZH+21] Nathan J. Szymanski, Yan Zeng, Haoyan Huo, Christopher J. Bartel, Haegyeom Kim, and Gerbrand Ceder. Toward autonomous design and synthesis of novel inorganic materials. *Materials Horizons*, 8(8):2169–2198, 2021.

[vNLH17] Evert PL van Nieuwenburg, Yi-Hong Liu, and Sebastian D Huber. Learning phase transitions by confusion. *Nature Physics*, 13(5):435–439, 2017.

[WSQ+25] Yixin Wei, Leyu Shan, Tong Qiu, Diannan Lu, and Zheng Liu. Machine learning-assisted retrosynthesis planning: Current status and future prospects. *Chinese Journal of Chemical Engineering*, 77:273–292, January 2025.

[ZIA+19] Alex Zhavoronkov, Yan A Ivanenkov, Alex Aliper, Maksim Veselov, Vladimir Aladinskiy, Anastasiya Aladinskaya, Victor A Terentiev, Daniil Polykovskiy, Mikhail Kuznetsov, Arman Asadulaev, et al. Deep learning enables rapid identification of potent ddr1 kinase inhibitors. *Nature Biotechnology*, 37(9):1038–1040, 2019.