

# 1 **Application of Self-Evolving AI Agents in Chemical Research: A** 2 **Novel Intelligent Assistance System**

3 Kangyong Ma<sup>1\*</sup>

4 1. College of Physics and Electronic Information Engineering, Zhejiang Normal  
5 University, Jinhua City 321000, China.

6 E-mail address: [kangyongma@outlook.com](mailto:kangyongma@outlook.com)(K.y.Ma); [kangyongma@gmail.com](mailto:kangyongma@gmail.com)

7 ORCID: 0000-0001-6948-2560

8

## 9 ***Abstract***

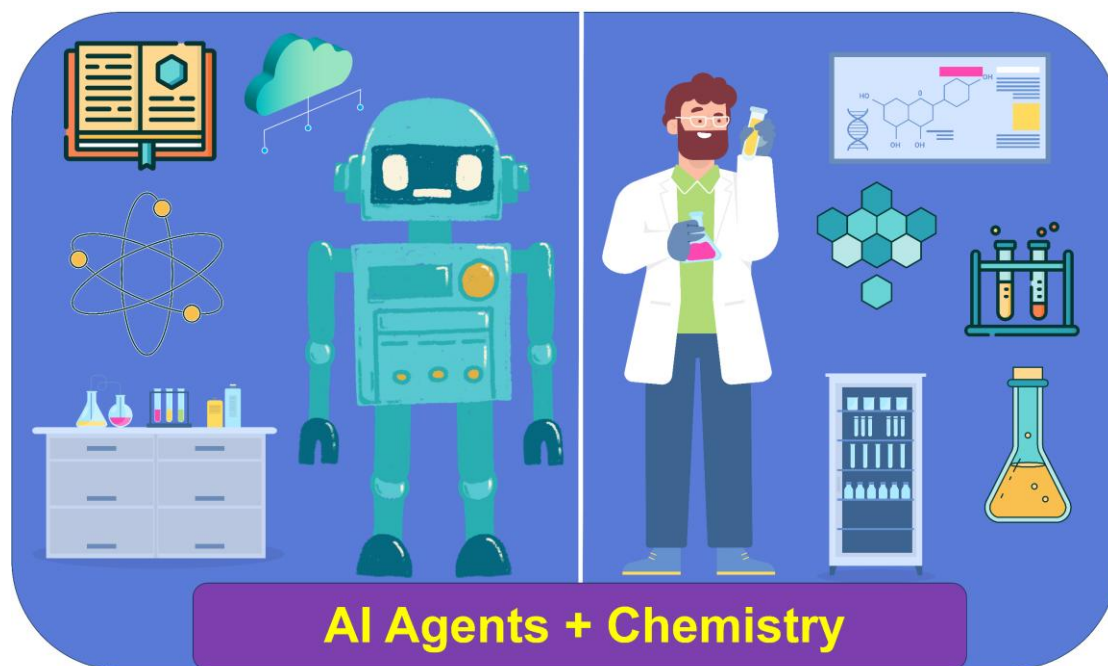
10 This work utilizes collected and organized instructional data from the field of  
11 chemical science to fine-tune mainstream open-source large language models. To  
12 objectively evaluate the performance of the fine-tuned models, we have developed an  
13 automated scoring system specifically for the chemistry domain, ensuring the  
14 accuracy and reliability of the evaluation results. Building on this foundation, we have  
15 designed an innovative chemical intelligent assistant system. This system employs the  
16 fine-tuned Mistral Nemo model as one of its primary models and features a  
17 mechanism for flexibly invoking various advanced models. This design fully  
18 considers the rapid iteration characteristics of large language models, ensuring that the  
19 system can continuously leverage the latest and most powerful AI capabilities. A major  
20 highlight of this system is its deep integration of professional knowledge and  
21 requirements from the chemistry field. By incorporating specialized functions such as  
22 molecular visualization, SMILES string processing, and chemical literature retrieval,  
23 the system significantly enhances its practical value in chemical research and  
24 applications. More notably, the system possesses autonomous evolution capabilities.  
25 Through carefully designed mechanisms for knowledge accumulation, skill  
26 acquisition, performance evaluation, and group collaboration, the system can  
27 continuously optimize its professional abilities and interaction quality. This dynamic  
28 adaptive feature enables the system to evolve autonomously, breaking through the

29 inherent static limitations of traditional AI systems.

30 **Key words :** large language models; fine-tuning; chemistry; autonomous evolution

31

32



33

34

**Fig TOC**

35

36

37

38

39

40

41

42

43

44

45

46

47

## 48 **1. Introduction**

49 Large Language Models (LLMs) stand out as one of the most noteworthy  
50 achievements in the field of artificial intelligence in recent years and represents a  
51 crucial direction for the development of Artificial General Intelligence (AGI)<sup>[1,2]</sup>.  
52 Since the introduction of ChatGPT and GPT-4o, Large Language Models (LLMs) and  
53 Multimodal Large Language Models (MLLMs) have attracted significant interest due  
54 to their versatile abilities in understanding, reasoning, and generating content<sup>[3]</sup>.  
55 However, the current state of this technology still presents significant deficiencies and  
56 imbalances, including persistent illusions, misaligned values, weak specialization, and  
57 the black box effect<sup>[2]</sup>. In this scenario, how to apply Large Language Models (LLMs)  
58 to different professional fields has become a current research hotspot.

59 Fine-tuning has a significant effect on improving the performance of LLM in  
60 specific application scenarios, which lays the foundation for LLM to further promote  
61 scientific progress in various fields<sup>[4,5]</sup>. For example, research by Ouyang et al. (2022),  
62 Wei et al. (2021), and Sanh et al. (2021) demonstrates that fine-tuning language  
63 models on a specific set of tasks significantly enhances their ability to understand and  
64 execute instructions<sup>[6-8]</sup>. This method not only reduces the reliance on large datasets  
65 but also improves the generalization capabilities of the models. Given the scale of  
66 LLMs, a common fine-tuning strategy currently involves adjusting a limited number  
67 of parameters while keeping the rest fixed<sup>[9]</sup>. This technique, known as Parameter-  
68 Efficient Fine-Tuning (PEFT), selectively tunes a small subset of parameters. PEFT  
69 has also gained interest beyond NLP, particularly in the CV community, for fine-  
70 tuning large-parameter visual models like Vision Transformers (ViTs), diffusion  
71 models, and visual-language models<sup>[4]</sup>.

72 However, fine-tuning large models still has some drawbacks. For example, this  
73 method requires substantial computational resources and data. Fine-tuning large  
74 models is also prone to overfitting on small-scale datasets and cannot accurately  
75 reflect potential risks (e.g., "hallucinations"), which may introduce latent hazards.  
76 Additionally, it cannot update its knowledge base in real time<sup>[10]</sup>. The primary reasons

77 for these drawbacks are that both pre-trained large models and fine-tuned large  
78 models use parameter memory to construct a parameterized implicit knowledge  
79 base<sup>[11]</sup>. Hybrid models that combine parametric memory and non-parametric (i.e.,  
80 retrieval-based) memory can address some of these issues<sup>[12-14]</sup>. The Retrieval-  
81 Augmented Generation (RAG) technique improves the accuracy and reliability of  
82 hybrid model generation by integrating knowledge from external databases (non-  
83 parametric memory), especially for knowledge-intensive tasks. This approach also  
84 allows for continuous knowledge updates and the integration of domain-specific  
85 information. RAG synergizes the intrinsic knowledge of large language models with  
86 the extensive dynamic repositories of external databases<sup>[15]</sup>.

87 Furthermore, with the continuous development of LLMs, they are seen as  
88 potential sparks for Artificial General Intelligence (AGI), providing hope for the  
89 construction of general AI agents<sup>[16]</sup>. Currently, AI agents are considered a crucial step  
90 towards achieving AGI, encompassing the potential for a wide range of intelligent  
91 activities<sup>[17-19]</sup>. In many real-world tasks, the capabilities of agents can be enhanced by  
92 constructing multiple cooperative agents<sup>[20]</sup>. Studies have shown that multi-agent  
93 systems help encourage divergent thinking<sup>[21]</sup> (Liang et al., 2023), improve factuality  
94 and reasoning abilities<sup>[22]</sup> (Du et al., 2023), and provide verification<sup>[23]</sup> (Wu et al.,  
95 2023). These features have garnered widespread attention. Currently, the general  
96 frameworks for constructing LLM applications with multiple agents include  
97 AutoGen<sup>[37]</sup>, crewAI<sup>[38]</sup>, Langchain<sup>[39]</sup> and others. Intelligent agents based on large  
98 language models (LLMs) are increasingly permeating various aspects of human  
99 production and daily life. However, designing artificial intelligence agents with self-  
100 evolution capabilities has become a current research hotspot. For example, Li et al.<sup>[24]</sup>  
101 proposed an evolutionary framework for agent evolution and arrangement called  
102 EvoluaryAgent. Qian et al.<sup>[25]</sup> proposed a general strategy for inter-task agent self-  
103 evolution based on Investigation-Consolidation-Exploitation(ICE).

104 These artificial intelligence technologies will provide a new paradigm for  
105 scientific research and open new avenues for scientific innovation, thereby  
106 significantly accelerating the pace of scientific discoveries. The close collaboration

107 between artificial intelligence technologies and scientists heralds the advent of a new  
108 era of scientific exploration and technological breakthroughs<sup>[26,27]</sup>.

109 In recent years, despite the rapid development of artificial intelligence  
110 technology, especially the emergence of large language models, its application in the  
111 field of chemistry has not yet been widely popularized. As an important productivity  
112 tool, artificial intelligence not only improves work efficiency but also provides a new  
113 paradigm for scientific research. For chemistry, a discipline with a long history, how  
114 to combine with this advanced productivity tool to breathe new life into the field has  
115 become an important topic facing the new generation of chemists. This research aims  
116 to address this challenge by developing a dedicated intelligent assistance system for  
117 the field of chemistry through the integration of cutting-edge AI  
118 technologies. Specifically, we first collected and organized a large amount of data  
119 from the field of chemical science to fine-tune mainstream open-source large  
120 language models. Secondly, we designed a set of evaluation systems specifically for  
121 the chemistry field to detect the performance of the fine-tuned models and select the  
122 best-performing model from them. On this basis, we developed an AI assistant for the  
123 chemistry field with autonomous evolution capabilities. This system integrates multi-  
124 agent architecture, retrieval-augmented generation (RAG) technology, online search  
125 functionality, dynamic learning and evolution mechanisms, and an interactive user  
126 interface. It not only provides an innovative platform for chemical research and  
127 education but also offers valuable research opportunities for exploring multi-agent  
128 collaboration and evolution mechanisms in complex systems. By fusing traditional  
129 chemical knowledge with cutting-edge AI technology, this system is expected to  
130 promote innovative development in the field of chemistry and provide new ideas and  
131 tools for solving current scientific and engineering challenges.

132

133

134



Fig1. Research Process

135

136

137

138

139

140

141

142

143

## 144 **2. Related Work**

### 145 **2.1 Fine-tuning LLMs for Applications in the Field of Chemistry**

146 In recent years, with the rapid development of artificial intelligence technology,  
147 Large Language Models (LLMs) have been increasingly applied in the field of  
148 chemical sciences. Through fine-tuning for specific chemical tasks, these models have  
149 demonstrated remarkable potential, bringing new perspectives and methods to  
150 chemical research. Currently, significant progress has been made in chemical science  
151 research using fine-tuned large language models, covering various aspects from  
152 material design to drug discovery. These studies not only showcase the exceptional  
153 ability of LLMs in handling complex chemical problems but also provide innovative  
154 approaches to addressing long-standing chemical challenges.

155 For example, Kevin Maik Jablonka et al. <sup>[45]</sup> fine-tuned the large language model  
156 GPT-3 to perform various tasks in chemistry and materials science, including  
157 properties of molecules and materials, as well as chemical reaction outcomes. Zikai  
158 Xie et al. <sup>[46]</sup> demonstrated the effectiveness of fine-tuned GPT-3 in predicting  
159 electronic and functional properties of organic molecules. Shifa Zhong et al. <sup>[47]</sup>  
160 developed quantitative structure-activity relationship (QSAR) models for water  
161 pollutant activity/properties by fine-tuning GPT-3 models. Seongmin Kim et al. <sup>[48]</sup>  
162 evaluated the effectiveness of pre-trained and fine-tuned large language models  
163 (LLMs) in predicting the synthesizability of inorganic compounds and selecting  
164 synthetic precursors. Results showed that fine-tuned LLMs performed comparably,  
165 and sometimes superiorly, to recent custom machine learning models in these tasks,  
166 while requiring less user expertise, cost, and time to develop.

167 These research findings conclusively demonstrate that fine-tuning LLMs can  
168 significantly enhance their application breadth and effectiveness in the field of  
169 chemical sciences. This approach not only provides powerful tools for chemical  
170 research but also promises to accelerate innovation in chemical sciences, offering new  
171 ideas and methods for solving complex chemical problems. As technology continues  
172 to advance, we can anticipate that fine-tuned LLMs will play an increasingly

173 important role in the field of chemical sciences, driving chemical research towards  
174 deeper and more precise directions.

## 175 **2.2 AI agents in the field of Chemistry**

176 Although large language models (LLMs) have demonstrated excellent  
177 performance in tasks across multiple domains, they face challenges in chemistry-  
178 related problems and lack the ability to access external knowledge sources, limiting  
179 their practicality in scientific applications. To address these deficiencies, researchers  
180 have conducted relevant explorations.

181 For example, Kevin Maik Jablonka et al. <sup>[49]</sup> developed ChemCrow, an LLM  
182 chemical agent designed to complete chemistry tasks such as organic synthesis, drug  
183 discovery, and materials design. By integrating multiple expert-designed chemical  
184 tools and using GPT-4 as the LLM, they enhanced the performance of LLMs in the  
185 field of chemistry and demonstrated new capabilities. Daniil A. Boiko et al. <sup>[50]</sup>  
186 reported on Coscientist, a GPT-4-powered artificial intelligence system capable of  
187 autonomously designing, planning, and executing complex scientific experiments.  
188 Coscientist leverages large language models combined with tools such as internet  
189 searches, document retrieval, code execution, and experimental automation. Andrew D.  
190 McNaughton et al. <sup>[51]</sup> introduced a system called CACTUS (Chemistry Agent  
191 Connecting Tool-Usage to Science), which is an intelligent agent based on large  
192 language models (LLMs) designed to enhance advanced reasoning and problem-  
193 solving capabilities in the fields of chemistry and molecular discovery by integrating  
194 cheminformatics tools.

195 These research findings demonstrate that AI Agents, by expanding the  
196 functionality of large language models, enable their more extensive application in the  
197 field of chemistry.

198

199

200

201



## 202 3. LLMS Fine-tuning Methods

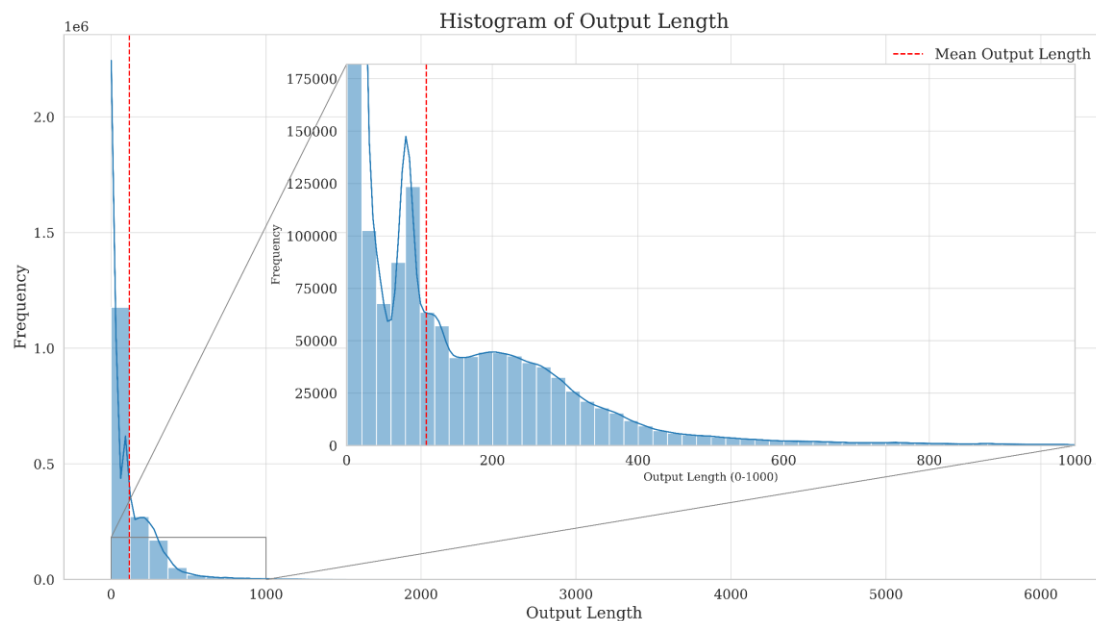
### 203 3.1 Dataset

204 During the adjustment of the LLM, various datasets related to chemical sciences  
205 were utilized. The fine-tuning data for this work comes from the datasets listed in  
206 Table 1, with a total of 1.72 million fine-tuning instructions collected and organized.  
207 Based on this, two different types of instructions were divided for the fine-tuning  
208 training of different large models.

Fine-Tuning Dataset Example	
JSON	Parquet
[ { "instruction": "Write a possible SMILES of given compound. ->", "input": " p-Cresol\n", "output": " Cc1ccc(O)cc1\n"} ]	[ {"from": "human", "value": "Given SMILES, write its compound. -> CCCCC(CC)C=O\n"}, {"from": "gpt", "value": " 2-Ethylhexanal\n"} ]

209  
210 **Figure 2. Example of fine-tuned data**  
211

212 Figures.3 and Figures.4 show the distribution of output character lengths for the  
213 instruction dataset and the usage frequency and types of the 20 most commonly used  
214 instructions in this work.



215

216

**Figure3. Histogram of Output Length**

217

218 Figure 3 illustrates the character count (output length) of the output text in the

219 dataset, which exhibits a wide distribution range, covering both short and long texts.

220 The distribution is concentrated in the 0 to 1000 character range. Short texts (texts

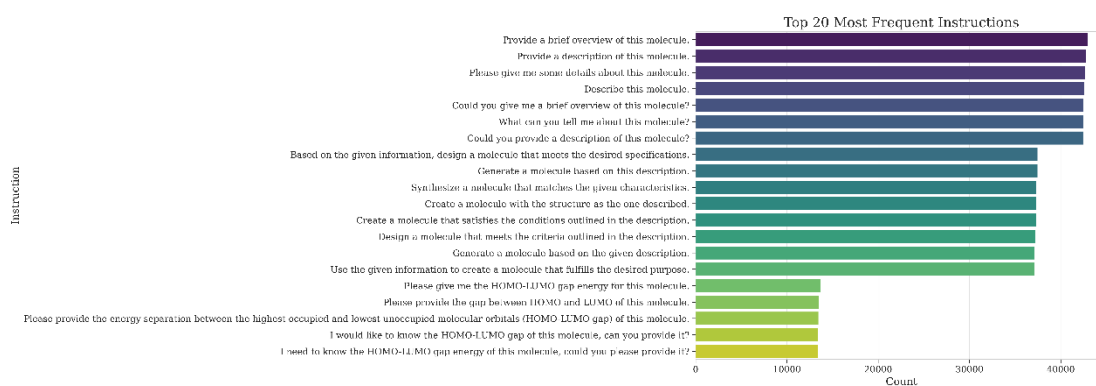
221 with fewer characters) appear more frequently, and as the output length increases, the

222 frequency decreases. Kernel Density Estimation (KDE), also known as Parzen's

223 window<sup>[28]</sup>, is one of the most renowned methods for estimating the underlying

224 probability density function of a dataset. The KDE curve provides a smooth estimate

225 of the distribution within this range, aiding in a more intuitive understanding of the



226

227

**Figure4. Top 20 Most Frequent Instructions**

228 The bar chart (Figure4) shows the frequency of the 20 most common instructions

229 in the dataset for this study. Among these, "Provide a brief overview of this

230 **molecule**" and **"Provide a description of this molecule"** appear significantly more  
 231 often than other instructions, indicating their prominent role in the dataset.  
 232 Nonetheless, other types of instructions also appear, demonstrating the diversity of  
 233 instruction types within the dataset.

234 **Table 1 List of datasets used in our study.**

Dataset	Url link	Data format
ESOL <sup>[43]</sup>	<a href="https://github.com/MasterAIEAM/Darwin/blob/main/dataset/ESOL/ESOL.json">https://github.com/MasterAIEAM/Darwin/blob/main/dataset/ESOL/ESOL.json</a>	Json
MoosaviCp <sup>[43]</sup>	<a href="https://github.com/MasterAIEAM/Darwin/blob/main/dataset/MoosaviCp/MoosaviCp.json">https://github.com/MasterAIEAM/Darwin/blob/main/dataset/MoosaviCp/MoosaviCp.json</a>	Json
MoosaviDiversity <sup>[43]</sup>	<a href="https://github.com/MasterAIEAM/Darwin/blob/main/dataset/MoosaviDiversity/MoosaviDiversity.json">https://github.com/MasterAIEAM/Darwin/blob/main/dataset/MoosaviDiversity/MoosaviDiversity.json</a>	Json
NagasawaOPV <sup>[43]</sup>	<a href="https://github.com/MasterAIEAM/Darwin/blob/main/dataset/NagasawaOPV/NagasawaOPV.json">https://github.com/MasterAIEAM/Darwin/blob/main/dataset/NagasawaOPV/NagasawaOPV.json</a>	Json
ChEMBL <sup>[43]</sup>	<a href="https://github.com/MasterAIEAM/Darwin/blob/main/dataset/chembl/chembl.json">https://github.com/MasterAIEAM/Darwin/blob/main/dataset/chembl/chembl.json</a>	Json
matbench_expt_gap <sup>[43]</sup>	<a href="https://github.com/MasterAIEAM/Darwin/blob/main/dataset/matbench_expt_gap/matbench_expt_gap.json">https://github.com/MasterAIEAM/Darwin/blob/main/dataset/matbench_expt_gap/matbench_expt_gap.json</a>	Json
matbench_glass <sup>[43]</sup>	<a href="https://github.com/MasterAIEAM/Darwin/blob/main/dataset/matbench_glass/matbench_glass.json">https://github.com/MasterAIEAM/Darwin/blob/main/dataset/matbench_glass/matbench_glass.json</a>	Json
matbench_is_metal <sup>[43]</sup>	<a href="https://github.com/MasterAIEAM/Darwin/blob/main/dataset/matbench_is_metal/matbench_is_metal.json">https://github.com/MasterAIEAM/Darwin/blob/main/dataset/matbench_is_metal/matbench_is_metal.json</a>	Json
matbench_steels <sup>[43]</sup>	<a href="https://github.com/MasterAIEAM/Darwin/blob/main/dataset/matbench_steels/matbench_steels.json">https://github.com/MasterAIEAM/Darwin/blob/main/dataset/matbench_steels/matbench_steels.json</a>	Json
Pei <sup>[43]</sup>	<a href="https://github.com/MasterAIEAM/Darwin/blob/main/dataset/Pei/pei.json">https://github.com/MasterAIEAM/Darwin/blob/main/dataset/Pei/pei.json</a>	Json
waterStability <sup>[43]</sup>	<a href="https://github.com/MasterAIEAM/Darwin/blob/main/dataset/">https://github.com/MasterAIEAM/Darwin/blob/main/dataset/</a>	Json

	<a href="#">waterStability/waterStability.json</a>	
description_guided_molecule_design <sup>[44]</sup>	<a href="https://huggingface.co/datasets/zjunlp/Mol-Instructions/tree/main/data">https://huggingface.co/datasets/zjunlp/Mol-Instructions/tree/main/data</a>	Json
forward_reaction_prediction <sup>[44]</sup>	<a href="https://huggingface.co/datasets/zjunlp/Mol-Instructions/tree/main/data">https://huggingface.co/datasets/zjunlp/Mol-Instructions/tree/main/data</a>	Json
molecular_description_generation <sup>[44]</sup>	<a href="https://huggingface.co/datasets/zjunlp/Mol-Instructions/tree/main/data">https://huggingface.co/datasets/zjunlp/Mol-Instructions/tree/main/data</a>	Json
reagent_prediction <sup>[44]</sup>	<a href="https://huggingface.co/datasets/zjunlp/Mol-Instructions/tree/main/data">https://huggingface.co/datasets/zjunlp/Mol-Instructions/tree/main/data</a>	Json
property_prediction <sup>[44]</sup>	<a href="https://huggingface.co/datasets/zjunlp/Mol-Instructions/tree/main/data">https://huggingface.co/datasets/zjunlp/Mol-Instructions/tree/main/data</a>	Json
Retrosynthesis <sup>[44]</sup>	<a href="https://huggingface.co/datasets/zjunlp/Mol-Instructions/tree/main/data">https://huggingface.co/datasets/zjunlp/Mol-Instructions/tree/main/data</a>	Json

235

## 236 3.2 Fine-tuning

237 In this work, we collected and curated 1,720,313 fine-tuning instructions from  
 238 the field of chemical science. Using the unsloth<sup>[29]</sup> tool, we fine-tuned open-source  
 239 large language models including llama-3-8B-Instruct-bnb-4bit, mistral-7B-instruct-  
 240 v0.3-bnb-4bit, gemma-7B-bnb-4bit, gemma-2-9b-bnb-4bit, Phi-3-mini-4k-instruct,  
 241 Mistral-Nemo-Instruct-2407-bnb-4bit and Llama-3.1-8B-Instruct-bnb-4bit. We  
 242 employed the PEFT (Parameter-Efficient Fine-Tuning) method to apply LoRA (Low-  
 243 Rank Adaptation) technique for fine-tuning the pre-trained models. The training  
 244 parameters were configured using SFTTrainer and TrainingArguments. By combining  
 245 quantization techniques, LoRA technology, and optimized training configurations, we  
 246 aimed to enhance performance and optimize resource utilization. Table 2 Parameter  
 247 settings for the fine-tuning process for LLMs.

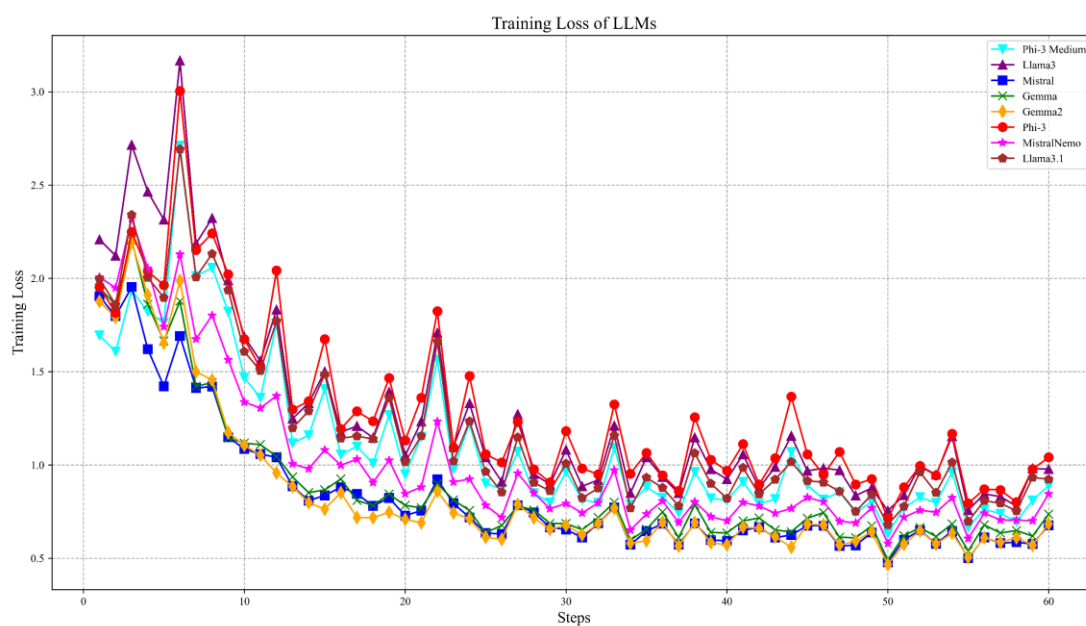
248

**Table 2.** Fine-tuning Process Parameter Settings

Parameter	Value	Description
-----------	-------	-------------

lora_alpha	16	LoRA alpha parameter
max_steps	60	Maximum training steps
learning_rate	2e-4	Learning rate
weight_decay	0.01	Weight decay parameter
seed	3407	Random seed

249 Fig5 represents the training loss curve during the training process of LLMs. In  
 250 the initial phase of training, the loss value is relatively high because the model  
 251 parameters have not yet been optimized, leading to a significant gap between the  
 252 predicted results and the actual values. As the training progresses, the model gradually  
 253 learns and continuously adjusts the parameters, making the predicted results  
 254 increasingly closer to the actual values. Consequently, the error decreases, and the  
 255 loss value gradually declines and tends to stabilize.



256

257

**Fig 5. Training Loss of LLMs**

### 258 3.3 Deployment of LLMs (Large Language Models)

259 After the fine-tuning step in Section 2.2 of the large language model, we  
 260 employed Ollama for the local deployment and testing of fine-tuned LLMs. Model  
 261 parameters were set using the Modelfile configuration file. Specifically, the model's

262 temperature was set to 0.8, and the context window size was configured to 8192  
263 tokens. Additionally, three stop markers were defined to control the boundaries of the  
264 generated text. The detailed configuration is shown in the Fig6. After fine-tuning, the  
265 four large language models were deployed on a local computer for testing. The four  
266 fine-tuned large language models(Llama3-8B,Phi-3-mini,Gemma-7B,Mistral-7B)  
267 were deployed on a local computer with an Intel(R) Core(TM) i5-10210U CPU @  
268 1.60GHz (up to 2.11 GHz) and an NVIDIA GeForce MX250 GPU for testing. The  
269 two fine-tuned models are tested using Google Colab, with Gemma2-9B tested on  
270 a T4 GPU, Phi-3Medium tested on an L4 GPU,Llama3.1-8B tested on Colab CPU  
271 and Mistral Nemo tested on L4 GPU.

```
FROM ./Name.gguf

TEMPLATE """"{{- if .System }}
<|system|>
{{ .System }}
{{- end }}
<|user|>
{{ .Prompt }}
<|assistant|>
""""

SYSTEM """"You are a helpful, smart,
kind, and efficient AI assistant. Your
name is (Set according to your
preferences). You always fulfill the
user's requests to the best of your
ability.""""

PARAMETER temperature 0.8
PARAMETER num_ctx 8192
PARAMETER stop "<|system|>"
PARAMETER stop "<|user|>"
PARAMETER stop "<|assistant|>"
```

272  
273 **Fig6. Model parameter specific settings**  
274

### 275 3.4 Methods for Evaluating the Quality of LLM Responses

276 Based on previous research, evaluation after fine-tuning large language models is  
277 crucial, as it serves as a key tool for identifying current system limitations and  
278 informing the design of more powerful models<sup>[30]</sup>. Therefore, in this work, to assess  
279 the performance of different large models after fine-tuning, 100 questions were  
280 randomly selected from the dataset for model testing. To evaluate the performance of

281 different models after fine-tuning more objectively, this study specifically designed  
282 OptimizedModelEvaluator, an automatic scoring program to evaluate the performance  
283 of different models.

284 Different scoring criteria were designed for different questions. Additionally, the  
285 evaluator considered some special cases in the field of chemical science, assigning  
286 higher weights to key words such as 'reaction', 'mechanism', 'synthesis', and 'catalyst'.  
287 It also recognizes specific chemical terms (e.g., 'alkane', 'alkene', 'alkyne'), considers  
288 conversions between different units when making numerical comparisons (such as kJ  
289 to kcal), and applies special processing for questions involving specific concepts like  
290 LUMO, HOMO, and orbital energies (comparing the signs (positive or negative) of  
291 the extracted answer value and the correct answer value; LUMO and HOMO energies  
292 are typically negative, so the correctness of the sign is important). For questions  
293 involving MOFs, it pays special attention to key concepts such as 'linker', 'node', and  
294 'topology'.

295 The system employs various methods to evaluate the quality of answers. For  
296 numerical problems, it calculates relative errors and assigns corresponding scores. It  
297 uses Levenshtein distance<sup>[31]</sup> or simple word set intersections to compute the  
298 similarity between answers and standard solutions. BLEU scores<sup>[32]</sup> and ROUGE  
299 scores<sup>[33]</sup> are used to assess the quality of generated text and summaries, respectively.  
300 The Flesch<sup>[34]</sup> Reading Ease Index is utilized to evaluate text readability. In addition  
301 to these methods, the system also incorporates evaluation criteria such as keyword  
302 relevance, coherence, conciseness, factual accuracy, and creativity.

303

<b>Scoring Criteria for Different Types of Questions</b>		
<b>Numeric</b>	<b>Descriptive</b>	<b>Generate</b>
<ul style="list-style-type: none"> <li>■ Numeric accuracy: weight 0.6</li> <li>■ Keyword relevance: weight 0.2</li> <li>■ Conciseness: weight 0.2</li> </ul>	<ul style="list-style-type: none"> <li>■ BLEU score: weight 0.2</li> <li>■ ROUGE scores: weight 0.2</li> <li>■ Keyword relevance: weight 0.2</li> <li>■ Readability: weight 0.2</li> <li>■ Coherence: weight 0.2</li> </ul>	<ul style="list-style-type: none"> <li>■ Creativity: weight 0.4 (Assess creativity based on the degree of difference between the answer and the standard answer)</li> <li>■ Coherence: weight 0.3</li> <li>■ Keyword relevance: weight 0.3</li> </ul>

**Fig7. Scoring Criteria for Different Types of Questions**

304

305

306

307

308

309

310

311

312

313

Through these detailed settings, the evaluator can better assess the model's understanding of concepts related to molecular orbital theory, rather than just simple numerical matching. This enables a comprehensive evaluation of AI models' performance in answering chemistry-related questions, covering multiple dimensions including accuracy, relevance, readability, and creativity. Figure 8 illustrates the scoring process. (See supporting information for details).



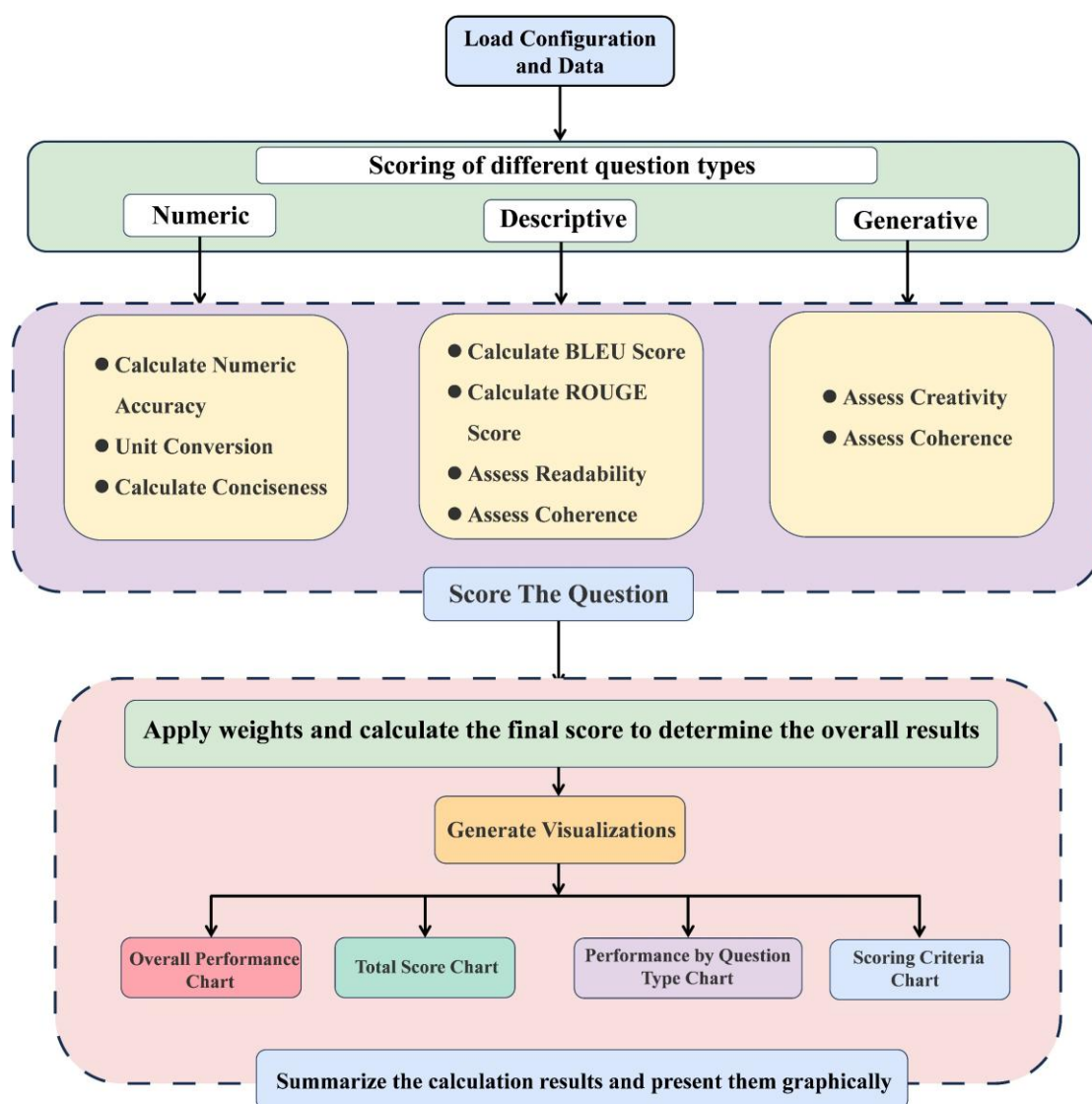
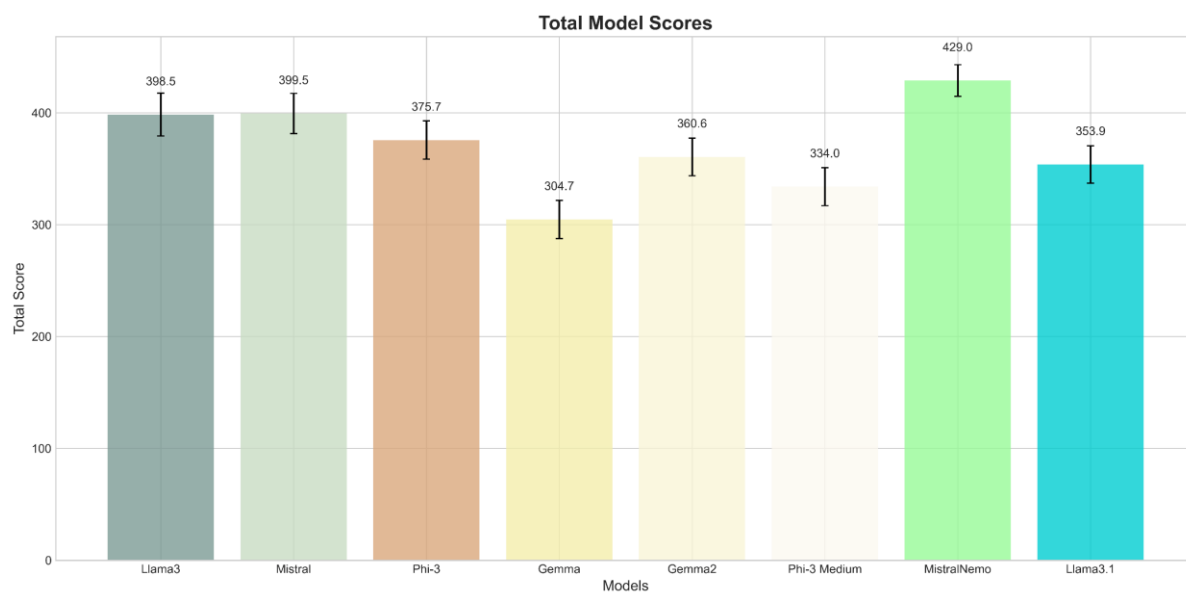


Fig8. Automatic Grading Program Process

### 3.5 LLMs Fine-Tuning Test Results and Discussion

This study conducts a comprehensive evaluation of six fine-tuned large language models: Llama3-8B, Mistral-7B, Phi-3 Mini, Gemma-7B, Gemma2-9B, Phi-3 Medium, Llama3.1 and MistralNemo. Through testing across multiple dimensions, we aim to gain a deep understanding of the performance differences of these models under various tasks and criteria, providing insights for model selection and future optimization directions. Using the automated scoring program introduced in Section 2.4, the fine-tuned models were evaluated with four main metrics: overall score, average performance, multi-dimensional criteria evaluation, and question type classification assessment. Each model was fine-tuned using the same strategy and tested on the same test set (details in the supporting information), ensuring the

328 comparability of the results.



329

330

**Fig9. Total scores for each model**

331

332

333

334

335

336

337

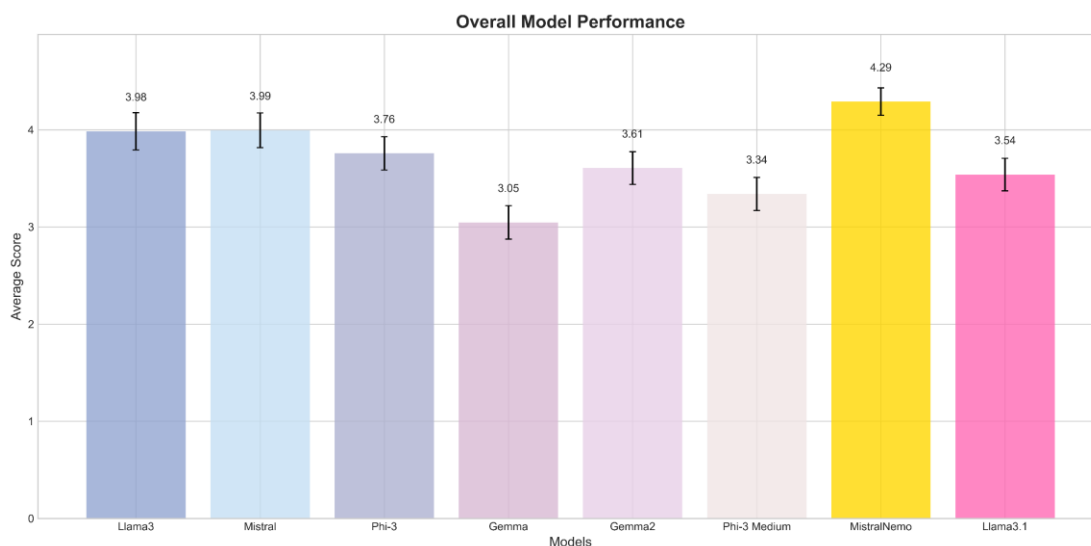
338

339

340

341

The overall model scores (Figure9) show the total scores of eight models, with MistralNemo performing the best, scoring 429.0. Llama3 and Mistral follow closely behind, scoring 398.5 and 399.5 respectively, with both performing very similarly. Phi-3 follows with a score of 375.7. Notably, Gemma2 (360.6 points) shows significant improvement compared to its predecessor Gemma (304.7 points). The iteration from Mistral to MistralNemo also demonstrates the effectiveness of model iterations. However, Phi-3 Medium scored lower than Phi-3 Mini, possibly due to the increased number of parameters making model optimization more challenging, requiring more complex training strategies and computational resources. Additionally, Llama3.1 scored lower than Llama3, indicating that not all model iterations contribute to improved performance after fine-tuning.



342

343

**Fig10. Average scores for each model**

344

345

346

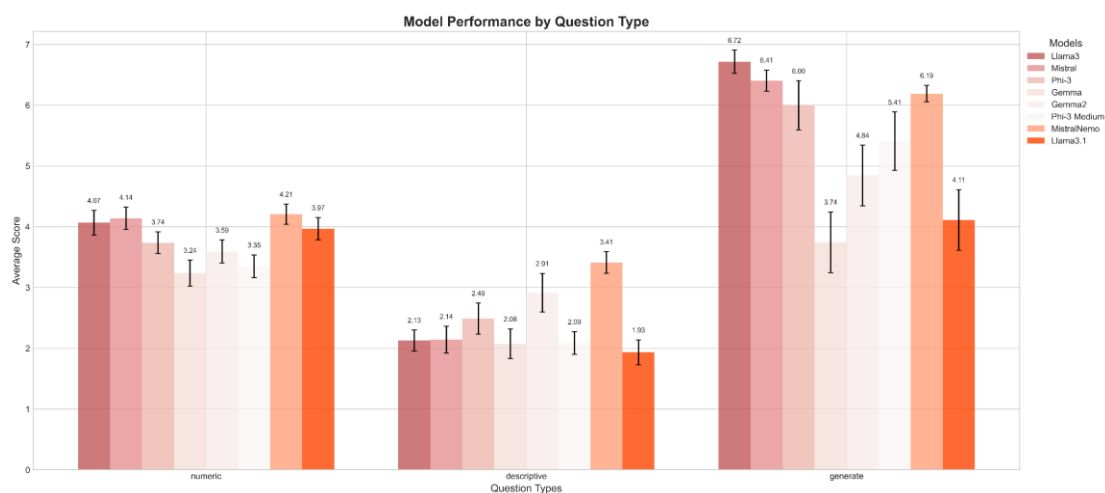
347

348

349

350

Overall Model Performance (Figure10): Figure 10 presents the same ranking trend using a normalized 0-5 scale. This normalization allows for a more intuitive comparison of relative performance differences between models. The 0-5 scale is closer to common rating systems, making performance evaluation more intuitive and relative performance clearer. On the normalized scale, differences after the decimal point become more meaningful, making subtle performance changes more apparent while maintaining the overall structure and relationships of the data.



351

352

**Fig11. Model Performance by Question Type**

353

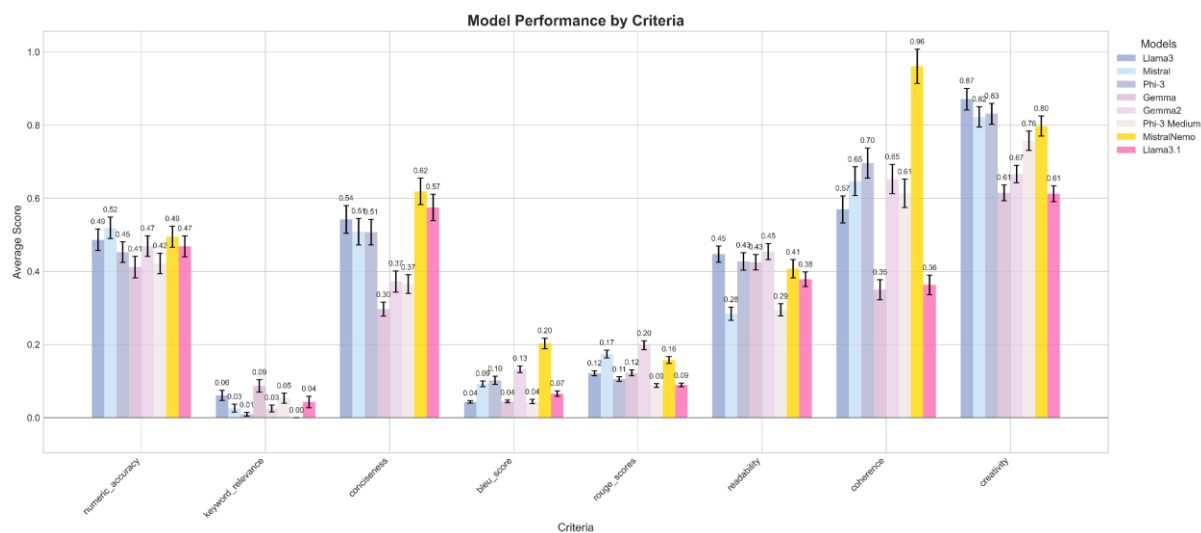
354

355

356

Model Performance Across Different Question Types (Figure11) categorizes questions into numerical, descriptive, and generative types, providing insights into model performance across different task natures. All models perform best on generative questions, with scores ranging from 3.74 to 6.72. This result aligns with

357 the high scores in creativity and coherence in Figure 11, further confirming the strong  
 358 capabilities of large language models in open-ended generation tasks. Descriptive  
 359 questions are the most challenging for all models, with scores ranging from 1.93 to  
 360 3.41, indicating room for improvement in precise description and information  
 361 extraction. Performance on numerical questions falls between the other two types  
 362 (3.24 to 4.21).



363 **Fig12.Model Performance by Criteria**

364 Model Performance Across Different Criteria (Figure12) provides a fine-grained  
 365 analysis of model performance across eight key criteria (numerical accuracy, keyword  
 366 relevance, conciseness, task score, response range, problem-solving ability, coherence,  
 367 and creativity). Models excel in creativity and coherence, generally scoring above 0.6,  
 368 reflecting the common strengths of large language models in generating fluent and  
 369 creative text. Keyword relevance and task score are common challenges for all  
 370 models, with scores generally below 0.2. This suggests that even after fine-tuning,  
 371 models still have room for improvement in accurately grasping task requirements and  
 372 key information. Mistral stands out in numerical accuracy, surpassing other models,  
 373 reflecting its optimization effect on specific tasks. MistralNemo maintains a lead in  
 374 most criteria, showcasing its comprehensive performance advantage.

375 Research findings reveal the significant impact of model iterations on  
 376 performance improvement, particularly evident in the evolution from Gemma-7B to  
 377 Gemma2-9B<sup>[35]</sup> and from Mistral-7B to Mistral-Nemo. However, the iteration from  
 378

379 Llama3-8B to Llama3.1-8B failed to achieve the expected performance leap, possibly  
380 due to different iteration priorities<sup>[36]</sup>. Notably, all tested models face common  
381 challenges, especially in keyword relevance and task scoring, highlighting the  
382 necessity of introducing additional technologies to address these shortcomings.

383 Nevertheless, the outstanding performance of these models in creative and  
384 generative tasks continues to demonstrate the inherent advantages of large language  
385 models in these domains. Test results indicate that fine-tuned large language models  
386 can meet researchers' needs to some extent, but still have many limitations, including  
387 the inability to update data in real-time, lack of online search capabilities, poor  
388 compatibility with specific domains, insufficient response accuracy, and limitations in  
389 decision-making for single large models.

390 Given these limitations exhibited by fine-tuned large language models, this study  
391 developed an artificial intelligence assistant for the chemical domain with  
392 autonomous evolution capabilities. This system cleverly integrates multi-agent  
393 architecture, Retrieval-Augmented Generation (RAG) technology, online search  
394 functionality, dynamic learning and evolution mechanisms, as well as a user-friendly  
395 interactive interface, aiming to comprehensively address the aforementioned  
396 shortcomings and provide researchers with a more intelligent, precise, and practical  
397 auxiliary tool.

#### 398 **4. Self-Evolving AI Agents for Chemistry**

399 This work builds upon the fine-tuning of the aforementioned large language  
400 models to design an AI assistant platform specifically tailored for the field of  
401 chemistry. The platform integrates multi-agent systems, retrieval-augmented  
402 generation, real-time web search, and chemical structure visualization. The system  
403 incorporates AI agents with diverse professional backgrounds (such as laboratory  
404 directors, senior chemists, safety officers, etc.), simulating a virtual chemistry  
405 research team environment. These agents can collaborate, continuously learn and  
406 evolve to provide researchers with comprehensive and professional support in  
407 chemical knowledge, experimental design suggestions, safety guidance, and data

408 analysis. Additionally, the system has the capability to convert chemical structure  
409 formulas (SMILES) into visualized images, greatly enhancing the efficiency and  
410 intuitiveness of chemical research, education, and team collaboration. The system  
411 primarily consists of the following components: Multi-agent system, Retrieval-  
412 augmented generation (RAG), Real-time web search, Chemical structure visualization,  
413 Agent evolution system and User-friendly interface design.

#### 414 **4.1 Multi-Agent System**

415 This system is the core architecture of the project, simulating a real chemical  
416 team. The system contains five specialized agents, each with a specific role and  
417 expertise, together forming a comprehensive and efficient virtual chemical research  
418 team. The Lab\_Director is responsible for overall task allocation and research  
419 direction guidance, ensuring the team's research direction aligns with the overall goals  
420 and coordinating work between agents. The Senior\_Chemist provides in-depth  
421 chemical knowledge and solutions to complex problems, possessing rich chemical  
422 theory and practical experience to handle challenging chemical issues and propose  
423 innovative research ideas. The Lab\_Manager is responsible for experiment planning  
424 and resource management, ensuring the feasibility of experimental plans, managing  
425 laboratory resources, optimizing experimental processes, and improving research  
426 efficiency. The Safety\_Officer ensures all discussions and suggestions comply with  
427 safety standards, focusing on experimental safety, reviewing potential risks of all  
428 experimental protocols, and providing safety operation guidance. The  
429 Analytical\_Chemist focuses on data analysis and instrument use, responsible for  
430 interpreting experimental data, providing instrument operation advice, and ensuring  
431 data accuracy and reliability. This design allows each agent to have its specific area of  
432 expertise, providing in-depth professional knowledge. Agents can complement each  
433 other to solve complex problems collaboratively. For example, when the  
434 Senior\_Chemist proposes an experimental protocol, the Safety\_Officer reviews its  
435 safety, while the Lab\_Manager considers its feasibility. This multi-perspective  
436 analysis allows agents with different backgrounds to analyze problems from various  
437 angles, providing comprehensive insights. The structure simulates the team dynamics

438 of a real chemistry research group, closely mimicking real team decision-making  
439 processes. Each agent in the system is based on a large language model but has  
440 specific system prompts to define its role and expertise, and different language  
441 models can be substituted to meet the needs of different tasks. AutoGen is used to  
442 manage interactions and dialogue flow between agents, adopting a round-robin  
443 approach to select speakers, ensuring each agent has the opportunity to contribute.  
444 The above multi-agent design allows the system to analyze and solve chemical  
445 problems from multiple perspectives, providing comprehensive insights (This  
446 research uses the fine-tuned and performance-tested MistralNemo as a main model for  
447 this system, and all fine-tuned large language models involved in this research have  
448 been uploaded to Hugging Face and can be set to call different large language models  
449 according to different needs KANGYONGMA/Chemistry).

#### 450 **4.2 Retrieval-Augmented Generation (RAG)**

451 RAG is a core functionality of the system, extending the knowledge base of  
452 agents by integrating preloaded chemical literature and experimental data. The RAG  
453 workflow includes document loading, text splitting, vector embedding, vector storage,  
454 similarity search, context enhancement, and answer generation. This process is  
455 implemented using the langchain library and RetrievalQA chain, significantly  
456 improving the accuracy and relevance of answers while reducing the possibility of AI  
457 generating false information. RAG technology enables agents to provide answers  
458 based on the latest chemical research, cite relevant literature to support views, and  
459 associate user queries with existing knowledge bases, thereby greatly enhancing the  
460 system's ability to handle complex chemical problems and provide more precise and  
461 relevant information.

#### 462 **4.3 Real-time Web Search**

463 Another important feature of the system is the ability to perform real-time web  
464 searches by integrating the Tavily search API<sup>[40]</sup> to supplement the preloaded  
465 knowledge base. The workflow of this feature includes query analysis, API calls,  
466 result processing, and information integration. The system uses the requests<sup>[41]</sup> library  
467 to send API requests and implements error handling and retry mechanisms to ensure

468 stability. This feature allows agents to access the latest chemical research and  
469 discoveries, supplement information that may be missing from the preloaded database,  
470 and significantly improve the system's ability to answer current affairs questions. By  
471 combining preloaded data and real-time search, the system can provide users with  
472 comprehensive, up-to-date, and accurate chemical information, excelling particularly  
473 in handling emerging research, latest discoveries, or real-time data-related issues.

#### 474 4.4 Chemical Structure Visualization

475 This feature greatly enhances the system's interactivity and intuitiveness when  
476 discussing chemical structures by converting SMILES<sup>[5]</sup> strings into 2D molecular  
477 structure images. The entire process involves SMILES parsing, molecular object  
478 creation, 2D coordinate generation, image rendering, and encoding, ultimately  
479 displaying on the Web interface. This functionality not only enhances the visual  
480 understanding of chemical concepts and improves the efficiency of discussing  
481 complex molecular structures but also makes the system more suitable for chemical  
482 education and research applications. Its implementation mainly relies on the RDKit<sup>[42]</sup>  
483 library for molecular manipulation and image generation, integrating it into the  
484 message processing flow to achieve automatic detection and conversion of SMILES  
485 strings, thereby providing chemistry researchers with a more intuitive and effective  
486 chemical structure interaction experience.

The screenshot displays the 'Chemistry Lab AI' interface. On the left, a sidebar lists 'Lab Agents' with roles and levels: Lab Director (Level 1), Senior Chemist (Level 1), Lab Manager (Level 2), Safety Officer (Level 2), and Analytical Chemist (Level 3). The main chat area shows a user message: 'Write water solubility of CCC(O)CC in room temperature'. The AI Assistant responds with the same text and a 2D skeletal structure of 2-butanol, with the hydroxyl group labeled 'HO'. Below the structure, the text 'in room temperature' is partially visible. At the bottom, there is a text input field 'Type your message here...', a 'Send' button, and two buttons: 'Set Literature Path' and 'Set Web URL Path'.

487



## Fig14. Dialogue Interface SMILES Visualization

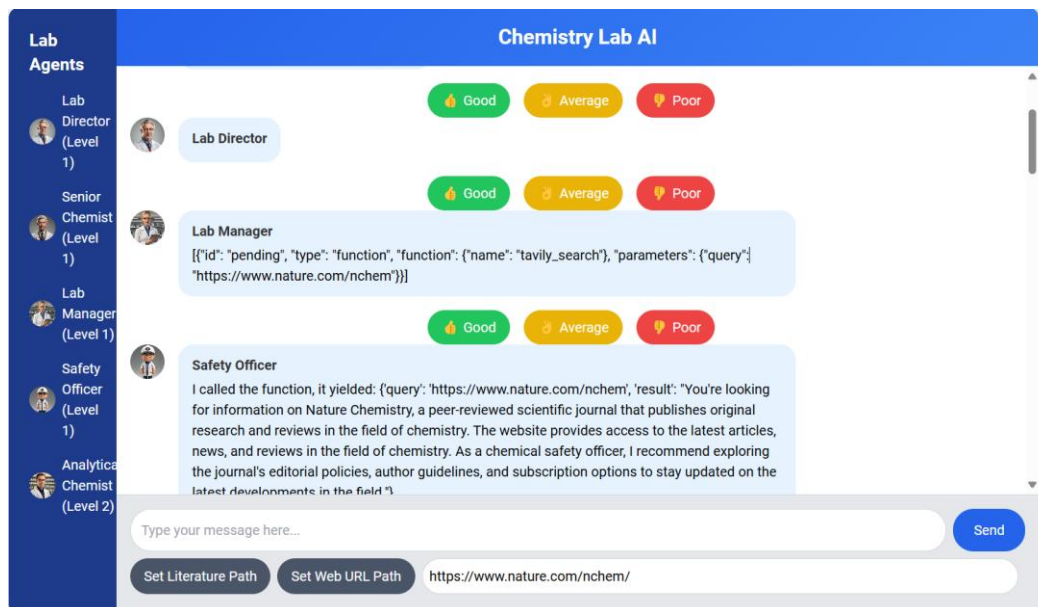
488

### 4.5 Agent Evolution System

489

490 Agent evolution is an innovative feature that allows agents to improve their  
491 performance and knowledge base over time. The system includes four main  
492 components: knowledge acquisition, skill development, performance evaluation, and  
493 adaptive adjustment. Through mechanisms such as knowledge base expansion, skill  
494 tree updates, feedback learning, and cross-learning, agents can learn new chemical  
495 concepts, acquire new problem-solving abilities, and adjust behaviors based on user  
496 feedback. The system uses dynamic data structures to store knowledge and skills,  
497 implements a scoring system to quantify performance, and adopts probability models  
498 to simulate the evolution process. This dynamic learning mechanism enables agents to  
499 continuously improve, adapt to user needs, and simulate human learning and  
500 professional development processes. Over time, agents continuously enhance their  
501 capabilities, providing increasingly relevant and useful information to users, thereby  
502 significantly improving the overall performance and user experience of the system.  
503 This system core is composed of two main classes: ChemistryAgent and  
504 ChemistryLab, implementing functions such as knowledge accumulation, skill  
505 acquisition, performance evaluation, and group evolution. The ChemistryAgent class  
506 stores knowledge and skills through the knowledge\_base and skills attributes,  
507 constantly expanding its capabilities using the learn() and acquire\_skill() methods.  
508 The performance evaluation mechanism records recent performance through  
509 performance\_history, and the evaluate\_performance() method assesses performance  
510 based on user feedback. The evolution mechanism is triggered by the evolve() method,  
511 determining whether to enhance or improve skills based on average performance. The  
512 improve() and refine\_skills() methods are responsible for acquiring new skills and  
513 optimizing existing skills, respectively. The system can also identify areas for  
514 improvement by analyzing interaction history. At the group level, the ChemistryLab  
515 class implements knowledge sharing among agents and multi-round evolution  
516 simulation. This design allows the system to continuously adjust and optimize based  
517 on actual interactions and feedback, continuously improving its professional

518 capabilities and interaction quality in the field of chemistry, forming a dynamically  
519 adaptive and self-improving intelligent ecosystem.



520

521 **Fig.15 User Feedback and Intelligent Agent Evolution Interface**

522

## 522 4.6 User-friendly interface design

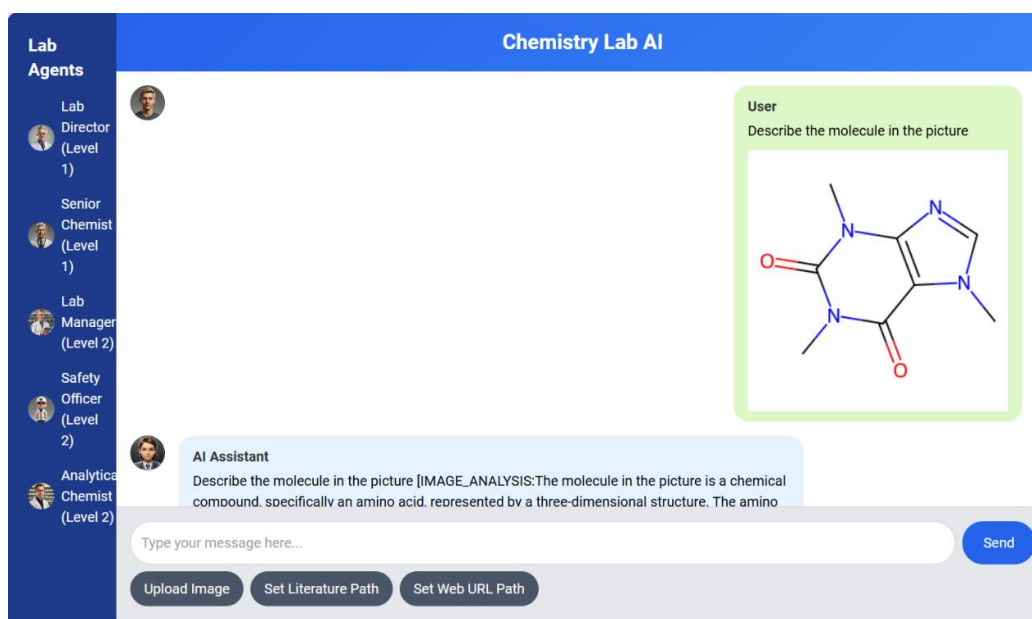
523

523 The project includes an intuitive web interface that can display real-time  
524 conversations between agents, agent status, and feedback mechanisms, providing a  
525 better interactive experience.

525

### 526 **Function Expansion:**

526



527

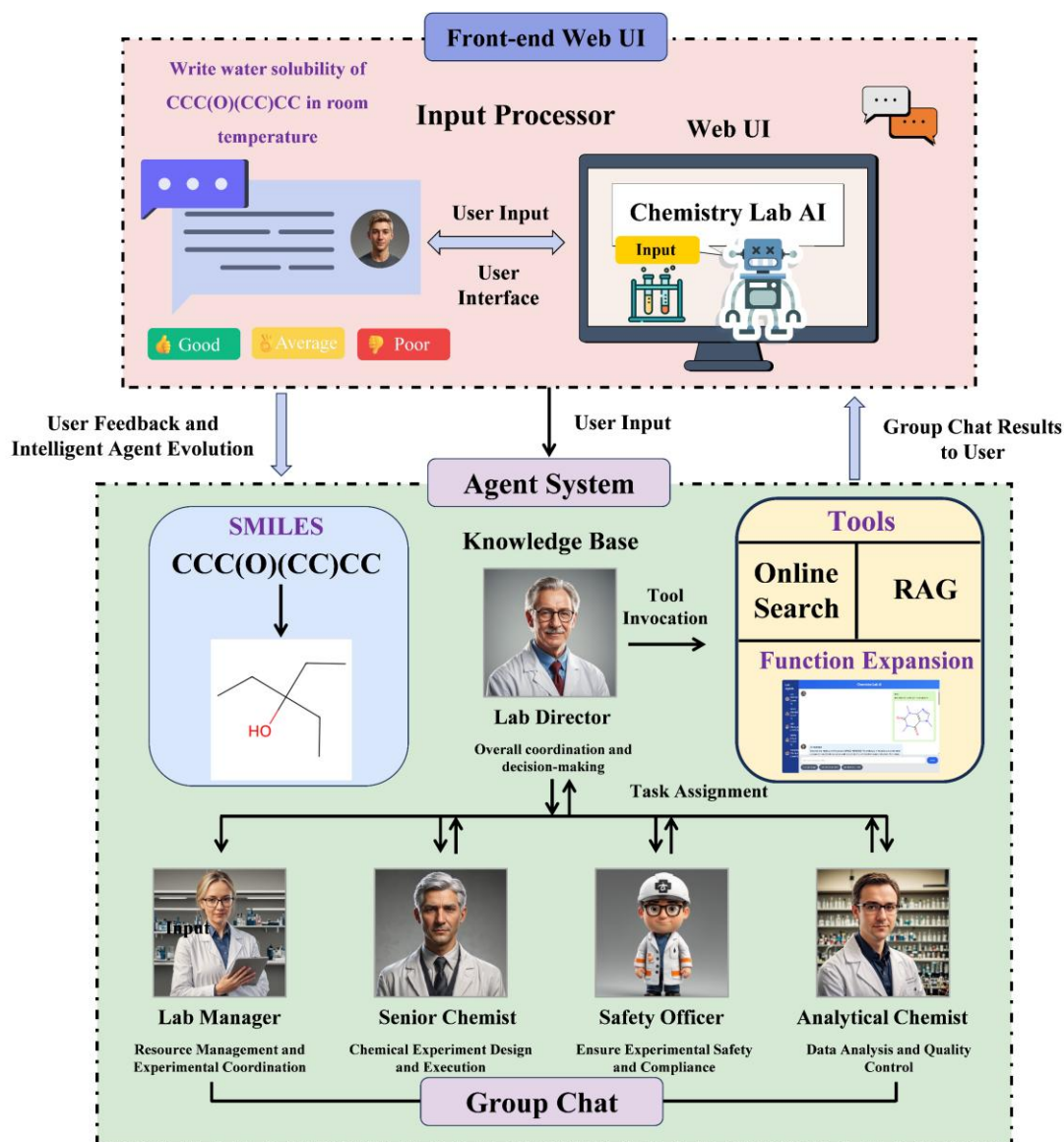
528 **Fig.16 Functionality Expansion—Multimodal Models**

529

## 529 4.7 Functionality Expansion

529

530 During the system design phase, the team fully considered the potential impact  
 531 of model update and iteration, and therefore reserved corresponding upgrade and  
 532 development space. Figure 16 demonstrates the image recognition capabilities after  
 533 the integration of multi-modal large models, which provides an important foundation  
 534 for expanding more functionalities in the future.



535

536

**Fig 17. The Structure of Self-Evolving AI Agents for Chemistry System**

537

538

539

540

541

The system's design fully considers the rapid iteration characteristics of large language models, implementing a flexible mechanism to call upon different advanced models. The system deeply integrates specialized functions in the field of chemistry, such as molecular visualization and SMILES string processing, precisely meeting the needs of chemical research.

542 The core advantage of the system lies in its autonomous evolution capability.  
543 Through knowledge accumulation, skill acquisition, performance evaluation, and  
544 group collaboration, it can continuously optimize its professional capabilities and  
545 interaction quality. This dynamic adaptive feature breaks through the static limitations  
546 of traditional AI systems, providing intelligent and efficient support for solving  
547 complex chemical problems.

548

## 549 **Conclusion**

550 This study utilized 1,720,313 instruction data points from the field of chemical  
551 science to fine-tune 8 mainstream open-source large language models, including  
552 Llama3-8B, Mistral-7B, Phi-3 Mini, Gemma-7B, Gemma2-9B, Phi-3 Medium,  
553 Llama3.1, and MistralNemo. Through an automatic scoring program specifically  
554 designed to evaluate the quality of responses from large language models in the  
555 chemistry domain, the MistralNemo model demonstrated the most outstanding  
556 performance, achieving a total score of 429 points, surpassing other models. Based on  
557 these results, an innovative chemical intelligent assistant system was designed. This  
558 system employs the fine-tuned Mistral Nemo model as its primary model and can call  
559 upon different large models according to task requirements. Furthermore, the system  
560 deeply integrates professional knowledge and requirements from the chemistry field,  
561 featuring specialized functionalities such as molecular visualization, SMILES string  
562 processing, and chemical literature retrieval. Benefiting from knowledge accumulation,  
563 skill acquisition, performance evaluation, and collaborative mechanisms, the system  
564 can continuously optimize its professional capabilities and interaction quality. This  
565 allows the system to learn and grow continuously, breaking through the inherent static  
566 limitations of traditional AI systems and opening up new possibilities for the  
567 application of artificial intelligence in the field of chemistry.

## 568 **Acknowledgements**

569 This work involves the following AI technologies: Llama3-8B, Mistral-7B, Phi-3  
570 Mini, Gemma-7B, Gemma2-9B, Phi-3 Medium, Llama3.1, and MistralNemo. The

571 aforementioned open-source large language models were used for fine-tuning tests in  
572 this work. Tavily Search AI was used for online searches, and sentence-  
573 transformers/all-mpnet-base-v2 was used for RAG (Retrieval-Augmented Generation).  
574 Additionally, Claude 3.5 Sonnet was used to address code issues encountered in this  
575 research, assist in developing the Web UI interface, optimize the multi-agent  
576 framework, and expand multi-agent tools. The Agent avatar in this work was  
577 generated by Stable Diffusion 3. The manuscript was polished using Claude 3.5  
578 Sonnet and ChatGPT-4o. We are grateful for the assistance of these AI technologies in  
579 completing this work.

### 580 **Author contributions**

581 Kangyong Ma was responsible for the conception and design of this study. He  
582 conducted the data analysis and interpretation. He wrote the original draft of the  
583 manuscript and created the visualizations.

### 584 **Conflict of interest**

585 The authors have no conflicts of interest to declare.

### 586 **Data availability**

587 The code for this work is available at <https://github.com/KangyongMa/GVIM>,  
588 while the data and models can be found at <https://huggingface.co/KANGYONGMA>.

### 589 **Funding**

590 This work has no financial support.

591

592

593

594

595

596

597

598

599 **References**

- 600 [1] Pavlick E. Symbols and grounding in large language models[J]. *Philos Trans A*  
601 *Math Phys Eng Sci*, 2023, 381(2251): 20220041.
- 602 [2] Chen W, Yan-Yi L, Tie-Zheng G, et al. Systems engineering issues for industry  
603 applications of large language model[J]. *Applied Soft Computing*, 2024, 151: 111165.
- 604 [3] Xiao H, Zhou F, Liu X, et al. A comprehensive survey of large language models  
605 and multimodal large language models in medicine[J]. *arXiv.org*, 2024.  
606 *arXiv:2405.08603*.
- 607 [4] Han Z, Gao C, Liu J, et al. Parameter-efficient fine-tuning for large models: a  
608 comprehensive survey[J]. *arXiv.org*, 2024. *arXiv:2403.14608*.
- 609 [5] Livne M, Miftahutdinov Z, Tutubalina E, et al. Nach0: multimodal natural and  
610 chemical languages foundation model[J]. *Chemical Science*, 2024.
- 611 [6] Ouyang L, Wu J, Xu J, et al. Training language models to follow instructions with  
612 human feedback[J]. *arXiv.org*, 2022. *arXiv:2203.02155*.
- 613 [7] Wei J, Bosma M, Zhao V Y, et al. Finetuned language models are zero-shot  
614 learners[J]. *arXiv.org*, 2022. *arXiv:2109.01652*.
- 615 [8] Sanh V, Webson A, Raffel C, et al. Multitask prompted training enables zero-shot  
616 task generalization[J]. *arXiv.org*, 2022. *arXiv:2110.08207*.
- 617 [9] Lester B, Al-Rfou R, Constant N. The power of scale for parameter-efficient  
618 prompt tuning[J]. *arXiv.org*, 2021. *arXiv:2104.08691*.
- 619 [10] Nori H, King N, Mckinney S M, et al. Capabilities of gpt-4 on medical challenge  
620 problems[J]. *arXiv.org*, 2023. *arXiv:2303.13375*.
- 621 [11] Lewis P, Perez E, Piktus A, et al. Retrieval-augmented generation for knowledge-  
622 intensive nlp tasks[J]. *arXiv.org*, 2021. *arXiv:2005.11401*.
- 623 [12] Guu K, Lee K, Tung Z, et al. Realm: retrieval-augmented language model pre-  
624 training[J]. *arXiv.org*, 2020. *arXiv:2002.08909*.
- 625 [13] Vladimir Karpukhin B O G S, Ledell Wu S E D C. Dense passage retrieval for  
626 open-domain question answering[J]. *arXiv.org*, 2020. *arXiv:2004.04906*.
- 627 [14] Petroni F, Lewis P, Piktus A, et al. How context affects language models' factual

628 predictions[J]. arXiv.org, 2020. <https://openreview.net/forum?id=025X0zPfn>

629 [15] Gao Y, Xiong Y, Gao X, et al. Retrieval-augmented generation for large language  
630 models: a survey[J]. arXiv.org, 2024. arXiv:2312.10997.

631 [16] Xi Z, Chen W, Guo X, et al. The rise and potential of large language model based  
632 agents: a survey[J]. arXiv.org, 2023. arXiv:2309.07864.

633 [17] Wooldridge, M. J., N. R. Jennings. Intelligent agents: theory and practice. *Knowl.*  
634 *Eng. Rev.*, 10(2):115–152, 1995.

635 [18] Shoham, Y. Agent oriented programming. In M. Masuch, L. Pólos, eds.,  
636 *Knowledge Representation and Reasoning Under Uncertainty, Logic at Work*  
637 [*International Conference Logic at Work, Amsterdam, The Netherlands, December*  
638 *17-19, 1992*], vol. 808 of *Lecture Notes in Computer Science*, pages 123–129.  
639 Springer, 1992.

640 [19] Hutter, M. *Universal artificial intelligence: Sequential decisions based on*  
641 *algorithmic probability*. Springer Science & Business Media, 2004.

642 [20] Wu Q, Bansal G, Zhang J, et al. Autogen: enabling next-gen llm applications via  
643 multi-agent conversation[J]. arXiv.org, 2023. arXiv:2308.08155

644 [21] Tian L, He Z, Jiao W, et al. Encouraging divergent thinking in large language  
645 models through multi-agent debate[J]. arXiv.org, 2023. arXiv:2305.19118

646 [22] Du Y, Li S, Torralba A, et al. Improving factuality and reasoning in language  
647 models through multiagent debate[J]. arXiv.org, 2023. arXiv:2305.14325

648 [23] Wu Y, Jia F, Zhang S, et al. An empirical study on challenging math problem  
649 solving with gpt-4[J]. arXiv.org, 2023. arXiv:2306.01337

650 [24] Shimin Li T S Q C. Agent alignment in evolving social norms[J]. arXiv.org, 2024.  
651 arXiv:2401.04620

652 [25] Cheng Qian S L Y Q, Sun M. Investigate-consolidate-exploit: a general strategy  
653 for inter-task agent self-evolution[J]. arXiv.org, 2024. arXiv:2401.13996

654 [26] Merz K M, Wei G, Zhu F. Editorial: harnessing the power of large language  
655 model-based chatbots for scientific discovery[J]. *Journal of Chemical Information and*  
656 *Modeling*, 2023, 63(17): 5395.

657 [27] Stokel-Walker, C., & Van Noorden, R. (2023). What ChatGPT and generative AI

658 mean for science. *Nature*, 614(7947), 214-216.

659 [28] Chen, Y. C. (2017). A tutorial on kernel density estimation and recent advances.  
660 *Biostatistics & Epidemiology*, 1(1), 161–187.

661 [29] Daniel Han and Michael Han. 2024. Unsloth, Unsloth AI. Company name:  
662 Unsloth AI Package/product name: Unsloth.

663 [30] Chang, Y., Wang, X., Wang, J., Wu, Y., Yang, L., Zhu, K., Chen, H., Yi, X., Wang,  
664 C., Wang, Y., Ye, W., Zhang, Y., Chang, Y., Yu, P. S., Yang, Q., & Xie, X. (2024). A  
665 Survey on Evaluation of Large Language Models. *ACM Transactions on Intelligent*  
666 *Systems and Technology*, 15(3), Article 39, 1-45. <https://doi.org/10.1145/3641289>

667 [31] Coates P, Breitinger F. Identifying document similarity using a fast estimation of  
668 the levenshtein distance based on compression and signatures[J]. *arXiv.org*, 2023.

669 [32] Ehud Reiter; A Structured Review of the Validity of BLEU. *Computational*  
670 *Linguistics* 2018; 44 (3): 393–401. doi: [https://doi.org/10.1162/coli\\_a\\_00322](https://doi.org/10.1162/coli_a_00322)

671 [33] Zhang M, Li C, Wan M, et al. Rouge-sem: better evaluation of summarization  
672 using rouge combined with semantics[J]. *Expert Systems with Applications*, 2024,  
673 237: 121364.

674 [34] Hershenhouse J S, Mokhtar D, Eppler M B, et al. Accuracy, readability, and  
675 understandability of large language models for prostate cancer information to the  
676 public[J]. *Prostate Cancer Prostatic Dis*, 2024.

677 [35] Gemma Team. (2024). Gemma 2: Improving Open Language Models at a  
678 Practical Size. In *International Conference on Learning Representations (ICLR)*.  
679 OpenReview.net.

680 [36] Llama Team, AI @ Meta. (2024). The Llama 3 Herd of Models. Retrieved from  
681 <https://llama.meta.com/>

682 [37] Wu Q, Bansal G, Zhang J, et al. Autogen: enabling next-gen llm applications via  
683 multi-agent conversation[J]. *arXiv.org*, 2023.

684 [38] crewAI from crewAI - Platform for Multi AI Agents Systems

685 [39] H Chase. 2023. LangChain LLM App Development Framework. Retrieved July  
686 10, 2023 from <https://langchain.com/>

687 [40] Tavily Search API from <https://tavily.com/>



- 688 [41] Kenneth Reitz. Requests: HTTP for Humans. Available at:  
689 <https://requests.readthedocs.io/en/>.
- 690 [42] Scalfani, V.F., Patel, V.D. & Fernandez, A.M. Visualizing chemical space  
691 networks with RDKit and NetworkX. *J Cheminform* 14, 87 (2022).
- 692 [43] Xie T, Wan Y, Huang W, et al. DARWIN Series: Domain Specific Large  
693 Language Models for Natural Science[J]. arXiv.org, 2023. arXiv:2308.13565
- 694 [44] Fang Y, Liang X, Zhang N, et al. Mol-Instructions: A Large-Scale Biomolecular  
695 Instruction Dataset for Large Language Models[C]. ICLR. OpenReview.net, 2024.
- 696 [45] Jablonka K M, Schwaller P, Ortega-Guerrero A, et al. Leveraging large language  
697 models for predictive chemistry [J]. *Nature Machine Intelligence*, 2024, 6: 161–169.  
698 <https://doi.org/10.1038/s42256-023-00788-1>.
- 699 [46] Xie Z, Evangelopoulos X, Omar ÖH, Troisi A, Cooper AI, Chen L. Fine-tuning  
700 GPT-3 for machine learning electronic and functional properties of organic molecules.  
701 *Chem. Sci.* 2024, 15, 500-510. DOI: 10.1039/d3sc04610a.
- 702 [47] Zhong, S., & Guan, X. (2023). Developing Quantitative Structure–Activity  
703 Relationship (QSAR) Models for Water Contaminants’ Activities/Properties by Fine-  
704 Tuning GPT-3 Models. *Environmental Science & Technology Letters*.
- 705 [48] Kim, S., Jung, Y., & Schrier, J. (2024). Large Language Models for Inorganic  
706 Synthesis Predictions. *Journal of the American Chemical Society*, 146(29), 19654-  
707 19659. doi: 10.1021/jacs.4c05840
- 708 [49] Bran, A. M., Cox, S., Schilter, O., Baldassari, C., White, A. D., & Schwaller, P.  
709 (2024). Augmenting large language models with chemistry tools. *Nature Machine*  
710 *Intelligence*, 6, 525–535. <https://doi.org/10.1038/s42256-024-00832-8>
- 711 [50] Boiko, D.A., MacKnight, R., Kline, B. *et al.* Autonomous chemical research with  
712 large language models. *Nature* **624**, 570–578 (2023). [https://doi.org/10.1038/s41586-](https://doi.org/10.1038/s41586-023-06792-0)  
713 [023-06792-0](https://doi.org/10.1038/s41586-023-06792-0)
- 714 [51] McNaughton, A. D., Ramalaxmi, G., Kruel, A., Knutson, C. R., Varikoti, R. A.,  
715 & Kumar, N. (2024). CACTUS: Chemistry Agent Connecting Tool-Usage to Science.  
716 Preprint. <https://arxiv.org/abs/2405.00972>

717

718

719

720

721