# ULSA: Unified Language of Synthesis Actions for Representation of Synthesis Protocols

Zheren Wang[1,2,a], Kevin Cruse[1,2,a], Yuxing Fei[1,2], Ann Chia[1,c], Yan Zeng[2], Haoyan Huo[1,2], Tanjin He[1,2], Bowen Deng[1,2], Olga Kononova[1,*,b], and Gerbrand Ceder[1,2,*]

[1]Department of Materials Science & Engineering, University of California, Berkeley, CA 94720, USA

[2] Materials Sciences Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA

[*]Corresponding author: olga_kononova@berkeley.edu and gceder@berkeley.edu

[a]Equal contribution

[b]Present address: Roivant Sciences, New York, NY 10036, USA

[c]Present address: Nanyang Technological University, Republic of Singapore, 639798

Applying AI power to predict syntheses of novel materials requires high-quality, large-scale datasets. Extraction of synthesis information from scientific publications is still challenging, especially for extracting synthesis actions, because of the lack of a comprehensive labeled dataset using a solid, robust, and well-established ontology for describing synthesis procedures. In this work, we propose the first *unified language of synthesis actions* (ULSA) for describing ceramics synthesis procedures. We created a dataset of 3,040 synthesis procedures annotated by domain experts according to the proposed ULSA scheme. To demonstrate the capabilities of ULSA, we built a neural network-based model to map arbitrary ceramics synthesis paragraphs into ULSA and used it to construct synthesis flowcharts for synthesis procedures. Analysis for the flowcharts showed that (a) ULSA covers essential vocabulary used by researchers when describing synthesis procedures and (b) it can capture important features of synthesis protocols. This work is an important step towards creating a synthesis ontology and a solid foundation for autonomous robotic synthesis.

# 1 Introduction

In the past decade, we have witnessed the growing success of data-driven and artificial intelligence (AI)-based methodologies promoting breakthroughs in predicting materials structure, properties, and functionality [1, 2, 3]. Nonetheless, adapting the power of AI to predict and control materials synthesis and fabrication is still challenging and requires substantial effort in gathering high-quality large-scale datasets. One approach to gather such datasets of synthesis parameters and conditions would be running high-throughput experiments. This requires a costly setup and substantial human labor and expertise, and is typically limited to a small part of chemical space. Another way to acquire the data or augment existing datasets is to extract information about materials synthesis from the wealth of scientific publications (e.g. papers, archives, patents) available online.

Scientific text mining has received its recognition in the past few years [4, 5, 6], providing the materials science community with datasets on a variety of materials and their properties [7, 8, 9] as well as synthesis protocols [10, 11, 12]. Nonetheless, a majority of these text mining studies have been focused on extracting chemical entities such as material names, formulas, properties, and other characteristics [13, 14, 15, 16, 17]. There have only been a few attempts to extract information about chemical synthesis and reactions and compile them into the flowchart of synthesis actions [18, 12, 19, 20]. This is largely due to the lack of comprehensive labeled datasets or annotation schema needed to train algorithms. Indeed, publicly available large-scale collections of standardized labeled data for named entities recognition (NER) tasks are well established in the biochemical and biomedical domains (GENIA [21], CHEMDNER [22]). Materials science datasets are less standardized and mainly task-specific [23, 24, 25]. To the best of our knowledge, the only publicly available annotated corpus of materials synthesis protocols was published by Mysore et al. [12]. It contains 230 labeled synthesis paragraphs with labels assigned to material entities, synthesis actions, and other synthesis attributes.

A major obstacle in annotating synthesis actions in the text corpora is the lack of a solid, robust, and well-established ontology for describing synthesis procedures in materials science [26]. Indeed, researchers prefer to vaguely sketch "methods" sections of the manuscript in general human-readable language rather than follow a specific protocol. This significantly impacts reproducibility of the results, not to mention

ambiguity in understanding even when read by a human expert [26]. While such ambiguity is inconvenient for human readers, the growing interest in automated AI-guided materials synthesis demands a robust and unified language for describing synthesis protocols in order to make them applicable to autonomous robotic platforms [27, 28, 29].

In this work, we propose a *unified language of synthesis actions* (ULSA) to describe solid-state, sol-gel, precipitation, and solvo-/hydrothermal synthesis procedures. We also present a labeled dataset of 3,040 synthesis sentences created using the proposed ULSA schema. To verify applicability of the ULSA and the dataset, we trained a neural network-based model that identifies a sequence of synthesis actions in a paragraph, maps them into the ULSA, and builds a graph of the synthesis procedure (Figure 1). Analysis of the graphs from thousands of paragraphs has shown that this ULSA vocabulary is large enough to obtain high-accuracy extraction of synthesis actions as well as to pick the important features of each of the aforementioned synthesis types. The dataset and the script for building such a synthesis flowchart is publicly available. We anticipate that these results will be widely used by the researchers interested in scientific text mining and will help to achieve a breakthrough in predictive and AI-guided autonomous materials synthesis.

## 2 Methodology

### 2.1 Unified Language of Synthesis Actions and annotation scheme

To unify terminology used to describe a synthesis procedure, we defined 8 *action terms* that unambiguously identify a type of synthesis action. Every action word (or multi-word phrase) in the dataset is mapped to the corresponding action term according to the following rule: the word (or multi-word phrase) is recognized as an action if it (a) results in modification of the state of the material or mixture during the synthesis or (b) carries a piece of information affecting the outcome of the synthesis procedure. The action terms used within the unified language are explained below. In each example, the text underlined is the word or phrase that is annotated.

- `Starting`: A word or a multi-word phrase that marks the beginning of a synthesis procedure. Specifically, this often indicates which materials will be produced. For example: *"PMN-PT was synthesized by the columbite precursor method"*, *"Solid-state synthesis was used to prepare the target material"*,

*"The powder was <u>obtained</u> after the aforementioned procedure"*.

- **Mixing**: A word or a multi-word phrase that marks the combination of different materials (in a solid or liquid phase) to form one substance or mass. For example: *"Precursors were weighted and <u>ball -milled</u>"*, *"Precursors were <u>mixed</u> in appropriate amounts"*, *"$Sb_2O_3$ is <u>added</u> to the solution"*, *"The solution was <u>neutralized</u>"*, *"The mixture was <u>stabilized</u> by the addition of sodium citrate"*.

- **Purification**: A word or a multi-word phrase that marks the separation of the sample phases. This also includes drying of a material. For example: *"Samples were <u>exfoliated</u> from substrates"*, *"The liquid was discarded and the remaining product was <u>filtered off</u> and <u>washed</u> several times with distilled water"*, *"The precursors were <u>heated</u> in order to remove the moisture"*, *"The precipitation was <u>collected</u> by washing the solution in distilled water"*.

- **Heating**: A word or a multi-word phrase that marks increasing or maintaining high temperature for the purpose of obtaining a specific sample phase or promoting a reaction rather than drying a sample. For example: *"The powder sample was <u>annealed</u> to obtain a crystalline phase"*, *"The mixture was subjected to <u>heating</u> at 240 °C for 24 h"*.

- **Cooling**: A word or a multi-word phrase that marks rapid, regular, or slow cooling of a sample. For example: *"The product was <u>cooled</u> down to room temperature in the furnace"*, *"The sample was <u>quenched</u> rapidly in the solid $CO_2$"*, *"The products was <u>left to cool</u> down to room temperature"*.

- **Shaping**: A word or a multi-word phrase that marks the compression of powder or forming the sample to a specific shape. For example: *"The powder was <u>pressed</u> into circular pellets"*, *"The powder was then <u>pelletized</u> with a uniaxial press"*.

- **Reaction**: A word or a multi-word phrase that marks a transformation without any external action. For example: *"The sample was <u>left to react</u> for 6 hrs"*, *"The temperature was <u>kept</u> at 1000 K"*, *"The solution was <u>maintained</u> at 200 K for 12 hrs"*.

- **Miscellaneous**: A word or a multi-word phrase that marks an action done on a sample that either does not induce any transformation of the sample or does not belong to any of the above classes. *"The*

pellets were _placed_ in a sealed alumina crucible", "The reaction vessel was _wrapped_ with aluminum foil", "The sample was _sealed_ in a tube", "The gel was _transferred_ to an oven".

## 2.2   Dataset annotation

To annotate synthesis paragraphs with the unified language of synthesis actions (ULSA), we selected 535 synthesis paragraphs from the database of 420K full-text publications acquired previously [11]. The paragraphs where chosen to proportionally represent four major types of ceramics synthesis: solid-state, sol-gel, solvo-/hydrothermal, and precipitation. The details of the content acquisition and synthesis type classification have been described in previous papers [11, 30].

The 535 paragraphs consisted of 3,781 tokenized sentences [14]. First, each sentence was classified as either related to synthesis or not related to synthesis. The latter case usually contains sentences about product characterization and other details. Next, we isolated 3,040 synthesis sentences and assigned labels to each word or multi-word phrase in the sentence on the basis of the ULSA protocol with annotation schema described in Section 2.1. Only words and phrases describing synthesis actions were annotated. The final dataset consists of these 3,040 labeled synthesis sentences. All annotations were performed using a custom Amazon Mechanical Turk-based server.

## 2.3   Annotation decisions and ambiguous cases

The ULSA was developed based on the authors' own experiences with the extraction of information from materials synthesis paragraphs [11] and extensive communication with experimentalists actively involved in various types of materials synthesis research. The annotation schema and the choice of action terms were designed to provide maximum flexibility to future users and allow them to adjust the schema according their preferences and tasks. For example, the annotated multi-word phrases such as "subjected to heating", "left to react", and "heated to evaporate" were handled as one entity. This way, they can be split into individual terms or modified later with a simple set of rules to make a customized labeled dataset.

It is important to keep in mind that we mapped words into the terms of synthesis action per sentence, meaning that we used only information in the context of a given sentence to make a decision about the annotation of a word, rather than the whole paragraph. The reason for this choice is the multiple and

diverse possibilities to combine and augment sentences leading to different meanings of the terms. The interpretation of the whole text or paragraph is an entirely separate field of research that is outside the scope of this work.

We chose to annotate those words that are characteristic of a synthesis procedure or result in the transformation of a substance. In other words, those actions which are usually performed by default are not annotated. For example, in the sentence "the solution was sealed in an autoclave", no terms would be annotated as actions since the sealing step for hydrothermal synthesis is considered a default step. Similarly, in the sentence "the precursors were weighed and mixed," the term "weighed" is not a synthesis action since it is to be expected in synthesis, while "mixing" is a synthesis action because it may have a specific condition and transform the sample, or can be preceded by calcination of the precursors in other syntheses.

The exclusion from this rule is the `Starting` action. Even terms belonging to this action do not bring any special information or explicit action to the synthesis, we chose to distinguish "starting" actions because in a substantial number of cases they can serve as flags to separate multiple synthesis procedures from one another. An illustration of this situation is when precursors are prepared prior to synthesizing a target material, as in sol-gel synthesis.

For the annotation of `Mixing` synthesis actions, we did not differentiate between powder mixing, ball milling (grinding), addition of droplets, or dissolving of substances. In many situations, this precise definition depends on the solubility of reactants and mixing environment, as well as on other details of the procedure that are never explicitly mentioned in the text. We leave it up to a user to create their own application-based definitions of these mixing categories. Nonetheless, in the application below we provide a rule-based example of how these types of synthesis actions can be identified in the text.

The `Miscellaneous` action term was introduced to make room for those synthesis actions that are not typical or do not fall into any other category but nevertheless appear as a synthesis action within our definitions. While `Miscellaneous` action terms can be easily confused with `Reaction` actions or non-actions, the decision depends on the sentence context and can be arbitrarily extended or removed by a user. Comparing "the sample was kept in the cruicible" and "the sample was kept overnight," the former is not a synthesis action while the latter should be considered an important synthesis step.

Ambiguous situations as in the ones mentioned above are ubiquitous in descriptions of syntheses. A substantial amount of these situations occur when authors try to be wordy or use flowery language when writing the synthesis methods. Unfortunately, this often presents a challenge for accurate machine interpretation of the text. We accounted for some of these cases when annotated the data as described below.

First, implicit mentions of synthesis actions (i.e. when a past participle form of a verb is used as a descriptive adjective referring to an already processed material) is the most frequent source of confusion. We chose to annotate these as synthesis actions. For example: "the <u>calcined</u> powder was <u>pressed</u> and <u>annealed</u>." In this sentence, the descriptive adjective *calcined* could be either a restatement of the fact that there was a calcination step or it could be additional information which had not been mentioned previously. These situations can be later resolved with a rule-based approach, hence we leave it as a task for users of the data.

The situation when a method is specified along with the synthesis action is also common. In a phrase of the form "transformed by a specific procedure," we consider only the key action (the transformation) as a synthesis action. For example: "the precipitates were <u>separated</u> by centrifugation." When required, the method can be retrieved with a set of simple rules.

Redundant action phrases are also abundant in many descriptions of the procedures. In a phrase of form "subjected to a process", we considered only the processing verb as a synthesis action. For example: "the samples were subjected to an initial <u>calcination</u> process."

Finally, phrases that attempt to reason the purpose of the action, such as "left to react", "brought to a boil", "heated to evaporate," are considered as one synthesis action. This is done for the purpose of providing flexibility to a user and to let them make a decision on how to treat these cases.

## 2.4 Synthesis terms mapping

We used lookup table (baseline) and neural network models to map synthesis sentences into the ULSA.

### 2.4.1 Baseline model

Two baseline models were implemented, both based on a lookup table. For the lookup table, we chose the most frequent words used to describe synthesis steps in the "methods" section of the papers. The first baseline model matches every token against the lookup table and assigns the corresponding action term if

any appear. The second baseline model uses information about the part of speech of a given word (assigned by SpaCy [31]) and matches only verbs against the lookup table.

### 2.4.2 Word embeddings

Word embeddings were used as a vectorized representation of the word tokens for the neural network model. To create an embedding, we trained a Word2Vec model [32] implemented in the Gensim library [33]. We used ∼420K paragraphs describing four synthesis types: solid-state, sol-gel, solvo-/hydrothermal and precipitation synthesis. The paragraphs were obtained as described in our previous work [11]. Prior to training, the text was normalized and tokenized using ChemDataExtractor [14]. Conjunctive adverbs describing consequences, such as "therefore", "whereas", and "next", were removed from the text. All quantity tokens were replaced with a keyword <NUM>, and all chemical formulas were replaced with keyword <CHEM>. All words that occur less than 5 times in the text corpus were replaced with the keyword <UNK>. We found that skip-gram with negative sampling loss (n = 10) performed best, and the final embedding dimension was set to 100.

### 2.4.3 Neural network model

We used a bi-directional long short-term memory (bi-LSTM) neural network model to map synthesis tokens into the aforementioned action terms. The model was implemented using the Keras library (https://keras.io/) with latent dimensionality 32 and dropout probability 0.2. Word embeddings were used as model input. The categorical cross-entropy was calculated as the loss function. The labeled dataset was split into training, test, and validation sets using a 70:20:10 split, respectively. Early stopping was used to obtain the best performance.

## 2.5  Data analysis

### 2.5.1  Reassignment of mixing terms

For data analysis, we separated `Mixing` synthesis action terms into `Dispersion Mixing` and `Solution Mixing` whenever there was enough information to distinguish between the two, otherwise they were left as `Mixing` action. Here, `Dispersion Mixing` is identified either by explicit "dispersion" action words or by words such as "grinding" or "milling" plus any liquid environment. `Solution Mixing` is identified by a

list of specific action words such as "dissolve", "dropwise added", and others. For this, we constructed and traversed the dependency trees of the sentences using SpaCy library [31] and used dictionaries of common solution and mixing terms.

### 2.5.2 Constructing synthesis flowchart for paragraphs

For every paragraph in the set, we then applied the bi-LSTM mapping model (Section 2.4) to extract the sequence of action terms from every sentence. Next, we merged all the synthesis actions obtained from all sentences within the paragraph into a synthesis flowchart. This was performed with a rule-based approach by traversing grammar trees and analysing the surrounding words of each action term and comparing them to the words and action terms of the previous sentence. Finally, the flowchart of synthesis actions for a given paragraph was converted into an adjacency matrix. For this, synthesis action terms were ordered and assigned to rows and columns of the matrix and initialized with zeros, resulting in a 10 by 10 matrix for every paragraph (8 action terms from vocabulary of ULSA plus two additional terms for `Mixing` term). Whenever there was a step from action $i$ to action $j$, the corresponding value in the matrix was incremented by 1. The matrices for all paragraphs were flattened and merged together for further principal component analysis.

## 3 Results

### 3.1 Code and data availability

The dataset of 3,040 annotated synthesis sentences as well as the processing scripts are available at CederGroupHub/synthesis-action-retriever at https://doi.org/10.5281/zenodo.5644302. In the dataset, each record contains the raw sentence tokens concatenated with a space between each token and a list of objects, each containing a token and the tag assigned to that token. For example:

```
{
    "annotations" :
        [
            {
                "tag" : token_tag,
                "token" : token
            }
```

```
        ],
    "sentence" : sentence
}
```

The repository also contains a script for training a bi-LSTM model that can be used to map words into action terms. Users are not limited to using only the provided dataset, but can augment their usage with other labeled data as long as they satisfy the data format described above. Finally, we also share scripts used for the inference of synthesis actions terms and for building synthesis flowcharts for a list of paragraphs. Examples of model application are available as well.

## 3.2   Dataset statistics

The quantitative characteristics of the set are provided in Table 1 and displayed in Figure 2. Briefly, 535 synthesis paragraphs resulted in 3,781 sentences of which 3,040 describe actual synthesis procedures. While we tried to maintain an even distribution of the action terms in the labeled set, it is still highly skewed toward `Mixing` and `Purification` actions. This is not surprising, since mixing of precursors occupies any synthesis procedure and purification is required in almost any non-solid-state method for ceramics synthesis. `Heating` is the next most prevalent synthesis action since it is also one of the basic operations in ceramic synthesis.

To probe the robustness of ULSA and our annotation schema, we asked 6 human experts to annotate the same paragraphs in our dataset and used Fleiss' kappa score to estimate the inter-annotator agreement between the annotations [34]. In general, the Fleiss' kappa score evaluates the degree from -1 to 1 to which different annotators agree with one another above the agreement expected by pure chance. A positive Fleiss' kappa indicates good agreement, scores close to zero indicates near randomness in categorization, and negative scores indicate conflicting annotations. This is a generalized reliability metric and is useful for agreement between three or more annotators across three or more categories.

Table 2 lists the Fleiss' kappa scores for agreement between human experts annotating the sentences according to the schema described in Section 2.1. The table shows good agreement on distinguishing synthesis sentences from non-synthesis sentences, as well as for all and for each individual synthesis action, including non-actions. The agreement across all action terms is 0.83. Among those, the action terms with lower scores

11

are `Shaping` and `Miscellaneous`. The low score for `Miscellaneous` is expected since a wide range of actions which do not induce a transformation in the sample could be mapped into this category. The `Shaping` action term can also be associated with many synthesis operations. For instance, granulating procedures that break a sample into smaller chunks could be considered a `Shaping` action; at the same time, a bench chemist could consider "granulation" to be `Mixing` action term since it requires performing a grinding operation to obtain the new shape. Less ambiguous actions terms, such as `Heating` and `Mixing`, showed higher agreement.

## 3.3 Mapping synthesis procedures into a unified language of synthesis actions

### 3.3.1 Mapping model

As a first approach for mapping of synthesis paragraphs into ULSA, we used dictionary lookup constructed as described in Section 2.4.1. We use the labeled dataset of 3,040 sentences to assess the performance of the model. We considered two options: mapping of all sentence words and mapping the verbs only. In both cases, the overall accuracy of the prediction (i.e. F1 score) is ∼60-70% (Table 3). Nonetheless, mapping of all words shows relatively good recall and poor precision, while mapping of only verbs improves the precision but diminishes recall.

These results moved us toward considering a recurrent neural network model for mapping paragraphs into ULSA. The bi-LSTM model combined with word embeddings (Section 2.4.3) was trained on the labeled dataset of 3,040 sentences. The bi-LSTM model significantly improves mapping accuracy, yielding >90% F1 score. It is important to notice here that all the metrics for baseline and neural network models were computed per sentence, i.e. we evaluated the whole sentence being mapped correctly rather than individual terms.

There are a few reasons why the bi-LSTM model outperforms plain dictionary lookup. First, researchers use diverse vocabulary to describe synthesis procedures, hence there are unlimited possibilities in constructing a lookup table. For instance, "heating" can be referred as "calcining", "sintering", "firing", "burning", "heat treatment", and so on. In this case, a word embedding model helps to significantly improve the score even for those terms that have never appeared in training set (e.g. "degas", "triturate"). Second, a given verb is defined as a synthesis action term largely based on the context. Prominent examples are "heating

rate", "mixing environment", "ground powder", etc. That is well captured by the recurrent neural network architecture. Lastly, synthesis actions are not only denoted by verb tokens, but also by nouns, adjectives, and gerunds. This can be also learnt by the neural network better than by a set of rules.

In summary, we designed a neural network-based model that maps any synthesis paragraph into ULSA with high accuracy and significantly outperforms a plain dictionary lookup approach.

### 3.3.2 Analysis of action embeddings

To analyse how well the ULSA represents the space of synthesis operations commonly used when describing ceramics synthesis processes, we plotted a 2D projection of the word embeddings calculated with a t-SNE approach. The results are shown in Figure 3. To achieve a clear representation, we only analysed those verbs that appear more than 10 times. We then mapped these paragraphs into ULSA by using the bi-LSTM model. Those verbs that were assigned with a ULSA label are color-coded in the figure correspondingly, the other non-synthesis action terms are colored in grey.

First, we observe that the verbs mapped into ULSA and hence representing synthesis actions are all grouped in the top-left corner of the projection. Indeed, analysis of the individual words in the rest of the space showed that those are the words that generally appear in synthesis paragraphs but do not carry any information about the synthesis procedure. For instance, these are verbs denoting characterization of a material ("detect", "quantify", "examine", "measure"), naming of a sample ("denoted", "referred", "named", "labeled") or referring to a table or figure. The blob of dots in the middle of the plot are all words that were either mis-tokenized during text segmentation or mistakenly recognized as verbs by the SpaCy algorithm. In the embeddings mapping, these words are replaced with the `<UNK>` token.

A second interesting observation is that the embeddings of firing (blue dots), pelletizing (purple dots) and grinding into powder (orange dots) are all located next to each other. This agrees well with the fact that those actions together describe solid-state synthesis processes. Oppositely, the verbs describing solution mixing (orange dots) are in close proximity with the verbs referring to purification or drying (green dots). Similarly, verbs indicating cooling processes (magenta dots) and the verbs referring to reaction processes (red dots) are clustered together. This agrees with the often encountered constructions of "left to cool" or "kept and then cooled" describing the final steps of a given synthesis.

13

Taken together, these results demonstrate that (a) the embeddings model we created reflects well the similarity of the verbs used for synthesis descriptions and (b) the vocabulary of ULSA covers all common synthesis actions used in ceramics synthesis.

### 3.3.3   Analysis of graphs clustering

As we showed above, ULSA can capture well the vocabulary commonly used for the description of synthesis and, further, we were able to design a high-accuracy model that maps arbitrary synthesis descriptions into ULSA. However, we want also to make sure that unification of synthesis actions still allows for distinguishing between ceramics synthesis types. For that purpose, we constructed synthesis flowcharts for 4,000 paragraphs (1,000 per each synthesis type) randomly pulled from the set of 420K ceramics synthesis paragraph (see Section 2.5.2 for procedure description). For constructing the flowchart for a synthesis (represented by an adjacency matrix), we used the synthesis action terms assigned to each sentence in a paragraph. Additionally, we augmented `Mixing` actions with two categories, `Dispersion Mixing` and `Solution Mixing`, by using heuristics and dictionary lookup (Section 2.5.1). It is important to note here that we assume a linear order of synthesis actions, i.e. that the sequence of sentences and synthesis actions in a paragraph corresponds to the true sequence of synthesis steps done during experiment. According to our estimation, this assumption is violated only in 2% of paragraphs in the 420K paragraphs set.

All the adjacency matrices were flattened and concatenated, resulting in a matrix of size $100 \times 4000$, i.e. $10 \times 10$ matrix per each of 4,000 paragraphs, where 10 is the size of the ULSA vocabulary with two additional mixing actions. Next, principal component analysis was used to perform dimensionality reduction of the matrix.

Figure 4 displays the projection of the 1st and 2nd principal components for each synthesis flowchart with different colors corresponding to different types of syntheses. A few observations can be made from the plot. First, the data points corresponding to solid-state synthesis are narrowly clustered along a line with negative slope unlike the other synthesis types which are spread widely and whose linear fittings have positive inclination. Second, the clusters of data points for precipitation and hydrothermal synthesis almost completely overlap and partially overlap with sol-gel synthesis, while overlapping with solid-state synthesis is negligible.

These two observations agree well with the standard procedures associated with each of the four synthesis types. Indeed, solid-state syntheses usually operate with mixing powder precursors, firing the mixture, and obtaining final products; sol-gel synthesis is considered as a solid-state synthesis with solution-assisted mixing of precursors; hydrothermal and precipitation syntheses usually involve preparation of the sample in solution, then filtering (purification) to separate the liquid and obtain the final product instead of including a firing step.

To get further insights, we sampled and compared synthesis procedures along each of the fitted lines. The results show that the 1st principle component correlates with the involvement of solution mixing for precursors. In other words, the larger and more positive the data point along the 1st principle component, the more steps of dissolving and mixing precursors in solution as well as purification that data point involves. This agrees well with the fact that solid-state synthesis mostly operates with powders while hydrothermal and precipitation procedures are solution-based procedures, and sol-gel syntheses exist in between.

The 2nd principal component corresponds to the level of complexity of the syntheses procedure. The larger and more positive the data point along the 2nd principle component, the more synthesis steps become involved in the process. Interestingly, all four synthesis types exhibit simple synthesis procedures (fewer steps) and complex synthesis procedures (many steps). Nonetheless, solid-state synthesis has the largest deviation compared to hydrothermal and precipitation synthesis since solid-state procedures can involve multiple heating and re-grinding steps for the sample to obtain the desired material phase while in solution synthesis this can often be achieved in one or two steps.

## 4   Discussion and Conclusions

In this work, we aim to fill the gap in automated synthesis information extraction from scientific publications by proposing a unified language for synthesis actions (ULSA). We used the ULSA on an annotated set of 3,040 sentences about ceramics synthesis including solid-state, sol-gel, precipitation and solvo-/hydrothermal syntheses. The dataset is publicly available and can be easily customized by researchers accordingly to fit their application. As an example of such application, we used a recurrent neural network and grammar parsing to build a mapping model that converts written synthesis procedures into a ULSA-based synthesis

flowchart. Analysis of the results demonstrates that the ULSA vocabulary spans the essential set of words used by researchers to describe synthesis procedures in scientific literature and that the flowchart representation of synthesis constructed using ULSA can capture important synthesis features and distinguish between solid-state, sol-gel, precipitation and solvo-/hydrothermal synthesis methods.

Despite these promising results, the ULSA scheme still suffers from imperfections and can be significantly improved in the future. First, we only demonstrated that it works for ceramics synthesis, and synthesis techniques such as deposition, crystal growth, and others may require extending the ULSA vocabulary or reconsidering the definitions of some terms. Second, the scheme and methodology will benefit from a robust approach to distinguish between various mixing procedures. This includes separation between, for example, dissolving precursors and dispersive mixing in a liquid environment, using ball-milling to homogenize the sample and using high-energy ball-milling to actually achieve the final product, adding reagents to promote reaction and adding precursors to compensate for loss due to volatility, and other cases. We have demonstrated that the details of mixing are important for distinguishing between ceramics synthesis methods using simple heuristics, however, the scheme will benefit from a high-fidelity approach. Nonetheless, we anticipate that our results and the ULSA schema will help researchers to develop a data-oriented methodology to predict synthesis routes of novel materials.

Efficient and controllable materials synthesis is a bottleneck in technological breakthroughs. While predicting materials with advanced properties and functionality has been brought to a state-of-the-art level with the development of computational and data-driven approaches, the design and optimization of synthesis routes for those materials is still a tedious experimental task. The progress in inorganic materials synthesis is mainly impeded due to (a) lack of publicly available large-scale repositories with high-quality synthesis data and (b) lack of ontology and standardization for communication on synthesis protocols. Indeed, the first matter arises from the fact that the vast majority of experimental data gets buried in lab notebooks and is never published anywhere. As a result, researchers are liable to perform redundant and wasteful experimental screenings through those parameters of synthesis that have already been performed by someone, but are not reported. Even published experimental procedures face the problem of ambiguity of the language used by researchers. This creates a major challenge in acquiring synthesis data from publications

by automated approaches including text mining.

The advantage of the paradigm we establish in this work is that it brings us closer to addressing important questions in materials synthesis: *"How should we think about the synthesis process?"*, *"What is the minimum information required to unambiguously identify a synthesis procedure?"*, and *"Can synthesis be thought of as a combination of fixed action blocks augmented with attributes such as temperature, time, and environment, or are there other important aspects that have to be taken into account?"*. These questions will become crucial when transitioning towards AI-driven synthesis.

Recent developments in autonomous robotic synthesis and the attempts to "close the feedback loop" in making decisions for the next synthesis step make the question of synthesis ontology and unification especially important [27, 35, 28]. Indeed, while theoretical decision-making and AI-guided systems can operate with abstract synthesis representations, implementation of this methodology to an autonomous robotic platform will require well-defined and robust mapping onto a fixed set of manipulations and devices available to the robot. The unified language we propose in this work can become a solid foundation for the future development in this direction.

## Author Contributions

Z.W., K.C. O.K. and G.C. conceived the idea, and drafted the manuscript. Z.W., K.C., A.C. and O.K. implemented the algorithms and analyzed the data. Z.W., Y.F., and H.H built the annotation tool. Z.W., K.C., Y.F., Y.Z., and O.K. defined the annotation schema. Z.W., K.C, Y.F., H.H., T.H., and B.D. prepared the annotation dataset. All authors discussed and revised the manuscript.

## Conflicts of interest

There are no conflicts to declare.

## Acknowledgements

# Tables

| | Amount |
|---|---|
| Paragraphs used for annotation | 535 |
| Per synthesis type: | |
| – solid-state synthesis | 199 |
| – sol-gel synthesis | 51 |
| – solvo-/ hydrothermal synthesis | 148 |
| – precipitation | 137 |
| Total sentences | 3781 |
| Synthesis sentences | 3040 |
| Action tokens | 5547 |
| Per action category: | |
| – starting | 619 |
| – mixing | 1853 |
| – purification | 1080 |
| – heating | 973 |
| – cooling | 259 |
| – shaping | 225 |
| – reaction | 232 |
| – miscellaneous | 306 |

Table 1: Quantitative characteristics of the dataset chosen for annotation with ULSA schema.

|  | Score |
|---|---|
| Identification of synthesis sentences | 0.69 |
| Action terms tagging | 0.83 |
| Per action terms: | |
| – starting | 0.82 |
| – mixing | 0.86 |
| – purification | 0.79 |
| – heating | 0.84 |
| – cooling | 0.88 |
| – shaping | 0.59 |
| – reaction | 0.66 |
| – miscellaneous | 0.45 |
| – no action | 0.87 |

Table 2: Fleiss' kappa score for inter-annotator agreement using ULSA scheme.

| Model | Precision | Recall | F1 score |
|---|---|---|---|
| Baseline 1 | 0.54 | 0.61 | 0.57 |
| – solid-state | 0.53 | 0.72 | 0.61 |
| – sol-gel | 0.57 | 0.75 | 0.65 |
| – hydrothermal | 0.54 | 0.53 | 0.54 |
| – precipitation | 0.55 | 0.50 | 0.53 |
| Baseline 2 | 0.84 | 0.50 | 0.63 |
| – solid-state | 0.84 | 0.54 | 0.66 |
| – sol-gel | 0.79 | 0.62 | 0.69 |
| – hydrothermal | 0.84 | 0.47 | 0.61 |
| – precipitation | 0.84 | 0.44 | 0.54 |
| bi-LSTM | 0.90 | 0.88 | 0.89 |
| – solid-state | 0.90 | 0.90 | 0.90 |
| – sol-gel | 0.88 | 0.86 | 0.87 |
| – hydrothermal | 0.90 | 0.86 | 0.88 |
| – precipitation | 0.90 | 0.91 | 0.91 |

Table 3: Performance of baseline and bi-LSTM models for mapping synthesis sentence into ULSA terms. In Baseline 1, all words in the sentence are matched against a lookup table. In Baseline 2, only verbs tagged by SpaCy are matched against the lookup table. The quantities are computed per sentence, i.e. the number of sentences with all the action tokens identified and assigned correctly.
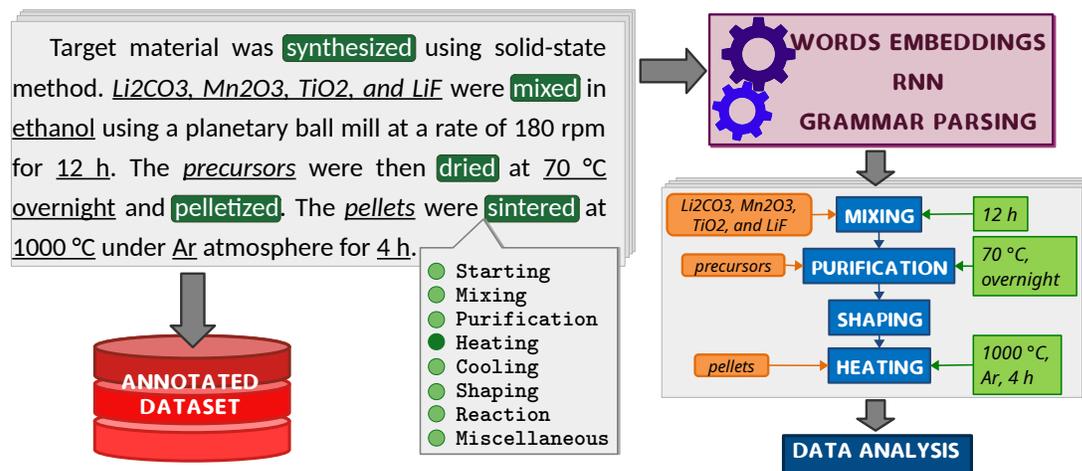
# Figures



Figure 1: **Schematic workflow of data annotation, extraction and analysis.** First, the set of paragraphs were annotated using an Amazon Mechanical Turk engine. Highlighted in green are the action token that were annotated and then extracted using a neural network model. Other highlighted tokens and phrases (i.e. synthesis action attributes and subjects) were obtained using rule-based sentence parsing solely for the purpose of data analysis and are not presented in the annotated dataset. The obtained labeled dataset is stored as single JSON file and is also used for training a neural network model to identify synthesis actions in the text. Obtained synthesis actions, attributes and subjects were converted into synthesis flowcharts that was further used for data analysis.
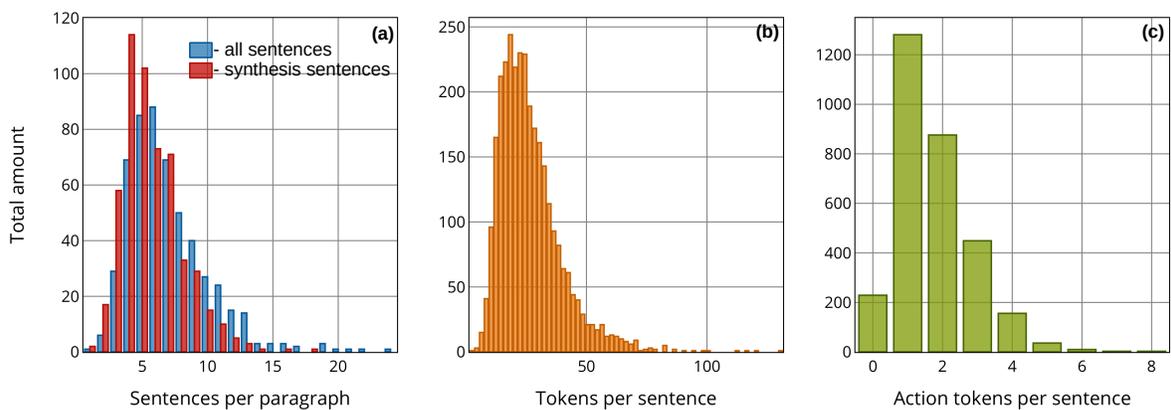
Figure 2: **Qualitative characteristics of the annotated dataset.** (a): Number of sentences per paragraph (blue), including sentences related to synthesis procedure (red). (b): Number of all tokens per sentence in the annotated set. (c): Number of tokens denoting a synthesis action per sentence in the annotated set.
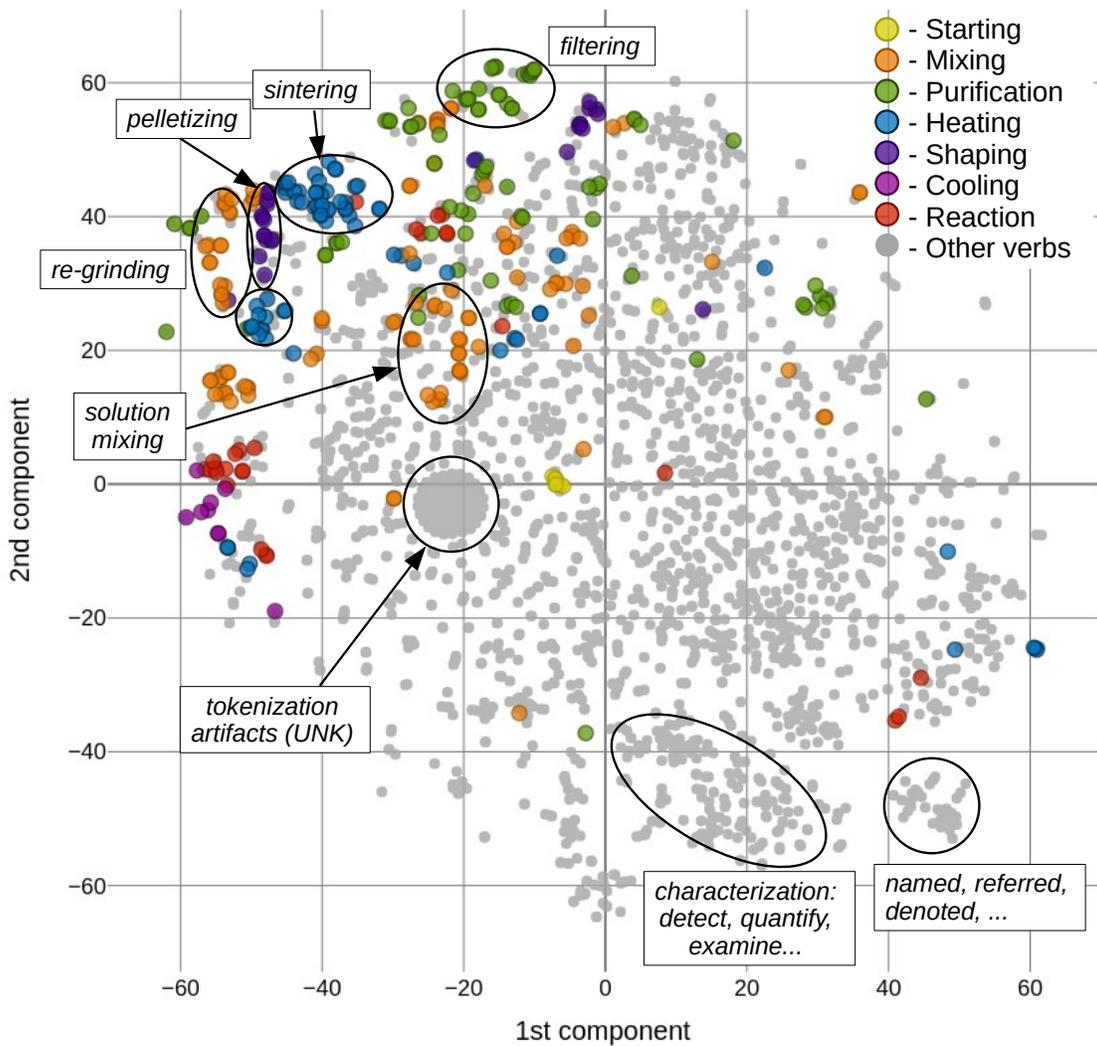
Figure 3: **2D projection of word embeddings vectors.** Shown are the most frequent verb tokens encountered in the set of 420K paragraphs describing a synthesis procedure. Highlighted in different colors are the vectors that correspond to the common verbs from the categories of synthesis actions used for annotation. Other prominent clusters of vectors are denoted with circles and labeled by a common term. Dimensionality reduction was performed using t-SNE approach.
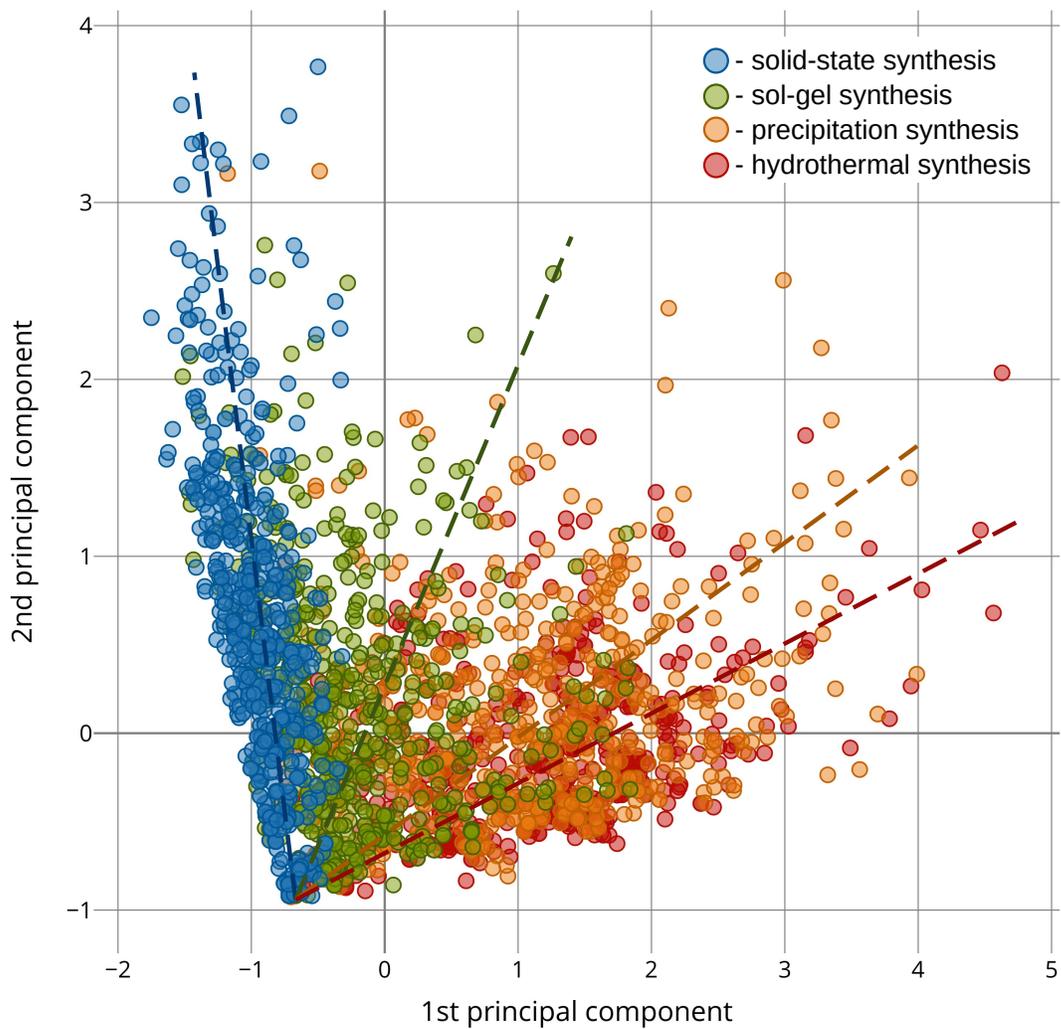
Figure 4: **Visualization of the first two principal components for the adjacency matrices of synthesis action graphs.** Each dot on the plot represent a synthesis graph colored according to its type. Dash lines display linear fitting of each data subset and show the overall direction for clustering of each synthesis graph. Note that the lines were shifted to have a common origin for representation purposes while preserving the slope.

# References

[1] Alberi, K. *et al.* The 2019 materials by design roadmap. *J. Phys. D: Appl. Phys* **52**, 013001 (2018).

[2] Himanen, L., Geurts, A., Foster, A. & Rinke, P. Data-driven materials science: Status, challenges, and perspectives. *Advanced Science* **6** (2019).

[3] Schmidt, J., Marques, M., Botti, S. & Marques, M. Recent advances and applications of machine learning in solid-state materials science. *npj Computational Materials* **5** (2019).

[4] Kononova, O. *et al.* Opportunities and challenges of text mining in materials research. *iScience* **24**, 102155 (2021).

[5] Olivetti, E. *et al.* Data-driven materials research enabled by natural language processing. *Appl. Phys. Rev.* **7**, 041317 (2020).

[6] Krallinger, M., Rabal, O., Lourenço, A., Oyarzabal, J. & Valencia, A. Information retrieval and text mining technologies for chemistry. *Chem. Rev.* **117**, 7673–7761 (2017).

[7] Huang, S. & Cole, J. M. A database of battery materials auto-generated using chemdataextractor. *Sci. Data* **7**, 1–13 (2020).

[8] Court, C. & Cole, J. M. Auto-generated materials database of curie and néel temperatures via semi-supervised relationship extraction. *Sci. Data* **5**, 180111 (2018).

[9] Court, C. & Cole, J. Magnetic and superconducting phase diagrams and transition temperatures predicted using text mining and machine learning. *npj Comput. Mater* **6**, 1–9 (2020).

[10] Kim, E. *et al.* Machine-learned and codified synthesis parameters of oxide materials. *Sci. Data* **4**, 170127 (2017).

[11] Kononova, O. *et al.* Text-mined dataset of inorganic materials synthesis recipes. *Sci. Data* **6**, 1–11 (2019).

[12] Mysore, S. *et al.* The materials science procedural text corpus: Annotating materials synthesis procedures with shallow semantic structures. *LAW 2019 - 13th Linguistic Annotation Workshop, Proceedings of the Workshop* 56–64 (2019). 1905.06939.

[13] Eltyeb, S. & Salim, N. Chemical named entities recognition: A review on approaches and applications. *J. Cheminform.* **6**, 1–12 (2014).

[14] Swain, M. C. & Cole, J. M. Chemdataextractor: a toolkit for automated extraction of chemical information from the scientific literature. *J. Chem. Inf. Model.* **56**, 1894–1904 (2016).

[15] Jessop, D. M., Adams, S. E., Willighagen, E. L., Hawizy, L. & Murray-Rust, P. Oscar4: a flexible architecture for chemical text-mining. *J. Cheminform.* **3**, 41 (2011).

[16] Weston, L. *et al.* Named entity recognition and normalization applied to large-scale information extraction from the materials science literature. *J. Chem. Inf. Model.* **59**, 3692–3702 (2019).

[17] Hiszpanski, A. *et al.* Nanomaterials synthesis insights from machine learning of scientific articles by extracting, structuring, and visualizing knowledge. *J. Chem. Inf. Model.* **60**, 2876–2887 (2020).

[18] Hawizy, L., Jessop, D. M., Adams, N. & Murray-Rust, P. Chemicaltagger: A tool for semantic text-mining in chemistry. *J. Cheminform.* **3**, 1–13 (2011).

[19] Kuniyoshi, F., Makino, K., Ozawa, J. & Miwa, M. Annotating and extracting synthesis process of all-solid-state batteries from scientific literature (2020). 2002.07339.

[20] Vaucher, A. *et al.* Automated extraction of chemical synthesis actions from experimental procedures. *Nat. Commun.* **11**, 3601 (2020).

[21] Kim, J.-D., Ohta, T., Tateisi, Y. & Tsujii, J. Genia corpus – a semantically annotated corpus for bio-textmining. *Bioinformatics* **19**, i180–i182 (2003).

[22] Krallinger, M. *et al.* The chemdner corpus of chemicals and drugs and its annotation principles. *J. Cheminform.* **7**, S2 (2015).

[23] Dieb, T., Yoshioka, M., Hara, S. & Newton, M. Framework for automatic information extraction from research papers on nanocrystal devices. *Beilstein J. Nanotechnol.* **6**, 1872–1882 (2015).

[24] Kulkarni, C., Xu, W., Ritter, A. & Machiraju, R. An annotated corpus for machine reading of instructions in wet lab protocols. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, 97–106 (Association for Computational Linguistics, Stroudsburg, PA, USA, 2018).

[25] Friedrich, A. *et al.* The SOFC-exp corpus and neural approaches to information extraction in the materials science domain. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 1255–1268 (Association for Computational Linguistics, 2020).

[26] Kim, E., Huang, K., Kononova, O., Ceder, G. & Olivetti, E. Distilling a materials synthesis ontology. *Matter* **1**, 8–12 (2019).

[27] Szymanski, N. J. *et al.* Toward autonomous design and synthesis of novel inorganic materials. *Mater. Horiz.* – (2021).

[28] Hammer, A. J. S., Leonov, A. I., Bell, N. L. & Cronin, L. Chemputation and the standardization of chemical informatics. *JACS Au* **0**, null (0).

[29] Mehr, S. H. M., Craven, M., Leonov, A. I., Keenan, G. & Cronin, L. A universal system for digitization and automatic execution of the chemical synthesis literature. *Science* **370**, 101–108 (2020).

[30] Huo, H. *et al.* Semi-supervised machine-learning classification of materials synthesis procedures. *npj Comput. Mater* **5**, 1–7 (2019).

[31] Honnibal, M. & Johnson, M. An improved non-monotonic transition system for dependency parsing. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, 1373–1378 (Association for Computational Linguistics, Lisbon, Portugal, 2015).

[32] Mikolov, T., Sutskever, I., Chen, K., Corrado, G. & Dean, J. Distributed representations of words and phrases and their compositionality (2013). `1310.4546`.

[33] Řehůřek, R. & Sojka, P. Software framework for topic modelling with large corpora. In *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, 45–50 (ELRA, Valletta, Malta, 2010).

[34] Fleiss, J. Measuring nominal scale agreement among many raters. *Psychological Bulletin* **76**, 378–382 (1971).

[35] Burger, B. *et al.* A mobile robotic chemist. *Nature* **583**, 237–241 (2020).