

# 1 **Fine-tuning ChatGPT Achieves State-of-the-Art Performance** 2 **for Chemical Text Mining**

3 Wei Zhang<sup>1,2,#</sup>, Qinggong Wang<sup>3,#</sup>, Xiangtai Kong<sup>1,2</sup>, Jiacheng Xiong<sup>1,2</sup>, Shengkun Ni<sup>1,2</sup>, Duanhua  
4 Cao<sup>1,4</sup>, Buying Niu<sup>1,2</sup>, Mingan Chen<sup>1,5,7</sup>, Runze Zhang<sup>1,2</sup>, Yitian Wang<sup>1,2</sup>, Lehan Zhang<sup>1,2</sup>, Xutong  
5 Li<sup>1,2</sup>, Zhaoping Xiong<sup>6</sup>, Qian Shi<sup>7</sup>, Feng Cheng<sup>8</sup>, Zunyun Fu<sup>1,\*</sup>, Mingyue Zheng<sup>1,2,3,\*</sup>

6 <sup>1</sup>Drug Discovery and Design Center, State Key Laboratory of Drug Research, Shanghai Institute of Materia  
7 Medica, Chinese Academy of Sciences, 555 Zuchongzhi Road, Shanghai 201203, China

8 <sup>2</sup>University of Chinese Academy of Sciences, No. 19A Yuquan Road, Beijing 100049, China

9 <sup>3</sup>Nanjing University of Chinese Medicine, 138 Xianlin Road, Nanjing 210023, China

10 <sup>4</sup>Innovation Institute for Artificial Intelligence in Medicine of Zhejiang University, College of  
11 Pharmaceutical Sciences, Zhejiang University, Hangzhou, Zhejiang 310058, China

12 <sup>5</sup>School of Physical Science and Technology, ShanghaiTech University, Shanghai 201210, China

13 <sup>6</sup>ProtonUnfold Technology Co., Ltd, Suzhou, China

14 <sup>7</sup>Lingang Laboratory, Shanghai 200031, China

15 <sup>8</sup>Department of Pharmaceutical Sciences, Taneja College of Pharmacy, University of South Florida, Tampa,  
16 Florida 33612, United States

17  
18 <sup>#</sup>Wei Zhang and Qinggong Wang contributed equally to this study.

19 <sup>\*</sup>Correspondence should be addressed to:

20 Mingyue Zheng: [myzheng@simm.ac.cn](mailto:myzheng@simm.ac.cn)

21 Zunyun Fu: [fuzunyun@simm.ac.cn](mailto:fuzunyun@simm.ac.cn)

## 22 **Abstract**

23           Extracting knowledge from complex and diverse chemical texts is a pivotal task for both  
24 experimental and computational chemists. The task is still considered to be extremely challenging  
25 due to the complexity of the chemical language and scientific literature. This study fine-tuned  
26 ChatGPT for five intricate chemical text mining tasks: compound entity recognition, reaction role  
27 labelling, metal-organic framework (MOF) synthesis information extraction, nuclear magnetic  
28 resonance spectroscopy (NMR) data extraction, and the conversion of reaction paragraph to action  
29 sequence. The fine-tuned ChatGPT demonstrated impressive performance, significantly reducing  
30 the need for repetitive and extensive prompt engineering experiments. It achieved exact accuracy  
31 levels ranging from 69% to 95% on these tasks with minimal annotated data. For comparison, we  
32 fine-tuned open-source pre-trained large language models (LLMs) such as Llama2, T5, and BART.  
33 The results showed that the fine-tuned ChatGPT excelled in all tasks. It even outperformed those  
34 task-adaptive pre-training and fine-tuning models that were based on a significantly larger amount  
35 of in-domain data. Given its versatility, robustness, and low-code capability, leveraging fine-tuned  
36 LLMs as toolkits for automated data acquisition could revolutionize chemical knowledge  
37 extraction.

38

## 39 **Main**

40           Chemical text mining is a crucial foundation in chemical research. It creates extensive  
41 databases that provide access to physicochemical properties and synthetic routes for experimental  
42 chemists. Additionally, it accumulates rich data and insights for computational chemists to use for  
43 modelling and predicting. More than just extracting information from chemical texts, the rule-  
44 based transformation of chemical text is particularly interesting. For instance, synthetic procedures

45 can be converted into action sequences<sup>1</sup> or programming languages<sup>2-4</sup>. This allows them to be  
46 understood and executed by robotic systems for automated syntheses.

47       Converting structured data from intricate scientific literature is a challenging task, especially  
48 due to the complexity and heterogeneity of chemical language. As a result, a number of text-mining  
49 tools have been developed. For instance, ChemDataExtractor<sup>5,6</sup> was created to extract chemical  
50 entities and their associated properties, measurements and relationships. ChemRxnExtractor<sup>7</sup> was  
51 designed to extract the product and label associated reaction roles such as reactant, catalyst, solvent,  
52 and temperature. Historically, the focus has been on designing models and algorithms specific to  
53 certain tasks, often using regular expressions with rule-based syntax or dictionary-matching. These  
54 tools require extensive domain knowledge and sophisticated data processing, which limits their  
55 versatility.

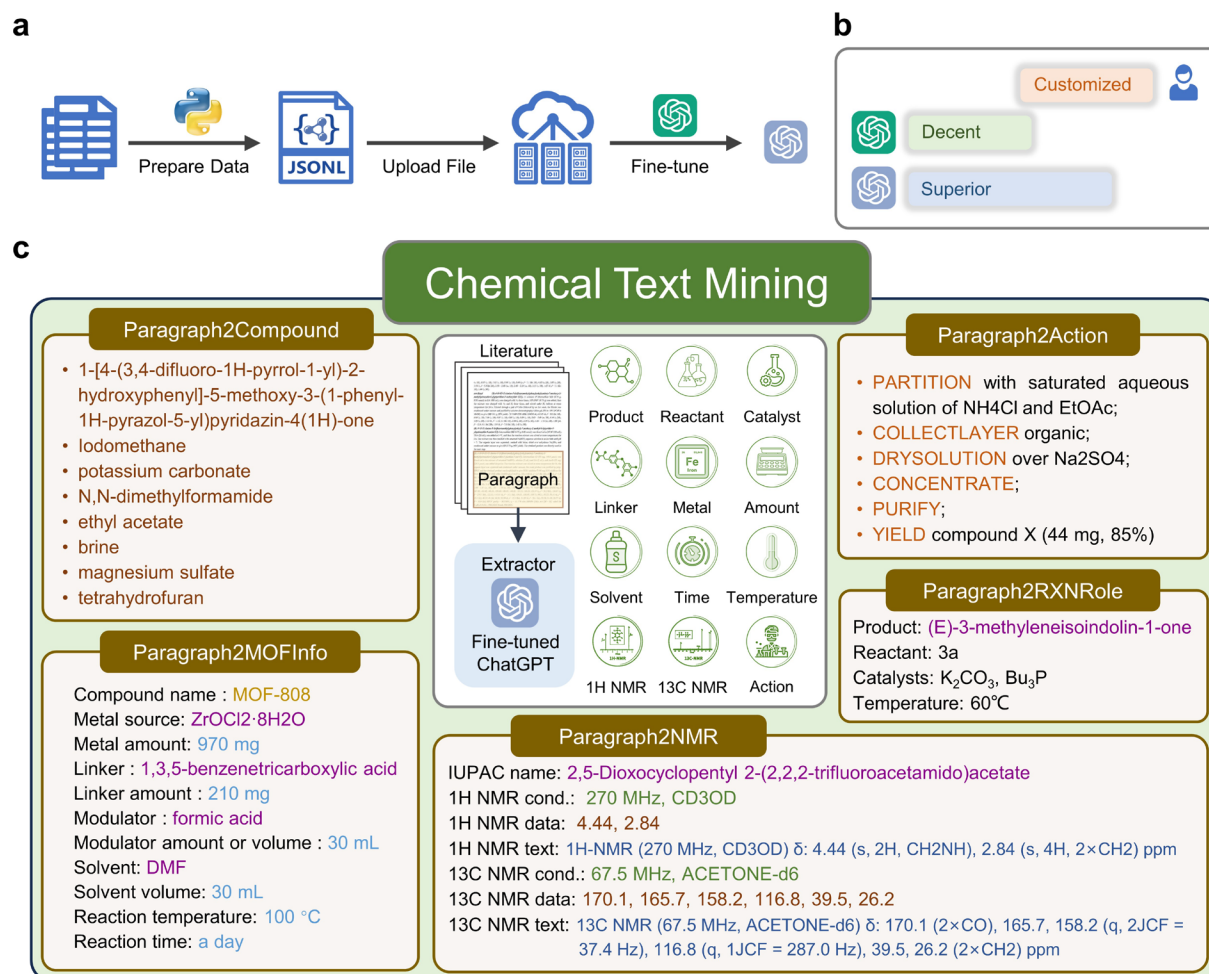
56       Recently, large language models (LLMs), represented by ChatGPT released in November  
57 2022, have shown the potential of Artificial General Intelligence (AGI). LLMs, such as GPT-3.5  
58 and GPT-4, can generate logical insights or content that meets requirements based on human  
59 instructions. We are entering a new era where AGI and medicinal chemists might work together.  
60 There have been assessments of ChatGPT's chemistry capabilities, including tasks like synonym  
61 transformation, property prediction, retrosynthesis, and molecule design<sup>8-10</sup>. However, LLMs tend  
62 to "hallucinate", meaning they generate unintended text that misaligns with established facts and  
63 real-world knowledge<sup>11,12</sup>. Moreover, objectively evaluating the results of open-ended questions  
64 remains a significant challenge.

65       At this juncture, LLMs may struggle to accurately answer factual and knowledge-based  
66 questions. However, using LLMs for knowledge extraction tasks should greatly alleviate  
67 hallucination and fully leverage their powerful text comprehension and processing capabilities,

68 making them promising universal tools for chemical text mining. For instance, Zheng et al.<sup>13</sup> used  
69 prompt engineering to guide ChatGPT in extracting information about metal-organic framework  
70 (MOF) synthesis. Yet, Chen et al.<sup>14</sup> reported that ChatGPT performed significantly worse on  
71 biomedical text mining compared to existing models. This finding contradicts the common belief  
72 in the LLMs' superior comprehension abilities. Either way, LLMs have limitations due to their  
73 model architecture and memory, including a maximum length of prompt tokens. Additionally,  
74 human expressions can be ambiguous, incomplete, vague, and difficult to refine. Outputs may not  
75 strictly adhere to formatting requirements, leading to misunderstanding and poor performance in  
76 mining complex text, such as patents or scientific literature. Therefore, zero-shot or few-shot  
77 prompts are often insufficient to address the diversity of scenarios.

78 In this study, we fine-tuned ChatGPT for five challenging tasks in chemical text mining:  
79 compound entity recognition, reaction role annotation, metal-organic framework (MOF) synthesis  
80 information extraction, nuclear magnetic resonance spectroscopy (NMR) data extraction, and  
81 conversion reaction paragraphs into action sequences. We found for the first time that fine-tuning  
82 ChatGPT significantly enhances performance in text mining tasks, compared to prompt-only  
83 version, while also reducing dependency on the repetitive and extensive prompt engineering  
84 experiments. Meanwhile, we also evaluated other newly emerged generative pre-trained language  
85 models, such as Llama2<sup>15</sup>, T5<sup>16</sup>, and BART<sup>17</sup>. Among these, the fine-tuned ChatGPT achieved  
86 state-of-the-art (SOTA) performance across all five tasks. Remarkably, it even outperformed  
87 models that have been trained specifically for each task and subsequently fine-tuned, based on a  
88 significantly larger amount of in-domain data. This study highlights the potential of fine-tuning  
89 LLMs, like ChatGPT, to revolutionize complex knowledge extraction with their versatility,  
90 robustness, and low code capability. Fine-tuned LLMs can be easily generalizable and can

91 optimize the labour-intensive and time-consuming data collection workflow, even when trained  
 92 with few data. This will accelerate the discovery and creation of novel substances, making them  
 93 powerful tools for universal use.



94  
 95 **Fig. 1 | Schematics of fine-tuning ChatGPT for chemical text mining.** **a**, The pipeline of fine-tuning ChatGPT on proprietary  
 96 data. The green OpenAI logo symbolizes official gpt-3.5-turbo, while the blue one symbolizes fine-tuned gpt-3.5-turbo. **b**,  
 97 Supervised fine-tuned ChatGPT outperforms prompt-only ChatGPT in some customized scenarios. **c**, Illustration of  
 98 cheminformatics insights to be extracted from paragraph. And illustration of the five practical tasks in chemical text mining with  
 99 respective example outputs, including Paragraph2Compound, Paragraph2RXNRole, Paragraph2MOFInfo, Paragraph2NMR, and  
 100 Paragraph2Action.

101

## 102 **Overview of Chemical Text Mining Tasks**

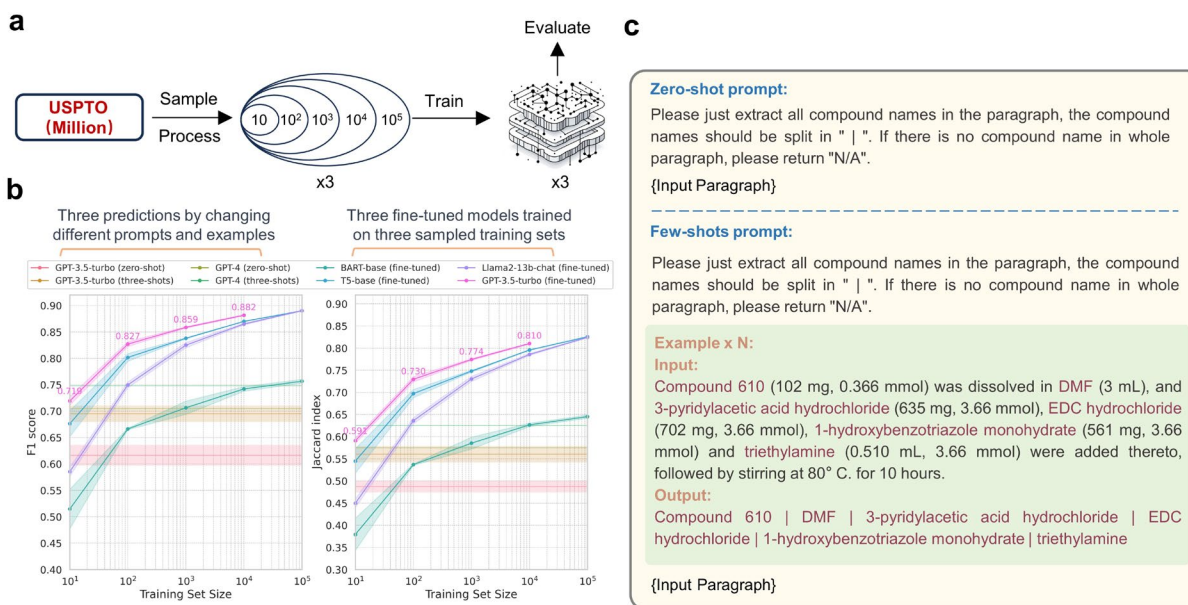
103        Given the complex and diverse information embedded in chemical literature, we designed  
104 five extraction tasks to demonstrate the potential and practicality of LLMs in chemical text mining.  
105 Paragraph2Compound task is to extract all chemical compound entities from the given chemical  
106 paragraph. Paragraph2RXNRole task is to label the reaction roles including product, reactant,  
107 catalyst, temperature, solvent, time, and yield in the paragraph. Paragraph2MOFInfo task is to  
108 extract all MOF synthesis conditions including compound name, metal source, metal amount,  
109 linker, linker amount, modulator, modulator amount or volume, solvent, solvent volume, reaction  
110 temperature and reaction time. Paragraph2NMR task is to extract the IUPAC name, experimental  
111 condition including frequency and solvent as well as chemical shift data for both <sup>1</sup>H NMR and  
112 <sup>13</sup>C NMR spectra. Paragraph2Action task is to convert experimental procedures to structured  
113 synthetic steps (action sequences). All tasks are unified to sequence-to-sequence formats to  
114 facilitate the uses of LLMs.

115

### 116 **Paragraph2Compound—Extract all chemical compound entities.**

117        Fig. 2a illustrates the sampling process, which narrows down from 100,000 to 10, out of  
118 millions of pre-processed annotations, followed by the training process. Fig. 2b demonstrates the  
119 performance of prompt-only models and fine-tuned models, which are evaluated on a consistent  
120 evaluation set of 1,000 samples across varying training data sizes. These results are obtained from  
121 three independent trials. In the case of prompt-only models, randomness is intentionally introduced  
122 by altering the prompt and examples (Fig. 2c). Given the task's straightforward nature and clear  
123 instructions, even the prompt-only language models achieved decent F1 scores over 0.6. For fine-  
124 tuned models, the sampling and training process for the training set is repeated three times, as

125 depicted in Fig. 2a. As shown in Fig. 2b, all fine-tuned models demonstrate a performance  
 126 improvement, especially in terms of the F1 score and Jaccard index, proportional to the increase  
 127 in dataset size. These models outperform the prompt-only models designed for this task. When the  
 128 training data size is substantial enough, the F1 scores of ChatGPT, Llama2, and T5 can reach close  
 129 to 0.9, and the Jaccard index can approach 0.8. Notably, gpt-3.5-turbo, when fine-tuned, showed  
 130 minimal fluctuations and superior performance. However, it is essential to emphasize that the cost  
 131 of fine-tuning gpt-3.5-turbo increased tenfold with each tenfold increase in data volume. Our  
 132 experimentation with gpt-3.5-turbo were capped at 10,000 training samples for 3 epochs due to  
 133 OpenAI's limitations, resulting in a nearly 90-dollar expense—a significant investment in  
 134 computational resources. In contrast, other fine-tuned language models have displayed notable  
 135 cost advantages in this simple task.



136  
 137 **Fig. 2. | Design and Performance for Paragraph2Compound task.** **a**, The workflow of sampling and training based on USPTO  
 138 dataset for Paragraph2Compound task. **b**, The performance of different models across varying size of training set. The data point  
 139 and the shaded areas represent respectively the mean values and standard deviations derived from three independent trials. **c**,  
 140 Example of the zero-shot and three-shots prompts utilized for Paragraph2Compound task.

141

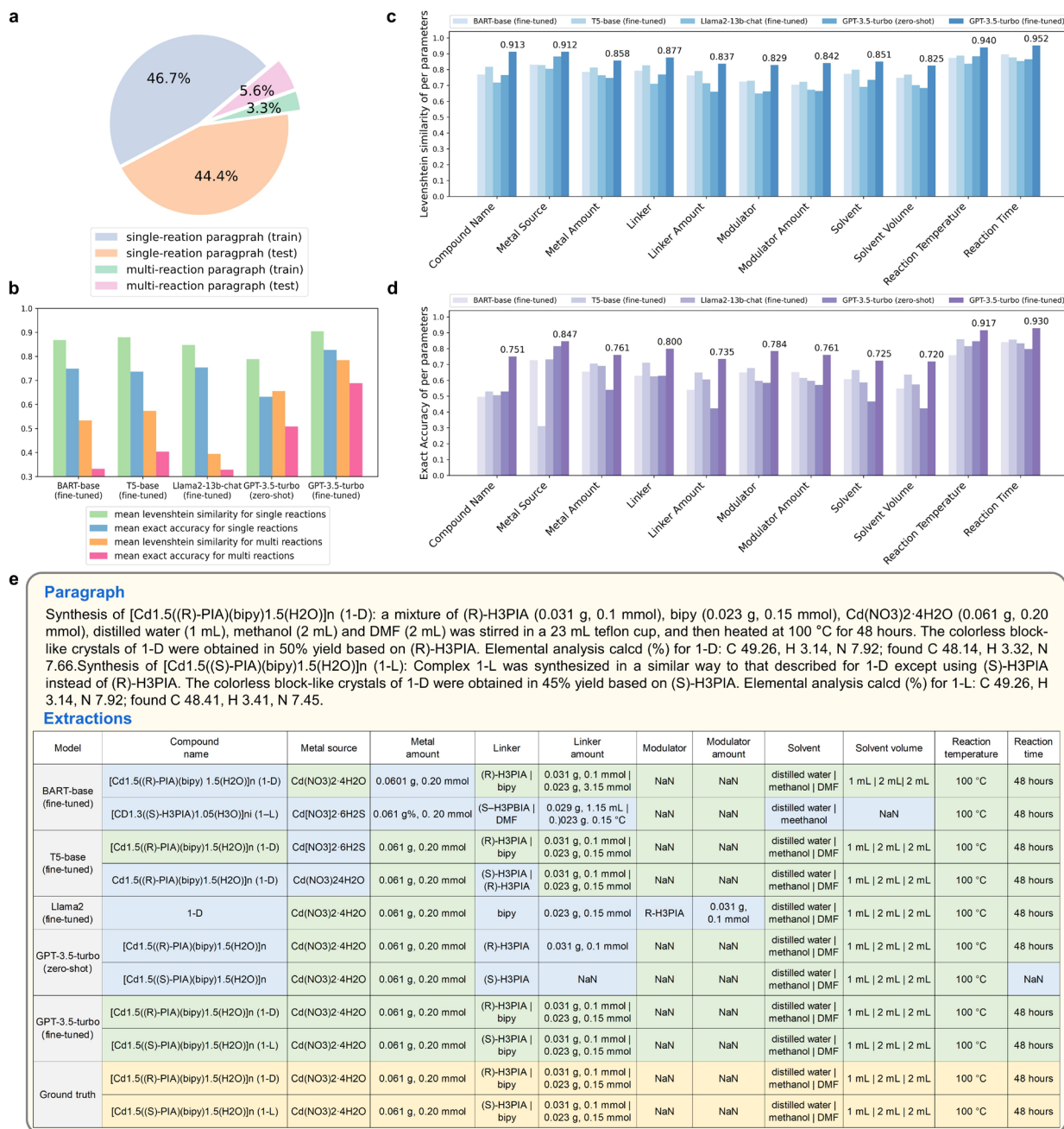
## 142 **Paragraph2RXNRole—Product Extraction and Reaction Role Labelling.**

143 According to Guo et al.<sup>7</sup>, the Paragraph2RXNRole task comprises two subtasks. The first is  
144 to extract the central product, and the second is to label the associated reaction roles within  
145 specified paragraphs (Fig. 3a). For these tasks, Guo et al. developed two-stage BERT-like token-  
146 multi-classification models. To enable a fair comparison with generative language models, we  
147 converted the data into sequence-to-sequence formats by adding <Role\*Compound\*Role>  
148 annotations to the input paragraphs. We then converted the language models' outputs back into  
149 lists of BIO-tags, followed by post-processing to align with the original BIO-tags labels for  
150 assessment. For product extraction, our fine-tuned ChatGPT (best over one epoch) achieved a F1  
151 score of 77.1%, slightly surpassing the previous SOTA approach, ChemBERT, which scored 76.2%  
152 (Fig. 3b). For reaction role labelling, our fine-tuned ChatGPT (best over five epochs) achieved a  
153 F1 score of 83.0%, significantly outperforming the previous SOTA approach, ChemRxnBERT,  
154 which scored 78.7% (Fig. 3c). It's notable that the fine-tuned ChatGPT models, which cost only  
155 \$1 and \$5 respectively, demonstrated extremely high cost-effectiveness with small training  
156 datasets. In contrast, ChemBERT was domain-adaptive pre-trained on 9,478,043 sentences from  
157 200,000 journal articles, and ChemRxnBERT was further task-adaptive trained on 944,733  
158 reaction-inclusive sentences. We should also mention that the outputs of fine-tuned ChatGPT and  
159 Llama2 align almost perfectly with the input text, with 100% and 99% post-processing-free ratios  
160 respectively. On the other hand, the outputs of fine-tuned T5 and BART require additional  
161 alignment due to their tokenization and vocabulary limitations, with a ratio of only 31% that does  
162 not require post-processing. Even after post-processing, the F1 scores of T5 and BART were  
163 significantly lower than those token-classification models or LLMs such as ChatGPT and Llama2.





179 Exact accuracy rates for single and multiple reactions are 82.7% and 68.8%, respectively (Fig. 4b).  
180 As depicted in Fig. 4c and Fig. 4d, while most models achieve high Levenshtein similarity across  
181 the 11 parameters, only a few maintain high exact accuracy, which is the golden metric that we  
182 mainly focus on. Considering that some MOF synthesis paragraphs may include multiple reactions,  
183 we provide an example of multi-reaction extraction by various models in Fig. 4e. The paragraph  
184 includes two reactions, the first with (R)-H3PIA and bipy as linkers, providing all reaction  
185 conditions explicitly, and the second with the substitution of (R)-H3PIA with (S)-H3PIA, keeping  
186 all other conditions unchanged. Most models successfully interpreted the semantics and extracted  
187 two reactions from the MOF synthesis paragraph. However, only the fine-tuned ChatGPT perfectly  
188 extracted information that matched our annotated ground truth. Other models showed varying  
189 degrees of incompleteness, particularly with items involving multiple components and their  
190 quantities.

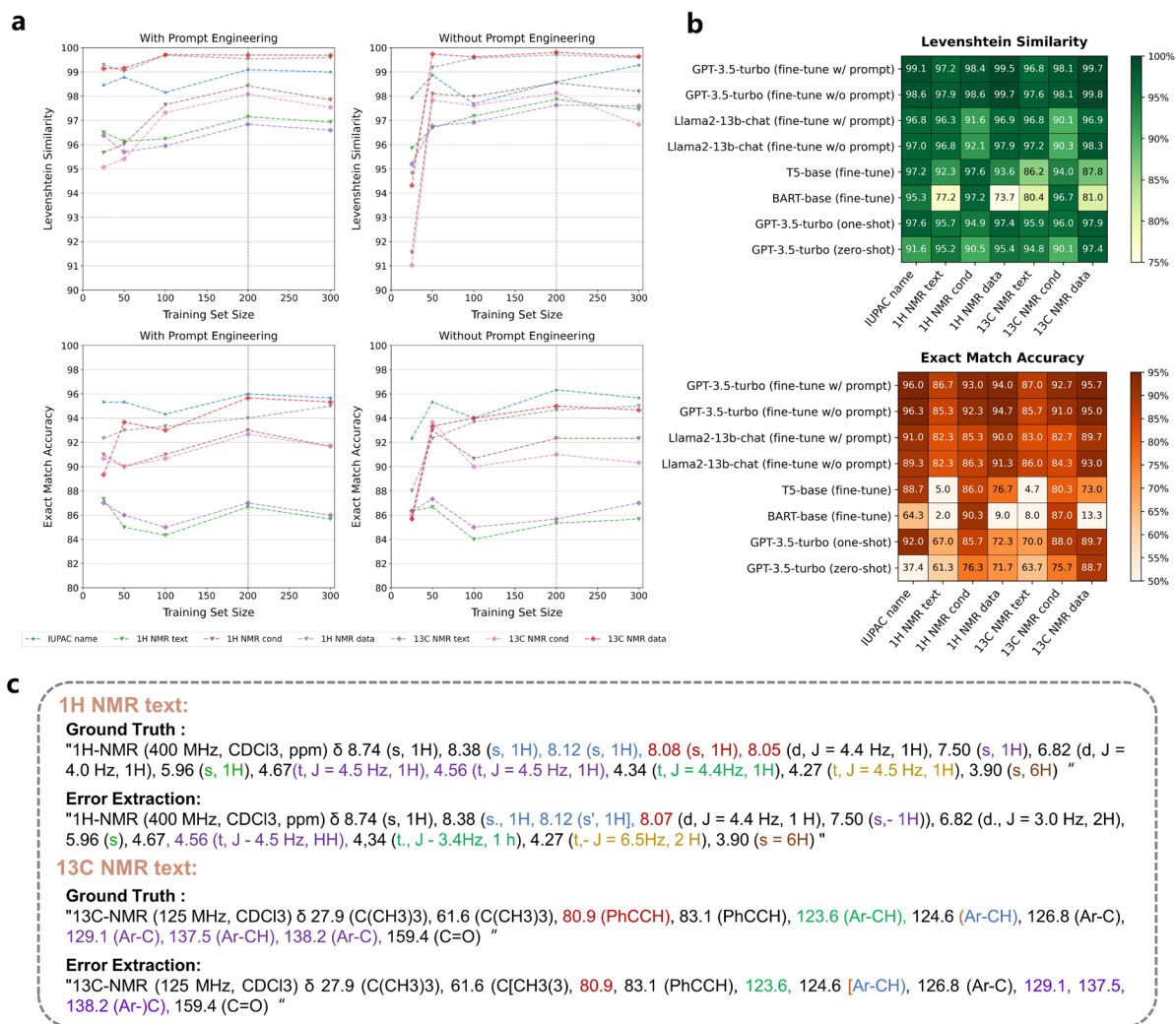


191  
 192 **Fig. 4. | Design and Performance for Paragraph2MOFInfo task.** **a**, A statistic of the dataset. **b**, Mean performance of  
 193 Levenshtein similarity and exact match accuracy by different models. **c**, Levenshtein similarity for 11 parameters in the  
 194 Paragraph2MOFInfo task. **d**, Exact match accuracy for 11 parameters in the Paragraph2MOFInfo task. **e**, An example of extractions  
 195 by different models from a multi-reaction MOF synthesis paragraph. The cells in yellow represented the ground truth. The cells in  
 196 green represented the exact match predictions. The cells in blue represented the incorrect predictions.

197

198 **Paragraph2NMR—Extraction of experimental conditions and NMR chemical shifts.**

199 The impact of training set sizes and the use of prompt engineering on the performance of fine-  
200 tuning ChatGPT in extracting NMR information is illustrated in Fig. 5a. Regardless of the training  
201 data size for fine-tuning (ranging from 25 to 300), or the presence of prompt engineering, there are  
202 hardly any significant fluctuations in performance. This holds true for metrics such as Levenshtein  
203 similarity and exact match accuracy of the fine-tuned ChatGPT when the numbers of training  
204 samples exceed 50. This demonstrates the strong learning capability and robustness of LLMs. Fig.  
205 5b illustrates the performance of different generative language models using the same 200 training  
206 data. In terms of Levenshtein similarity, a metric based on edit distance, almost all fine-tuned  
207 language models achieved impressive scores, outperforming ChatGPT that solely relies on prompt  
208 engineering (Fig. 5b). However, when considering the exact match accuracy metric, where each  
209 character must perfectly align with the ground truth count, LLMs such as ChatGPT and Llama2  
210 take the lead. While fine-tuned T5 and BART manage to extract the majority of the text, they often  
211 miss or mistakenly copy several characters. This contributes to a significant decrease in their exact  
212 match accuracy metric, as shown in Fig. 5c. In this context, the extraction of long complex text by  
213 LLMs is more standardized and high-quality, aligning more closely with human expectations. It is  
214 worth noting that fine-tuning Llama2 provides an alternative approach for deploying text mining  
215 locally, given its exceptional exact match accuracy.



216

217 **Fig. 5. | Performance for Paragraph2NMR task.** a, The performance of fine-tuning ChatGPT with and without prompt engineering

218 as it varies with training data size. b, Heat map illustrating Levenshtein similarity and exact match accuracy of various models in

219 extracting each NMR information. c, Examples of error extractions by T5 and BART, compared with the ground truth.

220

221 **Paragraph2Action: Action sequence extracted from an experimental procedure.**

222 The above-mentioned extraction tasks simply require the model to replicate specific

223 information from the paragraph. However, the Paragraph2Action task requires the model to

224 understand and transform the paragraph. Clearly, ChatGPT with prompt engineering has difficulty

225 with this task, especially when it involves multiple complex conversions and insufficient prompt

226 descriptions. To gauge the maximum potential of ChatGPT using only prompts, we incrementally

227 increased the number of transformation examples from 6 to 30. Despite encompassing all types of  
228 actions at least once and nearly reaching the token limit of 4,096 with 30 examples, ChatGPT's  
229 performance in a few-shot scenario remains disappointingly poor. It only achieved 13.9% full  
230 sentence exact accuracy, a BLEU score of 49.5, and a Levenshtein similarity of 66.0. In contrast,  
231 fine-tuning pre-trained language models with a small amount of data can yield decent results,  
232 achieving over 50% full sentence exact accuracy and over 80% for both BLEU Score and  
233 Levenshtein similarity. Remarkably, after 3 epochs of fine-tuning ChatGPT on 1,060 hand-  
234 annotated training data, we achieved 62.5% full sentence exact accuracy, an 84.8 Modified BLEU  
235 score, and an 87.6 Levenshtein similarity. This process took only 1 hour and cost \$3 for fine-tuning.  
236 These metrics surpass the SOTA results previously reported by Vaucher et al.<sup>1</sup>, which used an  
237 ensemble of three models, each task-adaptively pre-trained on 2 millions rule-based data and  
238 refined on 14,168 augmented data. Interestingly, further improvement was achieved by  
239 augmenting the training data size to 14,168. This resulted in 69.0% full sentence exact accuracy,  
240 an 86.4 Modified BLEU score, and an 89.9 Levenshtein similarity (Table 1). For tasks involving  
241 “fuzzy rules” or hard-to-define extraction, fine-tuning ChatGPT offers significant advantages in  
242 tailoring the transformation with a small amount of annotated data.

243  
244  
245

Tabel 1 | Performance on Paragraph2Action task.

Model	Training data strategy	100% accuracy	90% accuracy	75% accuracy	Modified BLEU score	Levenshtein similarity	Cost
GPT-3.5-turbo (6-shots)	No training	8.2	16.8	34.7	38.6	59.4	905 mean tokens
GPT-3.5-turbo (12-shots)	No training	8.8	19.3	42.3	43.1	62.3	1,374 mean tokens
GPT-3.5-turbo (18-shots)	No training	13.1	23.3	42.6	44.4	64.3	1,670 mean tokens
GPT-3.5-turbo (24-shots)	No training	14.8	25.9	45.5	47.0	65.8	2,598 mean tokens
GPT-3.5-turbo (30-shots)	No training	13.9	26.4	47.2	49.5	66.0	3,610 mean tokens
Transformer (single model) *	No task-adaptive pretraining, no augmentation (1,060)	13.1	15.1	21.9	22.5	45.9	-
BART-base (fine-tuned)	No task-adaptive pretraining, no augmentation (1,060)	51.1	65.9	77.6	73.2	83.9	-
T5-base (fine-tuned)	No task-adaptive pretraining, no augmentation (1,060)	57.7	71.6	83.2	81.8	86.8	-
Llama2-13b-chat (fine-tuned)	No task-adaptive pretraining, no augmentation (1,060)	56.8	66.8	80.7	80.3	86.0	40 min for training
GPT-3.5-turbo (fine-tuned)	No task-adaptive pretraining, no augmentation (1,060)	<b>62.5</b>	<b>72.7</b>	<b>82.9</b>	<b>84.8</b>	<b>87.6</b>	3 epochs, \$ 3, 1h
Transformer (single model) *	No task-adaptive pretraining, augmented unique (14,168)	37.8	47.7	62.8	64.7	76.4	-
BART-base (fine-tuned)	No task-adaptive pretraining, augmented unique (14,168)	52.0	68.5	80.1	74.4	84.8	-
T5-base (fine-tuned)	No task-adaptive pretraining, augmented unique (14,168)	59.7	74.1	82.4	84.1	87.1	-
Llama2-13b-chat (fine-tuned)	No task-adaptive pretraining, augmented unique (14,168)	60.2	70.4	83.5	81.6	87.9	9 hours for training
GPT-3.5-turbo (fine-tuned)	No task-adaptive pretraining, augmented unique (14,168)	<b>69.0</b>	<b>78.1</b>	<b>86.9</b>	<b>86.4</b>	<b>89.9</b>	5 epochs, \$ 92, 1.5 h
Transformer (single model) *	Task-adaptive pretraining (2 million), no augmentation (1,060)	56.8	67.3	80.4	81.5	85.7	-
Transformer (single model) *	Task-adaptive pretraining (2 million), augmented unique (14,168)	59.4	70.5	81.8	84.3	86.7	-
Transformer (ensemble models) *	Task-adaptive pretraining (200w+), augmented unique (14,168)	<b>60.8</b>	<b>71.3</b>	<b>82.4</b>	<b>85.0</b>	<b>86.6</b>	-

247 The symbol “\*\*” represented the result reported by Vaucher. The result in black bold is the best previous performance.

248 The result in red bold is the best new performance.

249

## 250 Discussion

251 Chemical text mining expedites scientific discovery in chemistry. Previously, tasks involving  
 252 complex chemical language and sophisticated processing required the development of specific  
 253 domain-focused models. Now, the fine-tuning of universal LLMs offers a highly generalized and  
 254 cost-effective solution. We have demonstrated the impressive efficacy and high exact accuracy of  
 255 fine-tuning LLMs, especially ChatGPT, across five tasks in text mining. An examination of  
 256 incorrect predictions revealed that only a small proportion were entirely incorrect, while most were  
 257 acceptable alternatives to the ground truth or even pointed out the incorrect labels (Supporting  
 258 Information). These errors can be attributed to inconsistent annotation standards and the inherent  
 259 ambiguity of terms with multiple interpretations or functions. Therefore, improving the formatted

260 data extraction requires continuous efforts, including the refinement of specific rules and the  
261 enrichment of annotations prone to misinterpretation during training and inference. With detailed  
262 specifications and high-quality formatted data, the fine-tuning method based on LLMs is highly  
263 reliable. It can be easily extended to tasks related to extracting information from scientific literature  
264 and transforming data into customized format. This approach will significantly contribute to the  
265 development of extensive databases like SciFinder<sup>18</sup> and Reaxys<sup>19</sup>, which gather comprehensive  
266 synthesis data through automated curation and expert verification.

267 In this work, we have scratched the surface of the vast potential of LLMs in chemistry and  
268 materials science by fine-tuning LLMs for chemical text mining. Technically, advancements like  
269 wider context windows, faster inference approaches, and improved model architectures in the era  
270 of LLMs are anticipated to further enhance text mining. However, it's essential to consider what  
271 else can be achieved with LLMs and how we can develop more effective LLMs for chemistry and  
272 materials science. For instance, LLMs have the potential to revolutionize predictive modelling by  
273 incorporating the extensive “fuzzy knowledge” encapsulated within scientific literature, especially  
274 in chemistry and drug discovery. By combining empirical results with documented knowledge,  
275 LLMs could assist chemists identify patterns in experiments that might otherwise be missed,  
276 predict properties of compounds and outcomes of reactions, and even generate new chemical  
277 hypotheses and theories. Furthermore, the integration of LLMs' comprehension with specialized  
278 tools could substantially lower the barrier of chemists to use these tools throughout the entire  
279 workflow, thanks to interactive interfaces in natural language. Future research could investigate  
280 how to merge formatted laboratory data with wealth of information in scientific literature and  
281 develop the multimodal capability to enrich specific domain knowledge for LLMs. This endeavour  
282 will require a sustained, long-term effort.



283 For the first time, we have demonstrated the effectiveness of fine-tuning ChatGPT and the  
284 potential of LLMs in chemical text mining. We conducted five complex tasks: compound entity  
285 recognition, reaction role labelling, MOF synthesis information extraction, NMR data extraction,  
286 and the transformation from reaction procedures to action sequences. Chemical text mining  
287 remains a challenging professional domain when leveraging language model mining, even with  
288 prompt engineering. However, LLMs that are fine-tuned with appropriate annotations can produce  
289 structured outputs that perfectly fulfil human requirements not easily expressed in natural language.  
290 This feature fully utilizes their natural language understanding and formatting capability. Using  
291 chemical text mining as an example, this study provides guidance on fine-tuning of LLMs to serve  
292 as universal knowledge extraction toolkits. These toolkits can be easily extended for automated  
293 extraction from documents and rule-based formatted transformations. Our work lays the  
294 groundwork for transformative applications of LLMs in knowledge extraction within the chemical  
295 domain.

296

297

## 298 References

- 299 1 Vaucher, A. C. *et al.* Automated extraction of chemical synthesis actions from  
300 experimental procedures. *Nature communications* **11**, 3601 (2020).
- 301 2 Mehr, S. H. M., Craven, M., Leonov, A. I., Keenan, G. & Cronin, L. A universal system  
302 for digitization and automatic execution of the chemical synthesis literature. *Science* **370**,  
303 101-108 (2020).
- 304 3 Steiner, S. *et al.* Organic synthesis in a modular robotic system driven by a chemical  
305 programming language. *Science* **363**, eaav2211 (2019).
- 306 4 Ha, T. *et al.* AI-driven robotic chemist for autonomous synthesis of organic molecules.  
307 *Science Advances* **9**, eadj0461 (2023).
- 308 5 Swain, M. C. & Cole, J. M. ChemDataExtractor: a toolkit for automated extraction of  
309 chemical information from the scientific literature. *Journal of chemical information and*  
310 *modeling* **56**, 1894-1904 (2016).
- 311 6 Mavracic, J., Court, C. J., Isazawa, T., Elliott, S. R. & Cole, J. M. ChemDataExtractor 2.0:  
312 Autopopulated ontologies for materials science. *Journal of Chemical Information and*  
313 *Modeling* **61**, 4280-4289 (2021).
- 314 7 Guo, J. *et al.* Automated chemical reaction extraction from scientific literature. *Journal of*  
315 *chemical information and modeling* **62**, 2035-2045 (2021).
- 316 8 Castro Nascimento, C. M. & Pimentel, A. S. Do Large Language Models Understand  
317 Chemistry? A Conversation with ChatGPT. *Journal of Chemical Information and*  
318 *Modeling* **63**, 1649-1655 (2023).
- 319 9 Clark, T. M., Anderson, E., Dickson-Karn, N. M., Soltanirad, C. & Tafini, N. Comparing  
320 the Performance of College Chemistry Students with ChatGPT for Calculations Involving  
321 Acids and Bases. *Journal of Chemical Education* (2023).
- 322 10 Guo, T. *et al.* What indeed can GPT models do in chemistry? A comprehensive benchmark  
323 on eight tasks. Preprint at <https://arxiv.org/abs/2305.18365> (2023).
- 324 11 Ji, Z. *et al.* Survey of hallucination in natural language generation. *ACM Computing*  
325 *Surveys* **55**, 1-38 (2023).
- 326 12 Zhang, Y. *et al.* Siren's Song in the AI Ocean: A Survey on Hallucination in Large  
327 Language Models. Preprint at <https://arxiv.org/abs/2309.01219> (2023).
- 328 13 Zheng, Z., Zhang, O., Borgs, C., Chayes, J. T. & Yaghi, O. M. ChatGPT Chemistry  
329 Assistant for Text Mining and the Prediction of MOF Synthesis. *Journal of the American*  
330 *Chemical Society* **145**, 18048-18062 (2023).
- 331 14 Chen, Q. *et al.* An Extensive Benchmark Study on Biomedical Text Generation and Mining  
332 with ChatGPT. *Bioinformatics*, btad557 (2023).
- 333 15 Touvron, H. *et al.* Llama 2: Open foundation and fine-tuned chat models. Preprint at  
334 <https://arxiv.org/abs/2307.09288> (2023).
- 335 16 Raffel, C. *et al.* Exploring the limits of transfer learning with a unified text-to-text  
336 transformer. *The Journal of Machine Learning Research* **21**, 5485-5551 (2020).
- 337 17 Lewis, M. *et al.* Bart: Denoising sequence-to-sequence pre-training for natural language  
338 generation, translation, and comprehension. Preprint at <https://arxiv.org/abs/1910.13461>  
339 (2019).
- 340 18 SciFinder. <https://scifinder-n.cas.org> (2023).
- 341 19 Reaxys. <https://www.reaxys.com> (2023).
- 342

## 343 **Methods**

344 More details will be released after the article is publicly available.

## 345 **Metrics for Evaluation**

346 Since fine-tuning ChatGPT does not allow for early stopping based on optimal validation loss,  
347 we report the performances of all models at the best epoch selected from the evaluation set for fair  
348 comparison. Given the task specifics, we use metrics including precision, recall, and F1 score for  
349 evaluating entity-level performance. For sentence-level performance assessment, we use  
350 Levenshtein similarity, exact match accuracy, partial accuracy, and a modified BLEU score.

## 351 **Data Availability**

352 All datasets used in this work are available from the authors upon request.

## 353 **Code Availability**

354 All scripts for training and evaluating can be found on GitHub at [https://github.com/zw-](https://github.com/zw-SIMM/SFTChatGPT_for_chemtext_mining)  
355 [SIMM/SFTChatGPT\\_for\\_chemtext\\_mining](https://github.com/zw-SIMM/SFTChatGPT_for_chemtext_mining).

## 356 **Acknowledgements**

357 We thank all contributions of the open-source community on LLMs. We appreciate Yaghi's  
358 group for guiding in ChatGPT prompt engineering for chemistry tasks.

359 This work was supported by National Natural Science Foundation of China (T2225002,  
360 82273855 to M.Y.Z., and 82204278 to X.T.L.), the National Key Research and Development  
361 Program of China (2022YFC3400504 to M.Y.Z.), SIMM-SHUTCM Traditional Chinese  
362 Medicine Innovation Joint Research Program (E2G805H to M.Y.Z.), and Shanghai Municipal  
363 Science and Technology Major Project.

364 **Contributions**

365 W.Z., J.C.X., Z.Y.F., and M.Y.Z. conceived the idea. M.Y.Z and Z.Y.F designed the research.  
366 W.Z., Q.G.W., F.C. implemented the codes. W.Z., Q.G.W., X.T.K, J.C.X, S.K.N., Z.Y.F collected,  
367 annotated, and processed training data. D.H.C., B.Y.N., Q.S., and X.T.L. checked the data. M.A.C.,  
368 R.Z.Z., Y.T.W., L.H.Z benchmarked the models. W.Z. wrote the initial draft. M.Y.Z., Z.Y.F. and  
369 Z.P.X. reviewed and refined the article. All authors contributed to the analysis of the results. All  
370 authors read and approved the final manuscript.

371 **Competing interests**

372 The authors declare no competing interests.