

ActionIE: Action Extraction from Scientific Literature with Programming Languages

Anonymous ACL submission

Abstract

001 Extraction of experimental procedures from human
002 language in scientific literature and patents
003 into actionable sequences in robotics language
004 holds immense significance in scientific do-
005 mains. Such an *action extraction* task is par-
006 ticularly challenging given the intricate details
007 and context-dependent nature of the instruc-
008 tions, especially in fields like chemistry where
009 reproducibility is paramount. In this paper, we
010 introduce ACTIONIE, a method that leverages
011 Large Language Models (LLMs) to bridge this
012 divide by converting actions written in natural
013 language into executable Python code. This en-
014 ables us to capture the entities of interest, and
015 the relationship between each action, given the
016 features of Programming Languages. Utilizing
017 linguistic cues identified by frequent patterns,
018 ActionIE provides an improved mechanism to
019 discern entities of interest. While our method is
020 broadly applicable, we exemplify its power in
021 the domain of chemical literature, wherein we
022 focus on extracting experimental procedures
023 for chemical synthesis. The code generated
024 by our method can be easily transformed into
025 robotics language which is in high demand in
026 scientific fields. Comprehensive experiments
027 demonstrate the superiority of our method. In
028 addition, we propose a graph-based metric to
029 more accurately reflect the precision of extrac-
030 tion. We also develop a dataset to address the
031 scarcity of scientific literature occurred in ex-
032 isting datasets.

033 1 Introduction

034 Recently, the integration of Natural Language Pro-
035 cessing (NLP) techniques into various scientific
036 fields has achieved significant success (Wang et al.,
037 2019; Soleimani et al., 2022; Song et al., 2023;
038 Lai et al., 2023). Among the applications, extract-
039 ing information from unstructured scientific litera-
040 ture has been one with growing significance (Guo
041 et al., 2022; Zhong et al., 2023a,b). For example,

Reaction Text

The residue is dissolved in EtOAc and washed sequentially with saturated Na₂CO₃ solution (2×), 10% aq. sodium dithionite (2×) and brine (1×), dried over Na₂SO₄, filtered and concentrated to give the title compound (7.47 g, 18.89 mmol, 90% purity) as a dark brown solid.

Chemical Reaction Actions	
No.	Action
1	ADD EtOAc
2	WASH with saturated Na ₂ CO ₃ solution 2 x
3	WASH with 10% aq. sodium dithionite 2 x
4	WASH with brine
5	DRYSOLUTION over Na ₂ SO ₄
6	FILTER keep filtrate
7	CONCENTRATE
8	YIELD title compound (7.47 g, 18.89 mmol, 90%)

Figure 1: An example of action extraction from literature that describes a sequence of chemical reaction actions. The text is drawn from Vaucher et al. (2020a).

chemists typically look through a wide range of publications to select candidate protocols for one organic synthesis scene, based on their own reading and repetitive trial-and-error procedures (Davies, 2019; Vaucher et al., 2021).

Therefore, structured chemical data, including reaction formulae, chemical entities, and experiment conditions, facilitates effective utilization and automatic analysis, such as indexing and searching by keywords; discovering and analyzing relations between entities; clustering related objects and discovering potential patterns; automatically executing protocols; and predicting and optimizing experiment conditions. Representatively, Figure 1 presents a case of structured chemical experimental procedure, essential for guiding practitioners in their laboratory work (Vaucher et al., 2020b; Zeng et al., 2023). This task involves extracting a sequence of chemical reaction actions from a scientific text passage, where each action is defined by

an operation and its corresponding attributes. For instance, in the example “*ADD EtOAc*” shown in Figure 1, “*ADD*” represents the operation, and the chemical “*EtOAc*” is the attribute.

However, the discovery of new chemical experimental procedures is scattered across unstructured scientific text and described in various writing styles, posing a significant challenge to the automatic creation of reaction action databases. Existing chemical databases, predominantly commercial ones such as Reaxys (Elsevier B.V., 2023), SciFinder (Chemical Abstracts Service (CAS), 2023), and Pistachio (NextMove Software, 2023), depend extensively on the manual contributions of domain experts. Analyzing, indexing, and utilizing scientific literature typically requires extensive and costly annotation or labeling by human experts. Moreover, this method is prone to errors due to the sheer volume of rapidly expanding scientific data. Despite the considerable manual effort, these databases prioritize storing information on the reactants, products, and reaction conditions, rather than the concrete sequences of chemical actions. This is primarily because manually designing these experiment procedures is both time-consuming and costly.

To tackle this issue, various studies have employed text mining techniques (Hawizy et al., 2011; Swain and Cole, 2016; Vaucher et al., 2020b; Wang et al., 2022b; Zeng et al., 2023) to automatically extract structured information on procedures from unstructured text, leveraging the advancements in NLP field. However, extracting experimental procedures remains a challenging task. One major hurdle is the complexity and variability of scientific language, which often features intricate sentence structures, domain-specific terminology, abbreviations, and acronyms. These elements pose substantial difficulties for sequential tagging-based approaches. For example, as shown in Figure 1, a text describing a series of chemical reaction actions includes the “*WASH*” operation followed by three chemicals. While sequential tagging-based methods might recognize the chemical compounds, they often struggle to accurately identify the operations and associate them with their corresponding attributes. Furthermore, the scarcity of large, annotated datasets poses an additional obstacle to training deep learning models on chemical experimental data effectively.

In this paper, we choose chemical experiment procedures as a case study, and explore the po-

tential of large language models (LLMs) to extract structured data from the complex and domain-specific language in chemical papers and patents. We propose a novel approach that frames the procedure extraction task as a code generation problem, where we express the experimental procedures as a series of pre-defined operations, and utilize the unique features of coding, such as classes, inheritance, and types, to structure this information. Our method leverages the capabilities of LLMs in few-shot in-context learning, reducing the need for large amounts of annotated data, and accelerating the preparation process. Moreover, our proposed framework also offers an easy solution to generate protocols for different automated platforms by applying different language configurations.

From the perspective of evaluation, we first pinpoint shortcomings within current evaluation metrics for the chemical action extraction task, and propose a novel metric based on graph-matching that substantially improves correlation with human judgments. Existing benchmarks largely concentrate on patent documents, which are inherently well-structured. To more accurately meet the real-world demands of practitioners, we meticulously annotate a test set derived from chemistry literature, which offers a more comprehensive evaluation of model performance. Notably, our new benchmark is considerably more extensive than previous ones, with an average length of 770.8 characters compared to 158.2 characters, providing a testing environment that mirrors realistic scenarios more closely. Experimentally, our method ActionIE demonstrates consistent superiority over strong baseline models, both against traditional benchmarks and our newly established testbed.

2 Related Work

The practice of using NLP in structured scientific data extraction has seen significant advancements, from utilizing traditional NLP techniques to integrating code generation methods into structure extraction, which is especially influenced by the growing capabilities of large language models (LLMs).

2.1 Action Extraction in Chemical Documents

The algorithms for action extraction in chemical texts evolve with the development of NLP. Earlier approaches, such as ChemDataExtractor (Swain and Cole, 2016) and ChemicalTagger (Hawizy et al., 2011), used part-of-speech tagging tech-

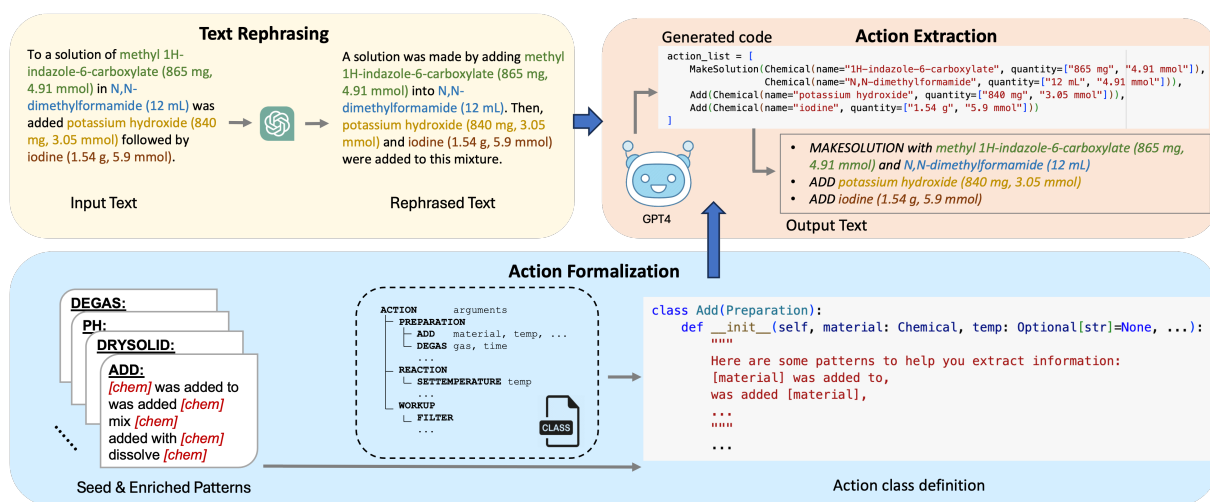


Figure 2: Overview of ActionIE.

163 niques to perform named entity recognition on
 164 chemical literature. These methods were fast and
 165 effective at extracting key information, but had lim-
 166 ited capabilities at handling more complex sentence
 167 structures in patent documents. In recent years,
 168 the transformer structure has also been introduced
 169 to action extraction. Paragraph2Actions (Vaucher
 170 et al., 2020b) used a transformer-based encoder-
 171 decoder architecture trained on human-annotated
 172 data to generate action sequences.

173 More recent advancements in NLP are led by
 174 pretrained LLMs. Wang et al. (2022b) finetuned
 175 a BERT model to perform named entity recog-
 176 nition on materials and extract synthesis actions
 177 on a dataset of solution-based inorganic materials
 178 synthesis. Zeng et al. (2023) finetuned both a T5
 179 model and a GPT-3.5 model on a human-annotated
 180 dataset. While these transformer-based models ex-
 181 cel in capturing the semantics of diverse scientific
 182 language, they rely on human-annotated datasets,
 183 which are created under extensive labor from do-
 184 main experts, and are prone to human errors. Also,
 185 these methods hard-code structure definitions ins-
 186 ide their framework, and have to infer structure
 187 semantics based on the training data, which could
 188 lead to inaccuracies if the training data is not rep-
 189 resentative enough.

190 2.2 Leveraging Programming Languages for 191 Structure Extraction Tasks

192 With the overwhelming success of very large
 193 decoder-only language models (such as GPT-3
 194 (Brown et al., 2020), GPT-3.5 and GPT-4 (Ope-
 195 nAI, 2023), PaLM 2 (Anil et al., 2023), Llama 2
 196 (Touvron et al., 2023), etc.) on a variety of NLP

197 tasks, recent research has increasingly focused on
 198 the application of LLMs for scientific structure ex-
 199 traction tasks. Agrawal et al. (2022) demonstrated
 200 the power of zero-shot learning on GPT-3 for ex-
 201 tracting information from clinical texts. Dunn et al.
 202 (2022) further performed chemical entities and re-
 203 lation extraction with a GPT-3 model finetuned on
 204 500 input-output pairs. Zhong et al. (2023b) uses
 205 GPT-4 to capture the roles of chemical entities in
 206 scientific text.

207 On the other hand, the large language models
 208 show noteworthy improvement in code generation.
 209 Codex (Tyers et al., 2023), finetuned from a GPT
 210 model, has shown remarkable abilities in code com-
 211 pletion. The recent year has seen the application of
 212 GPT-based agents (Hong et al., 2023; Zhou et al.,
 213 2023; Wang et al., 2024), which leverage the rea-
 214 soning and decision abilities of GPT models along
 215 with Chain-of-Thought approaches, in program-
 216 ming tasks.

217 Among these developments of structure extrac-
 218 tion and code generation, Code4Struct (Wang et al.,
 219 2022a) extracts structured event information from
 220 natural language using code generation. It aligns
 221 programming constructs, such as class definitions,
 222 inheritance, and functions with the entity and event
 223 types of interest, utilizing both the structural and
 224 semantic information of coding.

225 3 ActionIE Framework

226 3.1 Task Formulation

227 Given a text T , we aim to extract all procedures
 228 (actions) $P = \{(o_1, a_1), \dots, (o_n, a_n)\}$, $o_i \in S$ men-
 229 tioned in T in sequence, where S is a set of pre-

Module Name	Models
Pattern Mining	Flan-T5-Large
Text Rephrasing	GPT-4-0613
Code Generation	GPT-4-0613
Code to Natural Language	Pre-defined Rules

Table 1: Models used for each module in ActionIE.

defined operation types, and a_i is the pre-defined attributes of operation O_i . Note that rather than identifying the specific role a substitute plays within a reaction, our task focuses on the category of attribute to which it belongs. Following prior work (Vaucher et al., 2020a), we set the pre-defined operation types as follows: *Add*, *CollectLayer*, *Concentrate*, *Degas*, *DrySolid*, *DrySolution*, *Extract*, *Filter*, *MakeSolution*, *Microwave*, *Partition*, *PH*, *PhaseSeparation*, *Purify*, *Quench*, *Recrystallize*, *Reflux*, *SetTemperature*, *Sonicate*, *Stir*, *Triturate*, *Wait*, *Wash*, *Yield*, *FollowOtherProcedure*, *InvalidAction*, *OtherLanguage*, and *NoAction*. Definitions for each action are described in Appendix A.

3.2 Action Extraction with Programming Languages

Previous methods utilize a large amount of rules and patterns provided by human or train a model in a supervised way which require cost-sensitive labelled data. In addition, the definitions of actions may change based on the needs of scientists. Under certain circumstances, re-creating rules and patterns by human may be required for unsupervised methods; and relabelling data may be needed for supervised methods.

Driven by the aforementioned drawbacks, and with the emergence of Large Language Models (LLMs), we propose to use LLMs to tackle this action extraction task, as they have demonstrated promising capabilities in information extraction, particularly in data-scarce scenarios. Naively, one may directly input a paragraph along with all definitions of actions and ask LLMs to extract the action information. However, this approach poses some problems. The first is the well-known hallucination problem of the generation from LLMs (Huang et al., 2023). LLMs may generate actions that are not in the pre-defined action set since LLMs may directly output the verb found in the paragraph as an action or output an action not in the pre-defined set based on its summarization. Furthermore, LLMs may output detailed action sequences

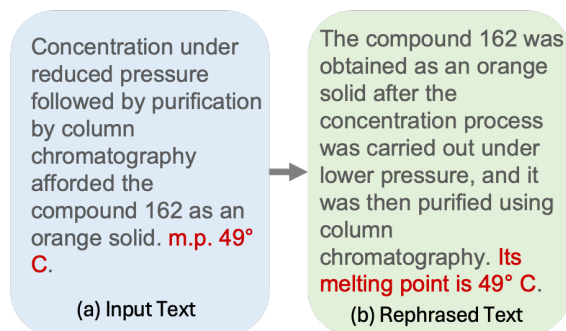


Figure 3: Rephrasing Example.

while it should summarize some of the actions. For instance, the ground-truth action for text “Add HCl to pH 5.” (adding HCl until the pH of the liquid is 5) is “PH with HCl to pH 5.”, while LLMs also include the “ADD” action which results in “ADD HCl; PH with HCl to pH5.” This demonstrates that LLMs fail to understand the relationship between actions.

In light of these limitations, we propose to reformulate the action extraction task as a code writing problem for LLMs that we transform each action type into a Python class. This has a few advantages. First, the abstract nature of class in programming languages and the relationship between classes including “Inheritance” and “Composition” relationships help LLMs better interpret the relationships between actions. Second, class variables in programming languages enable LLMs to understand what needs to be extracted for each action. Next, it is more suitable for an environment that needs changing the set of actions and the interested information for each action. The users can easily define the operations they want to extract and the attributes for each operation by simply modifying the Python class file which is fed to the LLMs. Finally, this minimizes the gap between natural language and robotics language as it is more convenient to transform the Python code produced by our method to the code that can be executed by robots. Figure 6 demonstrates the prompt we use for code generation.

3.3 External Information Guided Extraction

Text Rephrasing Scientific literature may have its own writing style that is different from ordinary writing, particularly in chemistry literature. We propose to first use LLMs to rephrase the given paragraph for two main reasons. First, rephrasing complex scientific texts into simpler language en-

hances their comprehensibility for large language models like GPT-4, which are pre-trained on general, non-scientific text sources. Second, it introduces domain knowledge encoded in the language model. This is exemplified in the case presented in Figure 3. The original paragraph (a) contains a phrase “m.p. 49 °C”, which is usually been misinterpreted as environment temperature. By leveraging LLMs for rephrasing, “m.p.” is rephrased as “melting point”, shown in Figure 3, and leads to a correct extraction. In practice, we prompt GPT-4 to rephrase the input text, as well as feeding in the mined patterns to keep the structure of the rephrased text as much as possible. Figure 5 shows the prompt we use for text rephrasing.

Pattern-guided Extraction For human, it is possible to identify certain action information, even lack of any prior domain knowledge. Consider an example “Partition between water (100 mL) and ethyl acetate (100 mL)” (Vaucher et al., 2020a). We can identify that “water (100 mL)” and “ethyl acetate (100 mL)” are the chemicals involved in the “Partition” action. This can be accomplished by the guidance of linguistic cues, including the semantics of phrases and the structure of sentences.

Motivated by this observation, we utilize frequent patterns in the text that indicate specific reaction actions as linguistic cues to guide LLMs to extract action information. Take “PH” action as an example, we first use a special token “[Chemical]” to replace all occurrences of the chemical with CHEMDATAEXTRACTOR (Swain and Cole, 2016). Several seed patterns are created, such as *pH [pH] with [Chemical]*. The red [pH] indicates a pH value, and the blue [Chemical] indicates the chemical for adjusting the pH. With a set of seed patterns for each action, we mine the enriched patterns through 1) labeling all occurrences in the corpus with seed patterns, 2) training a Flan-T5 model in a question-answering fashion, 3) re-labeling the corpus with the trained Flan-T5 model, and 4) selecting the most frequent patterns as the enriched patterns.

After merging enriched patterns with seed patterns as new seed patterns, we repeat the aforementioned process to mine more reliable patterns iteratively.

3.4 Extracted Action Evaluation

We observe that some actions are equivalent to each other, for instance, [MakeSolution] with A and B is

equivalent to [Add] A; [Add] B, and sometimes the order of actions does not matter. Previous evaluation metrics do not consider the order of actions nor the equivalence between actions, and penalize mismatches. In order to take the order of actions and their equivalences, we propose a graph-based metric called GRAPH MATCHING SIMILARITY. Given a sentence t with $n \in \mathbb{Z}$ actions a_1, a_2, \dots, a_n , and equivalent relations $f : A \rightarrow \{A\}$, where A is a set of actions and a_i is an arbitrary action, we first construct its corresponding graph G . Details can be found in Algorithm 1.

Algorithm 1 Algorithm for Action Graph Construction

Input: Sentence $t = (a_1, a_2, \dots, a_n)$ Equivalent Relations $f : A \rightarrow \{A\}$
Output: Graph $G = (V, E)$
procedure CONSTRUCTGRAPH(t, f)
 $V = \{a_1\}$
for $i \leftarrow 2$ **to** $n - 1$ **do**
 $V \cup \{a_i\}; E \cup \{(a_i, a_{i-1}), (a_{i+1}, a_i)\}$
if $a_i \in D(f)$ **then**
 $V \cup \{f(a_i)\}$
 $E \cup \{(f(a_i), a_{i-1}), (a_{i+1}, f(a_i))\}$
end if
end for
return $G = (V, E)$
end procedure

After constructing graphs for the ground truth sentence and the sentence to be evaluated, we compare the similarity of the two graphs with graph kernels, illustrated in Algorithm 2.

Algorithm 2 Algorithm for Evaluating Extracted Actions in Natural Language

Input: Ground Truth Sentence $t = (a_1, a_2, \dots, a_n)$, Sentence to be Evaluated $t' = (a'_1, a'_2, \dots, a'_m)$, Equivalent Relations f , Graph Kernel $k : G \times G \rightarrow \mathbb{R}$
Output: Similarity Score s
 $G \leftarrow \text{ConstructGraph}(t, f)$
 $G' \leftarrow \text{ConstructGraph}(t', f)$
 $s \leftarrow k(G, G')$

Graph kernels are widely used for evaluating the similarity between two graphs (Vishwanathan et al., 2010). The implementation is conducted with GraKeL (Siglidis et al., 2020). The evaluation with human judgements compared with other metrics

can be found at Section 4.3.

4 Experiments

4.1 Experimental Setup

Datasets We evaluate the effectiveness of our method on two datasets. One is the test set used in previous work (Vaucher et al., 2020a), which contains 352 texts related to an experimental procedure for chemical synthesis. We refer this dataset as **PATENTACTION** as all the paragraphs in it are from patent data. Dataset statistics can be found in Appendix B. The input is a paragraph from chemical literature which contains one or multiple actions. The output is a combination of pre-defined actions in natural language. This dataset is designed for evaluating action extraction in chemical literature setting based on chemist’s need.

Since extracting action information from scientific literature is of the same significance as from patent data, collaborating with chemists, we construct a dataset called **SCIENTIFIACTION**. 100 long paragraphs are collected from ChemRxiv (Cambridge Open Engage, 2023). The average length of paragraphs in ScientificAction is 770.77, while the average length in PatentAction is only 158.24. ScientificAction will be released upon the acceptance of the paper.

Baselines We compare ActionIE with several state-of-art methods: Paragraph2Actions (Vaucher et al., 2020a), ChemTrans (Zeng et al., 2023), GALACTICA-6.7b(Taylor et al., 2022), and GPT-4 (OpenAI, 2023).

Implementation Details We choose GPT-4-0613 (OpenAI, 2023) as the model for extraction, which supports up to 8,192 tokens. We use “google/flan-t5-large” (Raffel et al., 2020) for linguistic pattern extraction. GPT-4 (OpenAI, 2023) is accessed through OpenAI api. For the parameters of GPT-4 (OpenAI, 2023), we use sampling temperature $t = 0$, and set 500 as the maximum number of new tokens.

Evaluation Metrics for Natural Language Following previous work, we use BLEU score (Papineni et al., 2002) and Levenshtein Similarity (Levenshtein et al., 1966) to evaluate the quality of extracted actions in natural language. Following previous work, the BLEU score is modified since the original BLEU score does not consider short sentences which is common in the test data. The

proposed GRAPH MATCHING SIMILARITY is also used for evaluating in the natural language level.

Evaluation Metrics for Operation Level In order to verify the quality of the extracted action sequence in operation level, we use precision, recall, and F1 scores. The sets of operations in ground truth and output are compared, and the attributes are ignored. To better consider the order of operations, we employ SeqMatch-O (SM-O) proposed in Zeng et al. (2023), an evaluation metric for sequence matching in operation level. For details of SeqMatch-O, please refer to Zeng et al. (2023).

Evaluation Metrics for Attribute Level Following previous work, we leverage SeqMatch-A (SM-A) proposed in Zeng et al. (2023) for verifying the quality of attribute-level extraction. For each matched position in SeqMatch-O, the levenshtein similarity is calculated for each argument pair, and the average argument score is used rather than the original 1 in SM-O. Please refer to Zeng et al. (2023) for more details.

4.2 Experimental Results

Results for Extraction in Natural Language The first part of Table 2 represents the results of extraction in natural language in PatentAction dataset. ChemTrans cannot output natural language action sequences, hence, its scores are not calculated. Our proposed ACTIONIE significantly outperforms all baselines in levenshtein similarity, and outperforms all baselines in BLEU except Paragraph2Actions, but still get a very close score. GALACTICA-6.7 performs poorly as it is not designed for this task. GPT-4 demonstrates its promising performance given its comparable scores with just 10 demonstrations.

The second part of Table 2 represents the result of extraction in natural language in ScientificAction, a more complex and challenging dataset than PatentAction. Paragraph2Actions is trained on patent data and does not generalize well in scientific literature. Sometimes, Paragraph2Actions only outputs FollowOtherProcedure action and ignores other actions described in the input paragraph. Even GPT-4 receive higher scores in levenshtein similarity, demonstrating better generalization than Paragraph2Actions. Ablation study highlights the significance of using patterns as linguistic cues, in all cases, we gain much improvement by utilizing the patterns.

Models	BLEU	Levenshtein Similarity	Precision	Recall	F1	Graph Matching Similarity	SM-O	SM-A
<i>Results for PatentAction (Avg Length: 158.24)</i>								
Supervised Methods								
Paragraph2Actions	0.8511	0.8927	0.9017	0.9034	0.8985	0.8003	0.8893	0.8629
ChemTrans	-	-	0.5927	0.4325	0.4866	-	0.4401	-
Few-shot Methods (10-shot)								
Galactica-6.7b	0.0051	0.1336	0.3526	0.2705	0.2732	0.2921	0.1453	0.0534
GPT-4	0.4280	0.6822	0.7537	0.7758	0.7458	0.7923	0.7566	0.6633
ACTIONIE	0.8237	0.9018	0.9126	0.9198	0.9101	0.8136	0.8880	0.8521
- Patterns	0.6829	0.8070	0.8458	0.8220	0.8218	0.8074	0.8248	0.7583
<i>Results for ScientificAction (Avg Length: 770.77)</i>								
Supervised Methods								
Paragraph2Actions	0.4907	0.5380	0.8643	0.5933	0.6633	0.6391	0.5922	0.5118
ChemTrans	-	-	0.9212	0.4583	0.5982	-	0.4924	-
Few-shot Methods (10-shot)								
Galactica-6.7b	-	-	-	-	-	-	-	-
GPT-4	0.4571	0.6625	0.7858	0.7175	0.7312	0.7574	0.6670	0.5137
ACTIONIE	0.7808	0.8394	0.9236	0.8166	0.8584	0.8013	0.8277	0.7087
- Patterns	0.7193	0.8160	0.8942	0.8033	0.8444	0.7980	0.8099	0.6757

Table 2: Overall experimental results. ChemTrans does not support outputting natural language, only the operations are evaluated. Galactica-6.7b fails when the input is too long, therefore, the result is not reported.

<p>Input: The reaction of 1-(3,5-dichlorophenyl)-5-iodoimidazole 2o (0.3 mmol, 102 mg) with additional portion of Oxone® in H2SO4 (scaled down to 0.3 mmol) according to general procedure afforded 5,7-dichlorobenzo[d]imidazo[5,1-b][1,3]iodazol-4-ium hydrogen sulfate 3o as off-white solid, 110 mg, yield: 85%. mp = 205-208°C. 1H NMR (400 MHz, DMSO-d6) δ 9.18 (s, 1H), 8.45 (s, 1H), 7.80 (s, 1H), 7.56 (s, 1H). 13C NMR (100 MHz, DMSO-d6) δ 137.3, 137.1, 136.7, 133.9, 133.3, 126.8, 115.8, 99.5, 95.8. HRMS (ESI/Q-TOF, positive ionization): calcd for C9H4Cl2IN2 + (m/z: [M-HSO4] +): 336.8791, found: 336.8817.</p>
<p>Paragraph2Actions (Error: Only output FollowOtherProcedure action)</p> <ol style="list-style-type: none"> FollowOtherProcedure
<p>ChemTrans (Error: Only output YIELD action; Misclassify the product as reagent; Yield information only contains percentage)</p> <ol style="list-style-type: none"> YIELD Reagent: 1-(3,5-dichlorophenyl)-5-iodoimidazole 2o, Oxone®, H2SO4, 5,7-dichlorobenzo[d]imidazo[5,1-b][1,3]iodazol-4-ium hydrogen sulfate 3o YIELD Yield: 85% Reagent: 5,7-dichlorobenzo[d]imidazo[5,1-b][1,3]iodazol-4-ium hydrogen sulfate 3o
<p>GPT4 (Error: Contain an extra SetTemperature action (mp is melting point instead of environment temperature); The order of actions is incorrect; Missing 85% in YIELD action)</p> <ol style="list-style-type: none"> ADD 1-(3,5-dichlorophenyl)-5-iodoimidazole 2o (0.3 mmol, 102 mg) ADD Oxone® in H2SO4 (scaled down to 0.3 mmol) YIELD 5,7-dichlorobenzo[d]imidazo[5,1-b][1,3]iodazol-4-ium hydrogen sulfate 3o (off-white solid, 110 mg) SetTemperature 205-208 °C FollowOtherProcedure
<p>ActionIE (Perfect)</p> <ol style="list-style-type: none"> ADD 1-(3,5-dichlorophenyl)-5-iodoimidazole 2o (0.3 mmol, 102 mg) ADD Oxone® in H2SO4 (0.3 mmol) FollowOtherProcedure YIELD 5,7-dichlorobenzo[d]imidazo[5,1-b][1,3]iodazol-4-ium hydrogen sulfate 3o (110 mg, 85%)

Figure 4: Case Study.

Metric Name	Pearson	Spearman	Kendall's Tau
BLEU	0.1791	0.2427	0.2055
Levenshtein Similarity	0.1742	0.2603	0.2179
Graph Matching Similarity	0.3144	0.2976	0.3058

Table 3: Metric Correlations with Human Judgements.

Results for Operation-Level Extraction The middle columns of Table 2 represents the results of operational-level extraction. In the PatentAction dataset, ACTIONIE beats all baselines in precision, recall, and F1 scores, and have very close scores with Paragraph2Actions in SeqMatch-O (0.8880 vs 0.8893).

In the ScientificAction dataset, ACTIONIE outperforms all baselines. Both Paragraph2Actions and ChemTrans are trained on patent data, and achieve a high precision, but have a low recall and F1 scores.

As for the albatron study, ACTIONIE benefits significantly from the improvement provided by the patterns, which suggests that the patterns effectively help identify the actions.

Results for Attribute-Level Extraction As listed in the last column of Table 2, in PatentAction dataset, ACTIONIE outperforms all baselines except Paragraph2Actions, but still has a competitive score (0.8521 vs 0.8629).

In ScientificAction dataset, ACTIONIE surpasses all baselines by a substantial margin. Note that GPT-4 receives a slightly higher score than Paragraph2Actions, which further implies the limitation of supervised methods such as Paragraph2Actions and ChemTrans.

4.3 Evaluation Metric Analysis

To better understand how well our proposed GRAPH MATCHING SIMILARITY metric aligns with human evaluation, we randomly sample 100 outputs produced by Paragraph2Actions, GPT-4, and ACTIONIE, which are then given a score by chemists from 1 to 5. We calculate three correlation coefficients, Pearson, Spearman, and Kendall’s Tau. As the results shown in Table 3, the proposed GRAPH MATCHING SIMILARITY is better aligned with human judgements than BLEU and Levenshtein Similarity.

4.4 Case Study

We randomly sample an example from SCIENTIFIC-PATENT and study the output of different methods (see Figure 4). Paragraph2Actions only outputs FollowOtherProcedure action, and it has been noticed that it consistently does so whenever the input mentions another procedure. While the model is supervised to do so, this is an unwanted behavior since the output would ignore any other actions mentioned in the text. ChemTrans only captures the YIELD action, though it includes many details of the reagent. However, ChemTrans will fail if we are also interested in the melting point (mp) of the product given it is a supervised method. It also misclassifies the product as reagent. GPT-4 correctly extracts most of the actions and their attributes while missing the first ADD action, and the order of actions is wrong.

5 Conclusion and Future Work

In this paper, we propose ACTIONIE, a framework for extracting experimental action sequences from scientific literature. Our approach leverages the strength of LLMs by transforming the action extraction problem into a coding question for LLMs. Additionally, it incorporates text rephrasing and linguistic knowledge which further improve the overall performance. To more accurately evaluate the extraction quality, we introduce a graph-based metric, GRAPH MATCHING SIMILARITY. We have also developed a dataset, SCIENTIFICACTION, to offset the lack of scientific literature occurred in previous datasets. Experiments demonstrate that ACTIONIE outperforms state-of-the-art baselines and GRAPH MATCHING SIMILARITY is more aligned with human judgements than previous evaluation metrics. For future developments, one exciting yet challenging direction is to explore deeper into different aspects of the extraction process and integrating these parts into an automated workflow that transforms scientific papers into actionable experiments. This contains identifying relevant paragraphs from scientific papers that describe experimental procedures, creating a robotic system that runs the extracted chemical actions, and automated outcome validation.

Limitation

The limitations of this paper are stated as follows:

1. In our experiments, we use GPT-4 as the backbone model through OpenAI’s API. Although

567 ACTIONIE can be incorporated with other
568 causal language models, the performance may
569 change when using different language mod-
570 els. In addition, the performance might be
571 changed by the modification of GPT-4 since
572 its performance may be different over time
573 (OpenAI, 2023). Replacing the GPT-4 API
574 with a static large language model such as
575 Llama-2 (Touvron et al., 2023) could alleviate
576 this issue, but this may require considerable
577 computing resources, which are often limited.

- 578 2. Although the dataset proposed in this paper is
579 collected from scientific literature and is much
580 longer than previous datasets, it is still shorter
581 than a scientific paper. Extracting informa-
582 tion from a full paper may not be possible if
583 it is too long, given that current GPT-4 API
584 has token limits. Integrating a text segmenta-
585 tion module may be one direction to solve this
586 problem. Another direction may be deploy-
587 ing techniques that reduce the token limits
588 (Bertsch et al., 2023).

589 References

590 Monica Agrawal, Stefan Hegselmann, Hunter Lang,
591 Yoon Kim, and David Sontag. 2022. [Large language
592 models are few-shot clinical information extractors](#).
593 In *Proceedings of the 2022 Conference on Empirical
594 Methods in Natural Language Processing*, pages
595 1998–2022, Abu Dhabi, United Arab Emirates. Asso-
596 ciation for Computational Linguistics.

597 Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin John-
598 son, Dmitry Lepikhin, Alexandre Passos, Siamak
599 Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng
600 Chen, Eric Chu, Jonathan H. Clark, Laurent El
601 Shafey, Yanping Huang, Kathy Meier-Hellstern, Gau-
602 rav Mishra, Erica Moreira, Mark Omernick, Kevin
603 Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao,
604 Yuanzhong Xu, Yujing Zhang, Gustavo Hernandez
605 Abrego, Junwhan Ahn, Jacob Austin, Paul Barham,
606 Jan Botha, James Bradbury, Siddhartha Brahma,
607 Kevin Brooks, Michele Catasta, Yong Cheng, Colin
608 Cherry, Christopher A. Choquette-Choo, Aakanksha
609 Chowdhery, Clément Crepy, Shachi Dave, Mostafa
610 Dehghani, Sunipa Dev, Jacob Devlin, Mark Díaz,
611 Nan Du, Ethan Dyer, Vlad Feinberg, Fangxiaoyu
612 Feng, Vlad Fienber, Markus Freitag, Xavier Gar-
613 cia, Sebastian Gehrmann, Lucas Gonzalez, Guy Gur-
614 Ari, Steven Hand, Hadi Hashemi, Le Hou, Joshua
615 Howland, Andrea Hu, Jeffrey Hui, Jeremy Hur-
616 witz, Michael Isard, Abe Ittycheriah, Matthew Jagiel-
617 ski, Wenhao Jia, Kathleen Kenealy, Maxim Krikun,
618 Sneha Kudugunta, Chang Lan, Katherine Lee, Ben-
619 jamin Lee, Eric Li, Music Li, Wei Li, YaGuang Li,
620 Jian Li, Hyeontaek Lim, Hanzhao Lin, Zhongtao Liu,

Frederick Liu, Marcello Maggioni, Aroma Mahendru,
Joshua Maynez, Vedant Misra, Maysam Moussalem,
Zachary Nado, John Nham, Eric Ni, Andrew Nys-
trom, Alicia Parrish, Marie Pellat, Martin Polacek,
Alex Polozov, Reiner Pope, Siyuan Qiao, Emily Reif,
Bryan Richter, Parker Riley, Alex Castro Ros, Au-
rko Roy, Brennan Saeta, Rajkumar Samuel, Renee
Shelby, Ambrose Slone, Daniel Smilkov, David R.
So, Daniel Sohn, Simon Tokumine, Dasha Valter,
Vijay Vasudevan, Kiran Vodrahalli, Xuezhi Wang,
Pidong Wang, Zirui Wang, Tao Wang, John Wiet-
ing, Yuhuai Wu, Kelvin Xu, Yunhan Xu, Linting
Xue, Pengcheng Yin, Jiahui Yu, Qiao Zhang, Steven
Zheng, Ce Zheng, Weikang Zhou, Denny Zhou, Slav
Petrov, and Yonghui Wu. 2023. [Palm 2 technical
report](#).

Amanda Bertsch, Uri Alon, Graham Neubig, and
Matthew R Gormley. 2023. [Unlimiformer: Long-
range transformers with unlimited length input](#).
arXiv preprint arXiv:2305.01625.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
Neelakantan, Pranav Shyam, Girish Sastry, Amanda
Askell, Sandhini Agarwal, Ariel Herbert-Voss,
Gretchen Krueger, Tom Henighan, Rewon Child,
Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
Clemens Winter, Christopher Hesse, Mark Chen, Eric
Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,
Jack Clark, Christopher Berner, Sam McCandlish,
Alec Radford, Ilya Sutskever, and Dario Amodei.
2020. Language models are few-shot learners. In
*Proceedings of the 34th International Conference on
Neural Information Processing Systems, NIPS’20*,
Red Hook, NY, USA. Curran Associates Inc.

Cambridge Open Engage. 2023. [Chemrxiv: The
preprint server for chemistry](#). <https://chemrxiv.org>.
Accessed: 2023-12-15.

Chemical Abstracts Service (CAS). 2023. [Scifinder](#).
<https://scifinder.cas.org>. Accessed: 2024-01-
15.

Ian W. Davies. 2019. The digitization of organic syn-
thesis. *Nature*, 570(7760):175–181.

Alexander Dunn, John Dagdelen, Nicholas Walker,
Sanghoon Lee, Andrew S. Rosen, Gerbrand Ceder,
Kristin Persson, and Anubhav Jain. 2022. [Structured
information extraction from complex scientific text
with fine-tuned large language models](#).

Elsevier B.V. 2023. [Reaxys](#). <https://www.reaxys.com>.
Accessed: 2024-01-15.

Jiang Guo, A. Santiago Ibanez-Lopez, Hanyu Gao, Vic-
tor Quach, Connor W. Coley, Klavs F. Jensen, and
Regina Barzilay. 2022. [Automated chemical reaction
extraction from scientific literature](#). *Journal of Chem-
ical Information and Modeling*, 62(9):2035–2045.
PMID: 34115937.

Lezan Hawizy, David M. Jessop, Nico Adams, and Pe-
ter Murray-Rust. 2011. [Chemicaltagger: A tool for](#)

678	semantic text-mining in chemistry . <i>Journal of Cheminformatics</i> , 3(1):17.	733
679		734
680	Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Ceyao Zhang, Jinlin Wang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. 2023. Metagpt: Meta programming for a multi-agent collaborative framework .	735
681		736
682		737
683		738
684		739
685		740
686	Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2023. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions .	741
687		742
688		743
689		744
690		745
691		746
692	Tuan M. Lai, Chengxiang Zhai, and Heng Ji. 2023. Knowledge-enhanced biomedical language models. In <i>Journal of Biomedical Informatics</i> .	747
693		748
694		749
695	Vladimir I Levenshtein et al. 1966. Binary codes capable of correcting deletions, insertions, and reversals. In <i>Soviet physics doklady</i> , volume 10, pages 707–710. Soviet Union.	750
696		751
697		752
698		753
699	NextMove Software. 2023. Pistachio: Structure to name. https://www.nextmovesoftware.com/pistachio.html . Accessed: 2024-01-15.	754
700		755
701		756
702	OpenAI. 2023. Gpt-4 technical report .	757
703		758
704	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation . In <i>Proceedings of the 40th Annual Meeting on Association for Computational Linguistics</i> , ACL '02, page 311–318, USA. Association for Computational Linguistics.	759
705		760
706		761
707		762
708		763
709	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer .	764
710		765
711		766
712		767
713		768
714	Giannis Siglidis, Giannis Nikolentzos, Stratis Limnios, Christos Giatsidis, Konstantinos Skianis, and Michalis Vazirgiannis. 2020. Grakel: A graph kernel library in python. <i>Journal of Machine Learning Research</i> , 21(54):1–5.	769
715		770
716		771
717		772
718		773
719	Amir Soleimani, Vassilina Nikoulina, Benoit Favre, and Salah Ait Mokhtar. 2022. Zero-shot aspect-based scientific document summarization using self-supervised pre-training . In <i>Proceedings of the 21st Workshop on Biomedical Language Processing</i> , pages 49–62, Dublin, Ireland. Association for Computational Linguistics.	774
720		775
721		776
722		777
723		778
724		779
725		780
726	Yu Song, Santiago Miret, and Bang Liu. 2023. MatSci-NLP: Evaluating scientific language models on materials science language tasks using text-to-schema modeling . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3621–3639, Toronto, Canada. Association for Computational Linguistics.	781
727		782
728		783
729		784
730		785
731		786
732		787
	Matthew C. Swain and Jacqueline M. Cole. 2016. Chemdataextractor: A toolkit for automated extraction of chemical information from the scientific literature . <i>Journal of Chemical Information and Modeling</i> , 56(10):1894–1904. PMID: 27669338.	788
		789
		790
		791
	Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. 2022. Galactica: A large language model for science .	792
		793
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esioibu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models .	794
		795
	Francis Tyers, Robert Pugh, and Valery Berthoud F. 2023. Codex to corpus: Exploring annotation and processing for an open and extensible machine-readable edition of the florentine codex . In <i>Proceedings of the Workshop on Natural Language Processing for Indigenous Languages of the Americas (AmericasNLP)</i> , pages 19–29, Toronto, Canada. Association for Computational Linguistics.	796
		797
	Alain C. Vaucher, Philippe Schwaller, Joppe Geluykens, Vishnu H. Nair, Anna Iuliano, and Teodoro Laino. 2021. Inferring experimental procedures from text-based representations of chemical reactions . <i>Nature Communications</i> , 12(1):2573.	798
		799
	Alain C. Vaucher, Federico Zipoli, Joppe Geluykens, Vishnu H. Nair, Philippe Schwaller, and Teodoro Laino. 2020a. Automated extraction of chemical synthesis actions from experimental procedures . <i>Nature Communications</i> , 11(1):3601.	800
		801
	Alain C. Vaucher, Federico Zipoli, Joppe Geluykens, Vishnu H. Nair, Philippe Schwaller, and Teodoro Laino. 2020b. Automated extraction of chemical synthesis actions from experimental procedures . <i>Nature Communications</i> , 11.	802
		803
	S.V.N. Vishwanathan, Nicol N. Schraudolph, Risi Kondor, and Karsten M. Borgwardt. 2010. Graph kernels . <i>Journal of Machine Learning Research</i> , 11(40):1201–1242.	804
		805

792 Sheng Wang, Yuzhi Guo, Yuhong Wang, Hongmao
793 Sun, and Junzhou Huang. 2019. [Smiles-bert: Large
794 scale unsupervised pre-training for molecular prop-
795 erty prediction](#). In *Proceedings of the 10th ACM
796 International Conference on Bioinformatics, Compu-
797 tational Biology and Health Informatics, BCB '19*,
798 page 429–436, New York, NY, USA. Association for
799 Computing Machinery.

800 Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang,
801 Yunzhu Li, Hao Peng, and Heng Ji. 2024. Executable
802 code actions elicit better llm agents. In *arxiv*.

803 Xingyao Wang, Sha Li, and Heng Ji. 2022a.
804 Code4struct: Code generation for few-shot structured
805 prediction from natural language. *arXiv preprint
806 arXiv:2210.12810*.

807 Zheren Wang, Olga Kononova, Kevin Cruse, Tanjin
808 He, Haoyan Huo, Yuxing Fei, Yan Zeng, Yingzhi
809 Sun, Zijian Cai, Wenhao Sun, and Gerbrand Ceder.
810 2022b. Dataset of solution-based inorganic materials
811 synthesis procedures extracted from the scientific
812 literature. *Scientific Data*, 9.

813 Zheni Zeng, Yi-Chen Nie, Ning Ding, Qian-Jun Ding,
814 Wei-Ting Ye, Cheng Yang, Maosong Sun, Weinan
815 E, Rong Zhu, and Zhiyuan Liu. 2023. [Transcrip-
816 tion between human-readable synthetic descriptions
817 and machine-executable instructions: an application
818 of the latest pre-training technology](#). *Chem. Sci.*,
819 14:9360–9373.

820 Ming Zhong, Siru Ouyang, Minhao Jiang, Vivian Hu,
821 Yizhu Jiao, Xuan Wang, and Jiawei Han. 2023a. [Re-
822 actIE: Enhancing chemical reaction extraction with
823 weak supervision](#). In *Findings of the Association for
824 Computational Linguistics: ACL 2023*, pages 12120–
825 12130, Toronto, Canada. Association for Computa-
826 tional Linguistics.

827 Ming Zhong, Siru Ouyang, Yizhu Jiao, Priyanka Kar-
828 gupta, Leo Luo, Yanzhen Shen, Bobby Zhou, Xianrui
829 Zhong, Xuan Liu, Hongxiang Li, Jinfeng Xiao, Min-
830 hao Jiang, Vivian Hu, Xuan Wang, Heng Ji, Martin
831 Burke, Huimin Zhao, and Jiawei Han. 2023b. [Reac-
832 tion miner: An integrated system for chemical reac-
833 tion extraction from textual data](#). In *Proceedings of
834 the 2023 Conference on Empirical Methods in Nat-
835 ural Language Processing: System Demonstrations*,
836 pages 389–402, Singapore. Association for Compu-
837 tational Linguistics.

838 Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman,
839 Haohan Wang, and Yu-Xiong Wang. 2023. [Lan-
840 guage agent tree search unifies reasoning acting and
841 planning in language models](#).

842 A Action Type

843 We adopt the same action types as in the previous
844 study, including 26 pre-defined action types. We
845 include the detailed descriptions of action types in
846 (Vaucher et al., 2020a) as a reference in Table 4 to
847 help readers better understand the action types.

B Dataset statistics

The number of each action type mentioned in all
352 samples are summarized in Table 5.

C Prompt

Figure 5 demonstrates the prompt for text rephras-
ing. Figure 6 represents the prompt for code gener-
ation.

You are an expert in chemistry.

Rephrase the paragraph if you think it is difficult for general readers to understand. Keep the structure of the text as much as possible. Use the provided patterns when it is possible.

Here is the paragraph: [\[Input Text\]](#)

Here is the patterns your output should utilize: [\[Enriched Patterns\]](#)

Figure 5: Prompt for Text Rephrasing.

You are an expert in chemistry, programming, and extracting information.

The following python script describes chemical reaction procedure actions.

```
```python
\[User Defined Python Class File\]
```
```

Extract chemical reaction procedure actions from the following text: [\[Input Text\]](#)

To clearly explain the task, we provide the following example:
[\[Demonstrations\]](#)

Following the above examples, complete the following code:

```
```python
procedure =
```
```

Remember to strictly follow the output format.

Figure 6: Prompt for Code Generation.

| Action Type | Description |
|-----------------------------|--|
| Add | Add a substance to the reactor |
| CollectLayer | Select aqueous or organic fraction(s) |
| Concentrate | Evaporate the solvent (rotavap) |
| Degas | Purge the reaction mixture with a gas |
| DrySolid | Dry a solid |
| DrySolution | Dry an organic solution with a desiccant |
| Extract | Transfer compound into a different solvent |
| Filter | Separate solid and liquid phases |
| MakeSolution | Mix several substances to generate a mixture or solution |
| Microwave | Heat the reaction mixture in a microwave apparatus |
| Partition | Add two immiscible solvents for subsequent phase separation |
| PH | Change the pH of the reaction mixture |
| PhaseSeparation | Separate the aqueous and organic phases |
| Purify | Purification |
| Quench | Stop reaction by adding a substance |
| Recrystallize | Recrystallize a solid from a solvent or mixture of solvents |
| Reflux | Reflux the reaction mixture |
| SetTemperature | Change the temperature of the reaction mixture |
| Sonicate | Agitate the solution with sound waves |
| Stir | Stir the reaction mixture for a specified duration |
| Triturate | Triturate the residue |
| Wait | Leave the reaction mixture to stand for a specified duration |
| Wash | Wash (after filtration, or with immiscible solvent) |
| Yield | Phony action, indicates the product of a reaction |
| FollowOtherProcedure | The text refers to a procedure described elsewhere |
| InvalidAction | Unknown or unsupported action |
| OtherLanguage | The text is not written in English |
| NoAction | The text does not correspond to an actual action |

Table 4: Pre-defined action types used in this paper.

| Action Type | Total number of occurrences |
|-----------------------------|------------------------------------|
| Add | 255 |
| CollectLayer | 37 |
| Concentrate | 54 |
| Degas | 1 |
| DrySolid | 12 |
| DrySolution | 22 |
| Extract | 34 |
| Filter | 34 |
| MakeSolution | 62 |
| Microwave | 0 |
| Partition | 5 |
| PH | 47 |
| PhaseSeparation | 4 |
| Purify | 24 |
| Quench | 8 |
| Recrystallize | 2 |
| Reflux | 7 |
| SetTemperature | 60 |
| Sonicate | 0 |
| Stir | 118 |
| Triturate | 3 |
| Wait | 19 |
| Wash | 45 |
| Yield | 37 |
| FollowOtherProcedure | 15 |
| InvalidAction | 11 |
| OtherLanguage | 2 |
| NoAction | 25 |

Table 5: Dataset statistics.