

An Autonomous Large Language Model Agent for Chemical Literature Data Mining

Kexin Chen^a, Hanqun Cao^a, Junyou Li^b, Yuyang Du^a, Menghao Guo^b, Xin Zeng^b, Lanqing Li^b,
Jiezhong Qiu^b, Pheng Ann Heng^a, Guangyong Chen^b

^a The Chinese University of Hong Kong, New Territories, Hong Kong SAR

^b Zhejiang Lab, Zhejiang University, Hangzhou, China

Abstract—Chemical synthesis, which is crucial for advancing material synthesis and drug discovery, impacts various sectors including environmental science and healthcare. The rise of technology in chemistry has generated extensive chemical data, challenging researchers to discern patterns and refine synthesis processes. Artificial intelligence (AI) helps by analyzing data to optimize synthesis and increase yields. However, AI faces challenges in processing literature data due to the unstructured format and diverse writing style of chemical literature. To overcome these difficulties, we introduce an end-to-end AI agent framework capable of high-fidelity extraction from extensive chemical literature. This AI agent employs large language models (LLMs) for prompt generation and iterative optimization. It functions as a chemistry assistant, automating data collection and analysis, thereby saving manpower and enhancing performance. Our framework’s efficacy is evaluated using accuracy, recall, and F1 score of reaction condition data, and we compared our method with human experts in terms of content correctness and time efficiency. The proposed approach marks a significant advancement in automating chemical literature extraction and demonstrates the potential for AI to revolutionize data management and utilization in chemistry.

Index Terms—Chemical synthesis, literature mining, intelligent agent, large language models

I. INTRODUCTION

The discipline of chemistry, characterized by its immense potential and practical utility, is deeply intertwined with the synthesis of materials and the discovery of drugs. These sectors, propelled by the investigation of novel materials, contribute to advancements in energy, environmental science, and nanotechnology. They also lay the groundwork for the discovery of new pharmaceuticals and analyses in the life sciences. Such progressions are instrumental in enhancing therapeutic methods for diseases and promoting the growth of the health sector. The technological evolution has led to the accumulation of a wealth of data regarding chemical reactions, most of which is freely accessible. However, the challenge lies in efficiently harnessing this data to unearth intricate patterns, and to facilitate the discovery of novel reaction mechanisms. Addressing this issue could expedite the synthesis of materials, and the development of drugs, thereby fostering innovation in the field of chemistry.

Artificial Intelligence (AI) has the potential to identify salient features and patterns of reactions by learning from extant data and predicting outcomes of new reactions [1], [2]. This capability can aid chemists in rapidly screening and

evaluating diverse reaction conditions, optimizing synthetic routes, and enhancing synthetic efficiency. Moreover, when AI is combined with algorithms for predicting and optimizing reactions, it can generate a variety of synthetic paths and optimize them based on specific objectives and constraints. This process assists chemists in rapidly identifying efficient, sustainable synthetic routes, thereby improving the yield and purity of synthetic products.

Despite the successes of AI in these areas, gaining a deeper understanding of reaction rules remains essential for analyzing chemical reactions and discovering valuable chemical reactions. Unearthing associations and patterns concealed in data, revealing common characteristics and mechanisms between different reactions, aids chemists in 1) better understanding the underlying principles of reactions and 2) guiding the design of novel reactions.

To accomplish this, the integration and knowledge management of data regarding chemical reactions are of utmost importance, as they form the basis for discovering new reaction rules. Through automated data collection, organization, and annotation, AI can establish a comprehensive database of chemical reactions, enabling chemists to conveniently access and utilize these data. This aids in enhancing the discoverability and reproducibility of data, allowing researchers to better utilize extant knowledge of chemical reactions.

However, contemporary AI technologies encounter some challenges when dealing with data from chemical reaction literature. The data lacks uniform organization and processing, and extracting core reaction information from intricate and lengthy literature is a challenging task. This necessitates AI models to possess advanced context analysis capabilities and high standards for pattern recognition in text style and content.

The introduction of large language models (LLMs) like Chat-GPT enables efficient communication between humans and machines by converting instructions into text and integrating solvers for multiple subtasks [3], [4]. This provides vast opportunities for literature mining and opens new avenues for AI exploration in the field of chemistry. Traditional methods for extracting information from chemical literature include manual extraction, rule-based extraction, and model-based extraction. Manual extraction, however, relies on the intensive labor of chemists, which escalates the cost of the dataset and limits its size. Rule-based methods struggle to adapt to new chemical literature with different writing styles.

Revising old rules and designing new ones is also a labor-intensive task. Deep learning frameworks, such as the one proposed by Guo et al. [5], are hampered by the scarcity of annotated reactions, subsequently diminishing the model’s performance. The carefully designed prompts for information extraction make it difficult to adapt to the task of reaction yield extraction, especially when confronted with new challenges such as the coreference problem.

In response to these challenges, we propose an end-to-end framework based on a powerful AI agent that automatically extracts high-fidelity chemical data from the vast amount of literature, which is shown in Figure 3. The AI agent, through its automatic perception of the artificial environment and reasoning decision-making as a basis, has achieved efficient utilization of large language models, greatly saving manpower and improving model performance. By employing this technique, our approach achieves domain task-specific decision-making and tool-employing. Compared to traditional artificial prompt tuning, we develop a novel multi-task literature-mining scheme through several novel techniques. To enhance the language interaction between different LLMs, we employ Chat-GPT for highly efficient prompt building. We set the literature database as the interaction environment, boosting the automatic prompt refinement by pre-defined criteria. We link the prompt execution as a GPT-based function, and the text-perception agent enhances the mining tasks by iteratively optimizing the self-driven instructions. Moreover, we build the literature mining pipeline as an executable chemistry assistant. The well-constructed API serves as a downstream function of the agent, allowing for convenient and efficient calls.

To assess the effectiveness of content extraction, we propose a novel evaluation system based on accuracy, recall, and F1 score for valuable reaction-related information. The evaluation system not only provides quantitative rules for the effectiveness of AI agent, but also serves as the directional optimization criteria for its decision-making. Furthermore, to assess the permissibility of the self-driven AI agent, we compare the literature mining performance to human chemistry experts according to the correctness of the content as well as the time consumption. Finally, we ablate our framework according to the proposed techniques to further prove the validity of our design [5]–[10]. Our main contributions are summarised as follows:

- We propose a novel approach to create AI agents in chemistry literature information extraction. For the first time, we have linked the concept of AI agents with AI-based chemical research. The agent-based framework for extracting chemical literature has greatly saved manpower and achieved intelligent task automation.
- We design a novel assessment scheme for evaluating the agent’s intelligence in terms of literature text mining through accuracy, recall, and F1 score. This evaluation scheme is one of the most important links to ensure the efficiency of the agent’s task execution. Our task-oriented application of chemical professional knowledge can bring a more intuitive performance display of the agent.

II. RESULTS AND DISCUSSION

A. General Results

Acting as an efficient helper for chemists, our agent is expected to obtain higher reaction information retrieval quality and less time cost. Thus, quantitatively measuring the performance of AI-aided approaches, as well as effectively comparing its abilities with human experts, are necessary for exploring this new field. To investigate the efficacy of our framework, we devise a novel and comprehensive pipeline for evaluating the proficiency of GPT-based literature mining methods.

In our evaluation process, we place great emphasis on assessing the quality of reactants/reagents, solvents, products, and yields involved in each Suzuki reaction, which is the primary goal of effective retrieval.

To ensure precise quantification of each component, we have employed an evaluation scheme that includes precision, recall, and F1-score. By utilizing these metrics, we can gauge the model’s ability to accurately extract pertinent reaction information and conduct exhaustive searches of factors related to reactions.

Table I
THE PRECISION, RECALL, F1-SCORE RESULTS OF DATA MINING FOR YIELD, REACTANT/REAGENT, SOLVENT, PRODUCT.

	Precision	Recall	F1-score
Yield	92.19%	78.53%	84.81%
Reactant / Reagent	89.04%	76.00%	82.00%
Solvent	91.90%	75.77%	83.06%
Product	87.45%	78.22%	82.58%

Table II
THE QUANTITIES OF CORRECT, EXTRACTED, AND TOTAL PIECES OF REACTION INFORMATION.

	Correct data	Extracted data	Total data
Yield	236	256	326
Reactant / Reagent	203	228	300
Solvent	227	247	326
Product	223	255	326

After obtaining the generation results from ChatGPT, we restore the outcomes for later comparison with the ground truth collected by human experts. We annotate 17 literature and 326 reactions to validate the effectiveness of our agent. On average, the precision, recall, and F1-score reached 90.14%, 77.13%, and 83.11%, respectively. A detailed representation of the precision, recall, and F1-score results can be found in Table 1. We also provide the quantities of correct, extracted, and total pieces of reaction information in Table 2.

B. Comparison to Human Experts

As far as we know, there are currently no other open-source tools available for extracting chemical reaction data from academic journals. Thus, this paper primarily validates the effectiveness and performance of the agent in extracting chemical reaction data information through comparison with

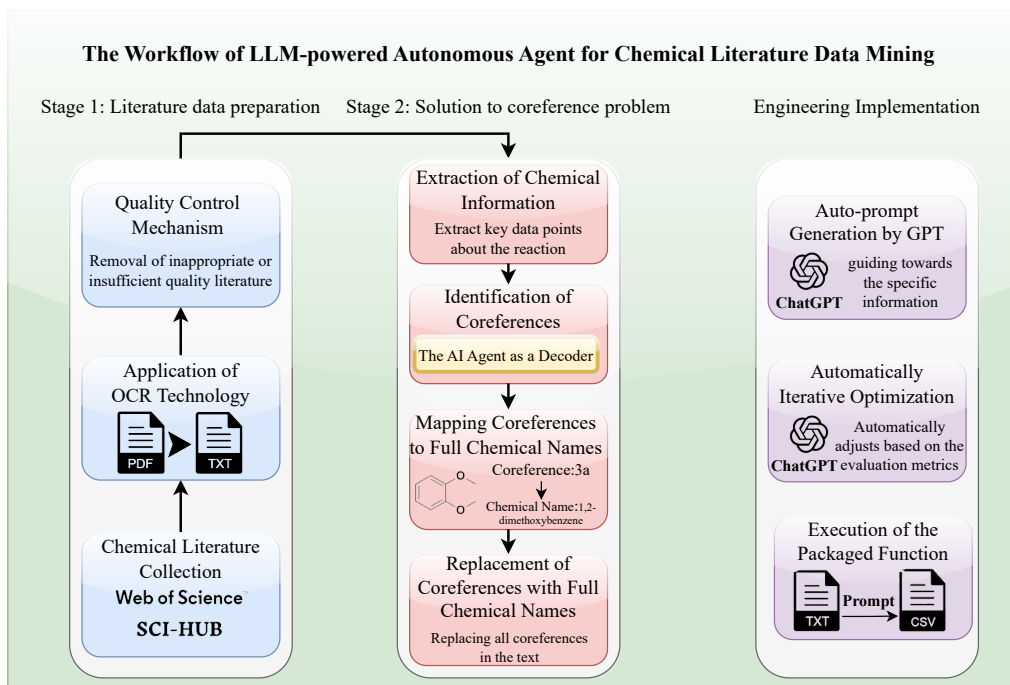


Figure 1. The framework of chemical literature analysis and reaction information extraction agent based on LLMs.

manual data collection by human chemists. The primary indicators for evaluation are precision, average cost, and average speed. To minimize uncertainty and randomness for human chemists, the comparative study selected ten graduate students specializing in chemistry at the master’s or doctoral level to perform manual data collection. The results from these ten chemistry professionals were then averaged to provide a comprehensive comparison with our agent. From Table 3, we can see that our AI agent reached competitive precision performance and much better performance in average cost and average speed.

III. METHOD

A. Data Preparation

Literature Collection: In the initial stage of our research, the acquisition of a high-quality literature dataset was paramount. To this end, we embarked on an extensive data collection process, leveraging the vast repository of chemical literature available on SciHub. Our specific focus was on articles about Suzuki reactions, a popular topic in organic chemistry. This search provides a substantial foundation for our subsequent analysis.

Application of OCR Technology: To convert these articles into a machine-readable format, we employed Optical Character Recognition (OCR) technology. This enabled us to transform the PDFs into text, thus making them amenable to further computational processing. However, it is noteworthy that the OCR process, while generally accurate, is not infallible. It occasionally struggles with complex layouts and low-quality scans, which can lead to errors in the resulting text.

Quality Control Mechanism: Given this potential for error, we instituted a quality control mechanism to ensure the reliability of our dataset. Specifically, we analyzed each article for the presence of key phrases such as "General Procedure", "Typical Procedure", or "General Experiment". These phrases are often indicative of the detailed methodological sections that are crucial for our purposes. Articles that did not contain these phrases were deemed to be of insufficient quality and were thus excluded from our dataset. Similarly, we also excluded articles where these phrases appeared more than five times, as this was suggestive of an overly complex or convoluted methodology that may not lend itself well to our extraction process.

Final Dataset: After this rigorous quality control process, we were left with a refined dataset of 1000 articles. These articles formed the basis of our subsequent extraction and performance measurement activities. Our dataset, while significantly reduced in size, was of high quality and well-suited to our research aims.

B. General Procedure

In our quest to extract chemical reaction conditions from the literature, we employed an AI Agent. This AI Agent, akin to a diligent chemist, navigates through the labyrinth of chemical literature, performing several tasks simultaneously to extract the desired information.

Task 1 (Extraction of Chemical Information): The final task that our AI Agent undertakes is the extraction of chemical information from the now standardized text. This is akin to a chemist analyzing the detailed notes of an experiment to extract key data points about the reaction.

Table III
THE PRECISION, AVERAGE COST, AND AVERAGE SPEED OF MANUAL DATA COLLECTION AND OUR AGENT.

	Precision	Average cost (USD)	Average speed (second)
Manual data collection	90%	1.41	288
AI agent	87%	0.0025	0.43

In this task, the AI Agent is required to extract information about the yield, reactant, catalyst, solvent, and product from each reaction. To achieve this, the AI Agent employs a multi-task framework and in-context learning technique.

From an engineering perspective, the AI Agent first identifies the sections of the text that describe the reaction conditions. It does this by searching for keywords and phrases that are commonly used in the chemical literature to describe these conditions. Once these sections are identified, the AI Agent then applies a series of extraction algorithms to pull out the required information. Figure 2 shows the prompt, example input, and example output of the in-context learning process, one of the important techniques for the agent.

Prompt:
Answer the question as truthfully as possible using the provided context.

Please summarize the following details with units in a json: yield(include%), reactant/reactent(s), solvent(s), product(s). Please note that the content usually includes a general procedure, followed by the specific description of the reaction. The general procedure provides the overall context, and the specific descriptions of each reaction offers unique details.

Please also note that the keys of a json object should be "yield", "reactant/reactent", "solvent", "product". If there do not exist such information, please tell me "NaN".

Example input:

General Procedure for the Preparation of Products. [Ni₂(dPr₂Im)₄(μ-COD)] (0.1 mmol, 83 mg), CsF (2 mmol, 304mg), Ar-Bneop (2 mmol), fluoroarene, and toluene (10 mL) were added to a Schlenk tube equipped with a magnetic stirring bar. The reaction mixture was heated at 100 °C for 18 h, and after that H₂O (5 mL) was added. The product was extracted with EtOAc (3 × 20 mL), and then the combined organic layers were dried over Na₂SO₄ and filtered, and the volatiles were removed in vacuo. The product was purified by column chromatography on silica gel using hexane as the eluent. The solvent of the product-containing fraction of the eluent was evaporated in vacuo. The yields provided are based on Ar-Bneop. Spectroscopic Data of the Products. 2,3,4,5,6-Pentafluoro-1,1'-biphenyl (3aa). Following the general procedure, a white solid in 72% yield (351 mg) was obtained from C6F6 (4 mmol, 462 μL) and C6H5-Bneop (2 mmol, 380 mg). 2,3,4,5,6-Pentafluoro-4'-methyl-1,1'-biphenyl (3ab). Following the general procedure, a white solid in 76% yield (390 mg) was obtained from C6F6 (4 mmol, 462 μL) and 4-CH₃-C6H₄-Bneop (2 mmol, 408 mg).

Example output:

yield	reactant/reactent	solvent	product
72% (351 mg)	C6F6(4 mmol, 462 μL),C6H5-Bneop(2 mmol, 380 mg),fluoroarene,[Ni ₂ (dPr ₂ Im) ₄ (μ-COD)](0.1 mmol, 83 mg),CsF(2 mmol, 304 mg)	toluene(10 mL)	2,3,4,5,6-Pentafluoro-1,1'-biphenyl
76% (390 mg)	C6F6(4 mmol, 462 μL),4-CH ₃ -C6H ₄ -Bneop(2 mmol, 408 mg),fluoroarene,[Ni ₂ (dPr ₂ Im) ₄ (μ-COD)](0.1 mmol, 83 mg),CsF(2 mmol, 304 mg)	toluene(10 mL)	2,3,4,5,6-Pentafluoro-4'-methyl-1,1'-biphenyl

Figure 2. The prompt, example input, example output in chemical reaction information extraction, demonstrating the in-context learning technique.

The result of this task is a structured dataset containing the yield, reactant, catalyst, solvent, and product information for each reaction described in the text. This dataset serves as the final output of our AI Agent, representing the culmination of its diligent and meticulous work, much like the final report of a chemist after a series of experiments.

Task 2 (Identification of Coreferences: The AI Agent as a Decoder): The task of identifying coreferences is the primary and most crucial step in our AI Agent's operation. In chemical literature, coreferences are typically denoted by a combination

of a number and a letter, serving as abbreviations for longer, more complex chemical names. This form of shorthand, while effective for human readers familiar with the context, presents a unique challenge for machine-reading and understanding.

To tackle this challenge, our AI Agent is equipped with the capability to recognize these specific patterns within the text. This process is not merely a superficial scan of the document; rather, it involves a deep, context-aware analysis of the text. The AI Agent uses the capabilities of GPT, which is designed to understand the context within the investigated text, to identify these coreferences accurately. It makes use of GPT's transformer-based architecture, which allows it to understand the dependencies between words in a sentence and across sentences. Figure 3 shows the prompt, example input, and example output of the in-context learning process in coreference extraction.

Prompt:
I am providing a paragraph from a piece of chemical literature. I would like you to help me identify instances of coreference, where a full chemical name is immediately followed by a shorthand label or alias.

Please provide the coreference in json format. Pay attention to direct aliases that come immediately after the chemical names. If there do not exist such coreference, please tell me "No coreference". Please check carefully about the full chemical name and shorthand label.

Example input:

Tetraethyl (E)-8,9-Bis(Z)-3-ethoxy-3-oxo-2-phenylprop-1-en-1-yl)hexadeca-1,8,15-triene-6,6,11,11-tetracarboxylate, 7c. It was obtained from 3n (25 mg, 0.06 mmol) following the general procedure for cycloisomerization reactions with Cp^{*}RuCl(cod) and purified by flash column chromatography (Hexane/AcOEt, 19:1). Colorless oil (22 mg, 0.03 mmol, 86%).

Example output:

Coreference	Full chemical name
7c	Tetraethyl (E)-8,9-Bis(Z)-3-ethoxy-3-oxo-2-phenylprop-1-en-1-yl)hexadeca-1,8,15-triene-6,6,11,11-tetracarboxylate

Figure 3. The prompt, example input, example output in the identification of coreferences, demonstrating the in-context learning technique.

Upon encountering a potential coreference, the AI Agent validates it against the patterns typically used for coreferences in the chemical literature. This step ensures that the identified coreferences are not false positives, such as a number and a letter appearing together by coincidence in the text.

Once a coreference is validated, the AI Agent records it for the subsequent steps. This record-keeping is meticulous and organized, ensuring that each coreference is accurately linked with its position in the text. This step is crucial as it sets the stage for the mapping of coreferences to their full chemical

names in the following task.

In summary, the task of identifying coreferences is a complex process that requires the AI Agent to combine pattern recognition with deep, context-aware text analysis. The accuracy and efficiency of this task are critical to the success of the subsequent steps in the AI Agent's operation. Through this task, the AI Agent demonstrates its ability to navigate and understand the complexities of chemical literature, setting the foundation for the extraction of chemical reaction conditions.

Task 3 (Mapping Coreferences to Full Chemical Names): The third task undertaken by our AI Agent is the mapping of coreferences to their corresponding full chemical names. This task is crucial, as it transforms the shorthand notations into their full forms, thereby enabling a more comprehensive understanding of the chemical reactions described in the text. In essence, this task serves as a bridge, linking the efficient but context-dependent coreferences with their context-independent full chemical names.

Building on the coreferences identified in Task 2, the AI Agent begins the process of mapping these coreferences to their full chemical names. To do this, the AI Agent makes use of GPT's context-understanding capabilities, scanning the text for instances where the coreference is defined, typically in proximity to its first mention.

The AI Agent is designed to handle the complexity of this task, which often involves navigating intricate sentence structures or piecing together information spread across multiple sentences. It employs advanced natural language processing techniques to understand the grammatical structure of the sentence, identify the subject and object, and distinguish between different clauses.

Once a full chemical name corresponding to a coreference is identified, the AI Agent meticulously records this mapping in a structured format. This data structure is designed for flexibility, allowing for updates if a more accurate or complete definition of the coreference is encountered later in the text.

In summary, Task 3 is a complex linguistic challenge that requires the AI Agent to act as a linguistic cartographer, drawing connections between different points of reference within the text. This task is vital for the subsequent steps in the AI Agent's operation, ensuring that the replacement of coreferences and the extraction of chemical information are based on accurate and complete data.

Task 4 (Replacement of Coreferences with Full Chemical Names): The fourth task is a critical juncture in the AI Agent's operation, where it begins to transform the raw text into a more analyzable form. This task involves replacing all instances of the identified coreferences in the text with their corresponding full chemical names.

This task is implemented by first creating a dictionary where the keys are the coreferences and the values are the corresponding full chemical names. The AI Agent then iterates over the text, and each time it encounters a coreference, it consults the dictionary for the corresponding full chemical name and replaces the coreference with it. This is accomplished using a string replacement function, which scans the text for the

coreference patterns and replaces them with the full chemical names.

This replacement process is akin to a search-and-replace operation in a text editor, but it is performed on a much larger scale and with a higher degree of complexity due to the intricacies of chemical nomenclature. The result is a text where the shorthand notations have been replaced with full chemical names, making the subsequent information extraction task more straightforward and accurate.

Engineering Implementation through GPT: From an engineering perspective, our approach to extracting chemical information from the literature is built upon a multi-task framework. This design was motivated by the complex and multi-faceted nature of the problem at hand. By breaking down the overall task into smaller, more manageable subtasks, we are able to tackle each aspect of the problem with specialized strategies, thereby improving the overall effectiveness and efficiency of our system.

The first stage in our process involves the use of GPT to generate prompts. GPT is a state-of-the-art language model that has been trained on a vast corpus of text from the internet. It has the ability to generate human-like text based on a given prompt. We leverage this capability to generate initial prompts for our tasks. These prompts act as the starting point for our AI Agent, guiding its focus toward the specific information we are interested in.

However, not all prompts are equally effective. Therefore, we employ an iterative optimization process to refine these prompts. In each iteration, the AI Agent evaluates the effectiveness of the current prompt in extracting the desired information from the text. This evaluation is based on a performance metric that we define, which could be accuracy, precision, recall, or any other metric that is relevant to the specific task. If the performance of the prompt is not satisfactory, the AI Agent modifies the prompt and evaluates its performance again. This process is repeated until we obtain a prompt that meets our performance criteria.

Once we have the optimized prompts, we map them to the API function that calls GPT. This mapping process involves translating the prompts into a format that the GPT API can understand. This is a crucial step as it ensures that our prompts are correctly interpreted by GPT, thereby maximizing the effectiveness of our extraction process.

The final step in our process is the execution of the packaged function. This function takes the text as input, applies the mapped prompts to extract the desired information, and returns this information as output. The function is packaged in a way that it can be easily integrated into other systems or workflows, making our AI Agent a versatile tool for chemical information extraction.

In summary, our engineering approach is a careful orchestration of several components, each designed with a specific purpose and all working together towards the common goal of effective and efficient chemical information extraction. By leveraging the power of GPT and the flexibility of our multi-task framework, we are able to tackle the complex task of

extracting chemical information from the literature, paving the way for new possibilities in the field of chemical informatics.

This AI Agent, through its systematic and rigorous approach, mirrors the meticulous nature of a chemist, making it an apt tool for the task at hand. By breaking down the complex task of chemical information extraction into manageable subtasks, we have created an AI Agent that is not only efficient but also scalable, paving the way for future advancements in the field.

IV. CONCLUSION

This research introduced an innovative AI agent that leverages LLMs to automate the extraction of high-fidelity chemical data from chemical literature. Our system has demonstrated superior performance in terms of accuracy, recall, and F1 score. The agent's ability to act as a chemistry assistant streamlines the data collection and analysis process, leading to significant savings in manpower and enhancements in performance. The agent's iterative optimization and prompt generation capabilities have shown to be particularly effective in dealing with the variegated and unstructured nature of literature. Additionally, the comparison with human experts has validated the AI agent's efficiency and correctness, showcasing its potential to revolutionize chemical data management and utilization.

Future research will focus on refining the AI agent's capabilities, expanding its applications, and integrating it with other advanced technologies to further elevate its performance and utility. The current work has laid a solid foundation for AI's role in chemical literature mining, promising to accelerate advancements in material synthesis, drug discovery, and a myriad of other areas within the field of chemistry.

REFERENCES

- [1] K. Chen, G. Chen, J. Li, Y. Huang, E. Wang, T. Hou, and P.-A. Heng, "MetaRF: attention-based random forest for reaction yield prediction with a few trails," *Journal of Cheminformatics*, vol. 15, no. 1, pp. 1–12, 2023.
- [2] K. Chen, J. Li, K. Wang, Y. Du, J. Yu, J. Lu, G. Chen, L. Li, J. Qiu, Q. Fang *et al.*, "Towards an automatic ai agent for reaction condition recommendation in chemical synthesis," *arXiv preprint arXiv:2311.10776*, 2023.
- [3] H. Cui, Y. Du, Q. Yang, Y. Shao, and S. C. Liew, "Llmind: Orchestrating ai and iot with llms for complex task execution," *arXiv preprint arXiv:2312.09007*, 2023.
- [4] Y. Du, S. C. Liew, K. Chen, and Y. Shao, "The power of large language models for wireless communication system development: A case study on fpga platforms," *arXiv preprint arXiv:2307.07319*, 2023.
- [5] J. Guo, A. S. Ibanez-Lopez, H. Gao, V. Quach, C. W. Coley, K. F. Jensen, and R. Barzilay, "Automated chemical reaction extraction from scientific literature," *Journal of chemical information and modeling*, vol. 62, no. 9, pp. 2035–2045, 2021.
- [6] D. T. Ahneman, J. G. Estrada, S. Lin, S. D. Dreher, and A. G. Doyle, "Predicting reaction performance in c–n cross-coupling using machine learning," *Science*, vol. 360, no. 6385, pp. 186–190, 2018.
- [7] D. Perera, J. W. Tucker, S. Brahmabhatt, C. J. Helal, A. Chong, W. Farrell, P. Richardson, and N. W. Sach, "A platform for automated nanomole-scale reaction screening and micromole-scale synthesis in flow," *Science*, vol. 359, no. 6374, pp. 429–434, 2018.
- [8] J. Schleinitz, M. Langevin, Y. Smail, B. Wehnert, L. Grimaud, and R. Vuilleumier, "Machine learning yield prediction from nicolite, a small-size literature data set of nickel catalyzed c–o couplings," *Journal of the American Chemical Society*, vol. 144, no. 32, pp. 14 722–14 730, 2022.

- [9] M. C. Swain and J. M. Cole, "Chemdataextractor: a toolkit for automated extraction of chemical information from the scientific literature," *Journal of chemical information and modeling*, vol. 56, no. 10, pp. 1894–1904, 2016.
- [10] Z. Zheng, O. Zhang, C. Borgs, J. T. Chayes, and O. M. Yaghi, "Chatgpt chemistry assistant for text mining and the prediction of mof synthesis," *Journal of the American Chemical Society*, vol. 145, no. 32, pp. 18 048–18 062, 2023, PMID: 37548379.