

CATATQA: A BENCHMARK FOR TOOL-AUGMENTED LLM QUESTION ANSWERING OVER HETEROGENEOUS CATALYSIS TABLES

Anonymous authors

Paper under double-blind review

ABSTRACT

Despite their success in general question answering, large language models (LLMs) struggle with hallucinations and inaccurate reasoning in scientific domains. A major challenge stems from experimental data, which are often stored in external sources like supplementary materials and domain-specific databases. These tables are large, heterogeneous, and semantically complex, making them difficult for LLMs to interpret. While external tools show promise, current benchmarks fail to assess LLMs’ ability to navigate this data—particularly in locating relevant tables, retrieving key columns, interpreting experimental conditions, and invoking tools. To address this gap, we introduce CataTQA, a new benchmark for **C**atalytic materials **T**able **Q**uestion and **A**nswering. CataTQA features an automated dataset framework and four auxiliary tools. We evaluate tool-enhanced LLMs across five dimensions: table location, column retrieval, condition analysis, tool calling, and question answering, identifying their strengths and weaknesses. Our work sets a new benchmark for evaluating LLMs in scientific fields and paves the way for future advancements. All data and code are publicly available on GitHub¹.

1 INTRODUCTION

Large language models (LLMs) have demonstrated remarkable success in a diverse range of fields Bogin et al. (2024); Boiko et al. (2023). Prominently, they have made significant inroads into the medical field Waisberg et al. (2024); Lamb et al. (2024), the materials field Zhang et al. (2024b); Kristiadi et al. (2024), and the industrial field Yang et al. (2023); Saka et al. (2024), revolutionizing traditional practices and driving innovation. Numerous methods, including domain data fine-tuning and chain of thought, enhance LLM question-answering capabilities. However, scientific Q&A applications face specific challenges: (1) Hallucination issue Sadat et al. (2023); Weidinger et al. (2021); Ji et al. (2023); Bubeck et al. (2023): When processing experimental data stored in external sources, LLMs are highly prone to the hallucination issue, which stems from the massive scale, heterogeneous structure, and semantic complexity of data tables. This not only misleads users in decision-making but also undermines the integrity of content in rigorous scenarios like scientific research. (2) Computational limitations Gao et al. (2023); Hendrycks et al. (2021); Lewkowycz et al. (2022); Madaan & Yazdanbakhsh (2022); Nogueira et al. (2021); Qian et al. (2022): Confronted with large-scale heterogeneous tabular data, the computational architectures of LLMs reveal significant limitation, a critical skill required in numerous real-world scientific applications. (3) Inability to identify key decision-making information Ziegler et al. (2025): When confronted with multiple pieces of decision-relevant data, models frequently fail to quickly pinpoint the key information that could enhance the quality and accuracy of their responses.

Regarding the above issues, some research uses LLMs as scheduling models and enhances their Q&A capabilities by leveraging external tools, such as Retrieval-Augmented Generation (RAG) dos Santos Junior et al. (2024); Izacard et al. (2022), math tools Wu et al. (2023); Schick et al. (2023); Lu et al. (2023); Yao et al. (2023), and code interpreters Chidambaram et al. (2024); Wang et al. (2022); Gao et al. (2023). Although LLMs have shown a decent performance in general domain Q&A, their

¹<https://github.com/kg4sci/CataTQA>

application in scientific domains reveals significant limitations. Most of the experimental data in the scientific domain exists in the form of structured tables, which are characterized by large data scale, high knowledge intensity, and highly complex content. Existing tool-used evaluation methods Zhuang et al. (2023) are still unable to accurately assess the Q&A performance of large models in the scientific field, which is specifically reflected in two aspects: firstly, there is a lack of effective evaluation of the ability of LLM to quickly locate key tables and fields in massive tabular data; Secondly, a systematic evaluation of the ability of LLM to accurately call external tools for Q&A has not been established.

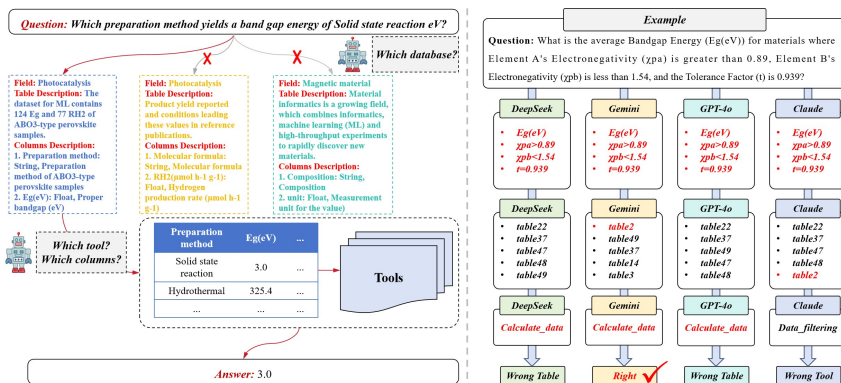


Figure 1: CataTQA has evaluated multiple capabilities of large language models, including the ability to locate key information, the ability to invoke tools, and the ability to answer questions.

To bridge this gap, we propose a Q&A benchmark named CataTQA for fairly evaluating the knowledge localization ability and tool invocation ability of LLM. CataTQA includes 68 catalytic datasets from fields like photocatalysis, with expert-annotated metadata and four defined tool types. Each instance in CataTQA consists of a question, an answer, a series of table requirements and tool for answering questions. As shown in Fig. 1, different from other benchmarks, CataTQA requires LLMs to locate the supporting dataset based on the question and invoke the correct tool for answering. This minimizes the possibility that LLM answers questions merely by recalling its internal knowledge and comprehensively evaluates its key knowledge localization ability and tool usage ability.

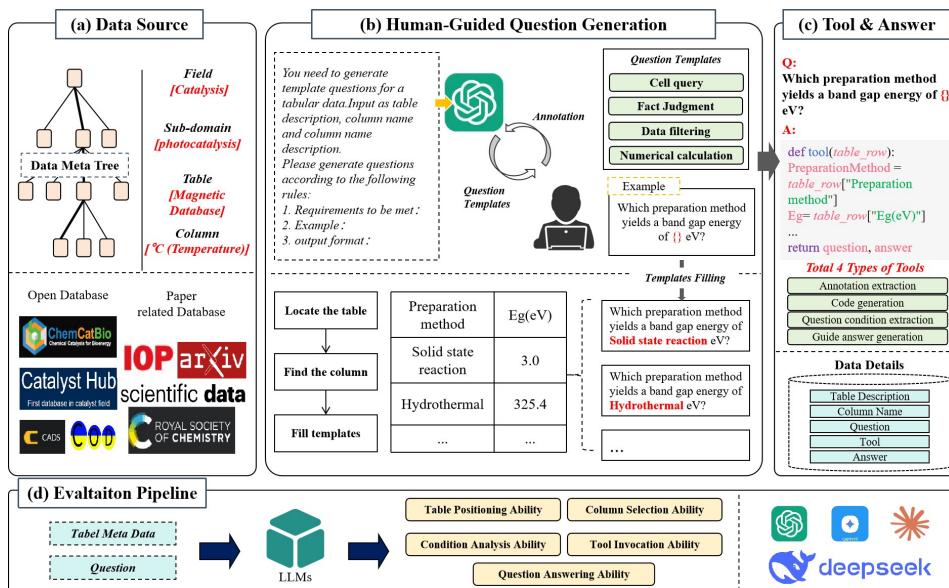


Figure 2: The Construction Process and Application Method of CataTQA, including: (a) Data Source, (b) Human-Guided Question Generation, (c) Tool & Answer and (d) Evaluation Pipeline.

The construction and application of CataTQA, as illustrated in Fig. 2, integrate expert annotation and follow a three-phase automated process: **(1) Reference Data Collection:** Based on catalytic field data systems, open-source datasets and paper-associated data are collected. Descriptive metadata is auto-gathered, raw data is preprocessed, and domain experts review/annotate the processed data. **(2) Human-Guided LLM Question Generation:** A template-based approach is used, involving human-guided template creation, validation, and question instantiation with tool attributes. Data table metadata is recorded during question generation to facilitate subsequent LLM localization ability verification. **(3) Tool-Driven Answer Generation:** Answers are derived from datasets via tool operations, ensuring external knowledge is needed. Domain expert annotations provide detailed data table descriptions to validate LLM key information localization capabilities.

Based on this benchmark, we conducted experiments using four proprietary LLMs, each with two prompting ways: Direct Prompting and Code Generation. Our findings indicate that current LLMs have limited proficiency in Q&A tasks based on multi-source tabular data, with significant performance variations across sub-tasks and question types. We decomposed the Q&A process into five key phases: Table Positioning, Column Selection, Condition Analysis, Tool Invocation, and Q&A Ability. While LLMs perform well on some sub-tasks, low accuracy in phases like tool invocation significantly reduces overall Q&A accuracy. The LLMs also show significant variability in handling different question types. For example, they struggle with Fact Judgment questions due to frequent misclassification of task types and poor tool selection.

Our research findings and analysis indicate that CataTQA is a challenging benchmark for existing tool-augmented large language model methods, especially when confronted with difficult questions in the scientific field that require quickly locating effective structured information and conducting more complex reasoning.

2 RELATED WORK

2.1 BENCHMARKS ON TOOL-AUGMENTED LLM

Recent advancements in augmenting LLMs with external tools and agents have spurred the development of benchmarks to evaluate their collaborative reasoning and tool-usage capabilities. Existing benchmarks primarily focus on either tool-augmented or table-based Q&A tasks. For instance, ToolBench Xu et al. evaluates LLMs’ ability to invoke Application Programming Interfaces (APIs) for task completion, while TAT-QA Zhu et al. (2021) and FinQA Chen et al. (2021) test numerical reasoning over structured tables. ToolQA Zhuang et al. (2023) evaluated LLMs’ ability to use external tools based on non overlapping pretraining data, and pointed out possible misuse. However, these benchmarks often operate in isolation, lacking scenarios where LLMs must dynamically coordinate multiple external agents to synthesize information across heterogeneous catalytic tables.

Efforts like WebArena Zhou et al. (2023) and ALFWorld Shridhar et al. (2020) simulate agent-environment interactions but are limited to specific domains (e.g., web navigation) and do not address cross-table reasoning. Meanwhile, benchmarks such as MMTB Koblitiz et al. (2021) introduce multi-modal table understanding but omit the integration of external tools. EHRQA Lee et al. (2022) is a practical benchmark on structured EHR data and take a step further towards bridging the gap between text-to-SQL research and its real-life deployment in healthcare. But, a critical gap remains in evaluating LLMs’ capacity to orchestrate diverse agents for tasks requiring joint analysis of interconnected tables, domain-specific knowledge retrieval, and multi-step operational workflows.

2.2 BENCHMARKS ON LLMs IN SPECIFIC SCIENCE FIELD

As the integration of artificial intelligence into different domains accelerates, LLMs are emerging as a dominant force in a growing number of applications Lu et al. (2024); Polak & Morgan (2024); Romera-Paredes et al. (2024); Birhane et al. (2023). Therefore, many benchmarks have emerged for evaluating the performance of LLMs in specified scientific fields Zhang et al. (2024a); O’Leary (2023); Singhal et al. (2023). For instance, MaScQA Zaki et al. (2024) as been classified according to the structure of the questions and subfields of materials science, with the aim of evaluating the understanding ability of LLMs regarding the key concepts of materials science. NASA-QA Bhattacharjee et al. (2024) is an extractive question answering task focused on the Earth science domain. UniDE Ye et al. (2025) represents an innovative framework. It consistently deals with multi-level

162 dialogue evaluation tasks through the utilization of data produced by LLM for training a small-scale
163 multitask evaluator.

164 165 166 3 CATATQA 167

168 169 3.1 EXPERT ANNOTATIONS 170

171 Before evaluating LLMs’ ability to answer questions by calling external tools via CataTQA, ensur-
172 ing their accurate key information localization is essential, as this capability is critical for subse-
173 quent complex tasks like data retrieval, inference, and computation. Thus, the benchmark test in-
174 corporates descriptive information from catalytic domain datasets (including source, structure, field
175 definitions, and research background) as the foundation for LLMs’ localization databases, enabling
176 them to understand and utilize the data effectively. Annotation of tables and fields in datasets is a
177 critical component. In scientific domains (e.g., materials science), specialized data and opaque col-
178 umn nomenclature pose significant comprehension challenges for non-experts, making tabular data
179 and column descriptor annotation paramount (Appendix A.2). Basic metadata annotation enhances
180 LLMs’ problem understanding and their ability to identify relevant datasets containing potential an-
181 swers. Notably, after metadata tagging, dataset expansion has minimal impact on QA efficiency—an
182 approach contrasting with systems like ChatMOF Kang & Kim (2024), which use a table searcher
183 to query all available data, significantly degrading QA efficiency.

184 During the collection of open datasets, we maintained a policy of minimal data intervention. This
185 approach differs significantly from traditional tabular question answering processing methodolo-
186 gies. Our primary objective was to preserve the authenticity of the original data, a crucial factor for
187 ensuring the reproducibility of table-based analyses. The only data processing performed was the
188 removal of entirely empty rows and columns, a procedure that further demonstrates our commitment
189 to data integrity. For dataset annotation, we engaged domain experts from the National Nanotechnol-
190 ogy Center’s catalysis field, specifically including an associate researcher and a doctoral student.
191 The annotation process covered three key aspects: 1) the research field classification of each dataset,
192 2) dataset description, and 3) descriptions of each data field.

193 194 3.2 DATASET DETAILS 195

196 CataTQA is mainly used to evaluate the comprehensive capabilities of LLMs in the vertical scientific
197 field, specifically in the field of catalysis. These capabilities include the ability to retrieve key data,
198 the ability to invoke tools, and the ability to answer questions. To this end, we have collected
199 and processed data to generate a dataset that meets the following criteria: 1) Ideally, the database
200 should not overlap with the pre-training data of the LLM (The experimental results are shown in
201 the appendix.E.1); 2) It should include accurate metadata information about the dataset, such as the
202 dataset name, dataset description, field names and their descriptions, etc.; 3) The LLM should be
203 able to obtain all the necessary information from the database to answer questions correctly.

204 CataTQA’s data comes from catalysis. To precisely evaluate table question-answering, each instance
205 includes six components: a question, answer, table, relevant columns, conditions, and tools used.
206 The dataset also contains table metadata to assess LLMs’ data location ability. It features four
207 question types with two difficulty levels based on the number of answer conditions.

208 The number of each type of problem in the CataTQA dataset and the average length of each type are
209 shown in the table 1. The four question categories are: 1) **Cell Query**, which addresses the task of
210 locating individual table cells, representing a fundamental table question answering scenario; 2) **Fact**
211 **Judgment**, judging the facts of the table, about 50% of the error conditions are generated.; 3) **Data**
212 **Filtering**, which require retrieving results that satisfy multiple specified conditions from the dataset,
213 and 4) **Numerical Calculation**, which involves questions demanding mathematical operations (e.g.,
214 median computation, summation, or maximum value identification) to derive correct answers.

215 Our study categorizes questions as simple or complex. Simple ones have 1-2 constraints and require
single-step reasoning. Complex questions contain ≥ 3 constraints, demanding multi-step LLM rea-

soning, which significantly increases difficulty. The metadata¹ and Q&A datasets² developed in this study are publicly available through the Hugging Face platform.

Table 1: CataTQA dataset information.

Question Type	Level	Average Length of Questions	Count
Cell Query	simple	12.6	4425
	complex	20.6	4339
Fact Judgment	simple	15.3	5038
	complex	21.4	5756
Data Filtering	simple	14.9	3462
	complex	19.4	1223
Numerical Calculation	simple	13.7	3325
	complex	20.6	2041

Table 2: Comparison of dataset dimensions.

Dataset	#. Tables	Size (GB)	#. Rows	#. Columns	#. Questions
ToolQA Zhuang et al. (2023)	6	5.29	4,356,045	112	1,530
ERATTA Roychowdhury et al. (2024)	7	~5	-	-	100
DataBench Grijalba et al. (2024)	65	-	3,269,975	1,615	1,300
CataTQA(ours)	68	11.62	9,021,397	1,190	29,609

In addition, we compared our previous work and found that we used more data dimensions and generated more practical problems than our previous work. Table 2 shows the detailed statistical information of CataTQA.

3.3 TOOLS DETAILS

CataTQA has systematically categorized thirteen specialized tools into four primary classes, based on characteristic features of catalysis domain data and problem typologies. The architectural design of these tools is implemented as follows:

Annotation Extraction: We implemented an annotation extraction framework consisting of three specialized tools to facilitate LLM reasoning: 1) a table data reader, 2) A labeled data extractor, and 3) a field information extractor. This toolset was specifically designed to provide structured annotation data to large language models for enhanced reasoning capabilities.

Code Generation: To evaluate the capability of LLMs in generating executable code and answering queries based on given table conditions, we developed two specialized tools: 1) a code generation module and 2) a code execution module.

Question Condition Extraction: Condition analysis is essential for accurate answers. We use four components to systematically derive problem conditions. 1) table sorting, 2) column name analysis, 3) condition extraction, and 4) entity analysis.

Guide Answer Generation: To guide answer generation while accounting for LLMs input constraints, we employ an LLM-driven tool-calling framework comprising four specialized modules: 1) data filtering, 2) data calculation, 3) tool selection, and 4) answer retrieval.

3.4 EXPERT-GUIDED QUESTION GENERATION

The aim of the question generation stage is to create questions answerable using tools in the reference corpus. There are two common approaches: human experts crafting questions based on the corpus, and LLM generating relevant questions. However, both approaches have drawbacks. Human-generated questions demand extensive manpower and time, and lack scalability. LLM-generated questions may be unanswerable, contain fabricated content, or be too simplistic, relying solely on

¹https://huggingface.co/datasets/kg4sci/CataTQA_Metadata

²<https://huggingface.co/datasets/kg4sci/CataTQA>

the LLM’s internal knowledge instead of external tools. This problem is especially prominent in scientific areas like catalysis and medicine.

To address these challenges, we introduce a hybrid question generation approach combining human guidance with LLMs, using structured templates to integrate human input and automated LLMs generation Carpuat et al. (2022); Zhang et al. (2022). Our framework leverages LLMs’ natural language understanding and expert-annotated table metadata in prompts to generate high-quality questions from tabular data. LLMs not only create meaningful problem templates aligned with table structures but also generate column names for template variables, enhancing practical question creation and experimental design. Expert verification ensures data quality in the final generated questions.

3.5 ANSWER GENERATION

To verify the impact of different question types on table Q&A, we distinguished four types of questions: Cell Query, Fact Judgment, Data Filtering, and Numerical Calculation. The question generation process in CataTQA follows a structured methodology based on expert-annotated table data and supplementary table information. We developed a standardized prompt template (Appendix C.2) that incorporates both the annotated table metadata and predefined question types. Using the Deepseek-671B DeepSeek-AI (2024) model, we generated diverse question templates along with their corresponding required field variables. All generated question templates underwent rigorous expert review and validation. Through programmed template variable instantiation, we systematically converted these templates into functional table-based questions. To ensure answer accuracy, we implemented specialized programming routines for each question type that automatically retrieve correct answers by processing the template variables.

This framework provides precise ground truth labels for all sub-tasks. While CataTQA’s main goal is to test LLMs’ question-answering ability, we also evaluate their performance in various aspects (Fig. 3), including code generation accuracy and execution results.

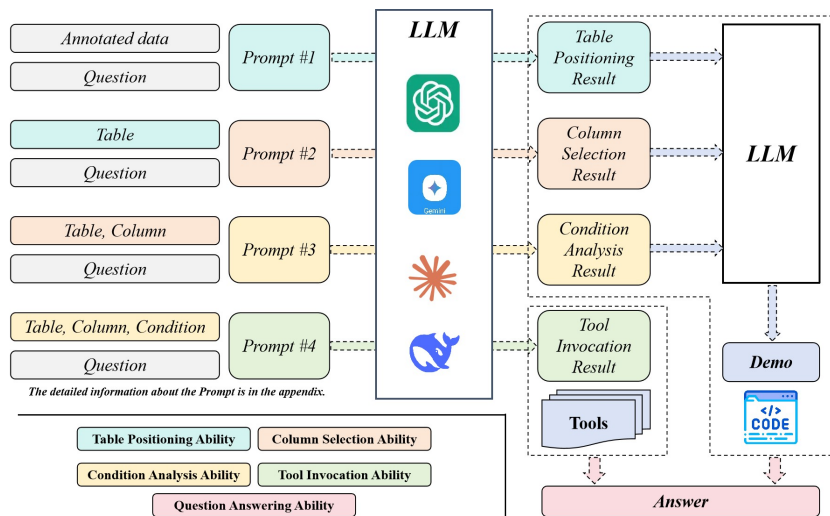


Figure 3: The evaluation process of CataTQA. We have focused on evaluating five capabilities of CataTQA, including table positioning ability, column selection ability, condition analysis ability, tool invocation ability and question answering ability.

4 EXPERIMENTAL SETUP

We experimented with four proprietary LLMs on CataTQA, GPT-4o (gpt-4o-2024-11-20) Brown et al. (2020), DeepSeek V3 (deepseek-v3-250324), Claude-3 (claude-3-haiku-20240307) and Gemini-2.5 (gemini-2.5-flash-preview-04-17). For all experiments, we used the same hyperparameters and perform 0-shot prompting via the APIs. The prompts are included in Appendix C.7. We evaluated the LLMs under two different prompting frameworks:

Direct Prompting: Given a task input, LLM is directly prompted to generate the corresponding result. We provide the corresponding prompt instructions for the four sub-tasks in Appendix C.3 to C.6. To ensure experimental fairness, we employed random sampling to select 2,000 questions per problem type, thereby maintaining consistent data conditions across all evaluated models.

Code Generation: Given a task input, LLM is prompted to generate the appropriate program in one pass. We use this framework to evaluate the basic code generation capabilities of each LLM, which is provided in Appendix C.7. For the code generation experiment, we randomly sampled 200 questions from the original set of 2,000 questions used in Direct Prompting.

We evaluated the performance of LLMs on CataTQA in four types of predefined tasks. The CataTQA differs from other benchmarks with a unique set of research challenges. We decomposed the scientific question answering process of LLMs for multiple sources into five key phases: Table Positioning Ability, Column Selection, Condition Analysis, Tool Invocation Ability and Q&A Ability. Following this pipeline, we evaluated the performance of all methods regarding accuracy scores for each sub-task and the final question answering task. The results are summarized in Table 3 and 4. To evaluate which sub-task contributes the most, we conducted a comparative study by independently removing all antecedent tasks for each sub-task. For example, in the case of condition analysis, we provide correct table positioning and column selection to individually compare the ability of all models to analyze answer conditions.

5 RESULTS AND ANALYSIS

5.1 MAIN RESULTS

Table 3 presents a comparative analysis of the performance of four LLMs (GPT, Deepseek, Claude, Gemini) across five different question types and multiple evaluation metrics. In general, Gemini appears to be the top-performing model across most question types and evaluation metrics, while other models like GPT, Deepseek, and Claude also show varying levels of competence in different aspects. Detailed analysis will be described in the discussion section.

Table 3: Performance of current LLMs.

Question Type	Model	Table Positioning Ability			Column Selection		Condition Analysis		Tool Invocation	Q&A Ability		
		ACC	ACC@3	ACC@5	ACC	F1	ACC	F1	ACC	ACC	ACC@3	ACC@5
Cell	gpt-4o-2024-11-20	0.46	0.64	0.71	0.92	0.95	0.91	0.94	0.27	0.09	0.14	0.15
	deepseek-v3-250324	0.44	0.61	0.70	0.84	0.90	0.93	0.96	0.35	0.10	0.15	0.17
Query	claude-3-haiku-20240307	0.38	0.57	0.66	0.86	0.90	0.94	0.95	0.39	0.08	0.14	0.15
	gemini-2.5-flash-preview-04-17	0.50	0.69	0.77	0.95	0.96	0.85	0.91	0.50	0.18	0.26	0.29
Fact	gpt-4o-2024-11-20	0.43	0.63	0.70	0.94	0.95	0.80	0.88	0.19	0.03	0.06	0.06
	deepseek-v3-250324	0.44	0.63	0.72	0.95	0.94	0.84	0.90	0.14	0.03	0.04	0.05
Judgment	claude-3-haiku-20240307	0.38	0.56	0.65	0.91	0.92	0.80	0.88	0.27	0.04	0.06	0.09
	gemini-2.5-flash-preview-04-17	0.50	0.71	0.78	0.95	0.96	0.75	0.85	0.25	0.06	0.09	0.09
Data	gpt-4o-2024-11-20	0.53	0.73	0.78	0.85	0.91	0.72	0.83	0.96	0.18	0.22	0.25
	deepseek-v3-250324	0.51	0.71	0.77	0.93	0.95	0.61	0.75	0.96	0.18	0.22	0.26
Filtering	claude-3-haiku-20240307	0.49	0.68	0.72	0.78	0.86	0.61	0.75	0.94	0.12	0.17	0.18
	gemini-2.5-flash-preview-04-17	0.54	0.74	0.84	0.92	0.95	0.37	0.54	0.96	0.12	0.17	0.21
Numerical	gpt-4o-2024-11-20	0.44	0.68	0.73	0.94	0.96	0.79	0.88	0.44	0.04	0.09	0.09
	deepseek-v3-250324	0.37	0.62	0.72	0.95	0.97	0.63	0.77	0.43	0.04	0.10	0.12
Calculation	claude-3-haiku-20240307	0.33	0.56	0.65	0.91	0.94	0.58	0.73	0.07	0.002	0.006	0.008
	gemini-2.5-flash-preview-04-17	0.47	0.72	0.81	0.94	0.96	0.33	0.50	0.50	0.06	0.11	0.12

We evaluated the answering ability of the LLM in two ways. The first ability is to assess whether the LLM accurately invokes the correct tool and fills in the correct fields and conditions (see Table 3); the second ability is to directly ask the LLM to generate relevant code, and then determine whether the code can be executed and whether the execution result is accurate. Table 4 shows the evaluation results of the LLM directly generating code to answer questions.

Table 4: The performance of generating code and obtaining answers.

Question Type	Assessment Task	gpt-4o	deepseek-v3	claude-3	gemini-2.5
		2024-11-20	250324	haiku-20240307	flash-preview-04-17
Cell Query	Code Executable	0.64	<u>0.70</u>	0.63	0.84
	Correct Answer	0.61	0.59	<u>0.60</u>	0.52
Fact Judgment	Code Executable	<u>0.81</u>	0.79	0.55	0.97
	Correct Answer	<u>0.47</u>	0.52	0.35	0.22
Data Filtering	Code Executable	<u>0.91</u>	0.94	0.83	0.87
	Correct Answer	0.67	<u>0.64</u>	0.52	0.47
Numerical Calculation	Code Executable	0.87	0.81	0.72	<u>0.82</u>
	Correct Answer	0.41	<u>0.38</u>	0.35	0.31

5.2 TABLE POSITIONING ABILITY

From the Table 3 and Fig. 4(a), we draw the following observations: (1) In categories such as Electrocatalysis, Magnetic materials, and Perovskite, LLMs can accurately locate the correct background knowledge data. However, in the "Others" category, LLMs are more likely to mis-locate data from other fields. This may be because the tables in this category have low specificity, causing interference to the LLMs. (2) Regardless of the type of question, the data-location capabilities of LLMs do not vary significantly. This indicates that LLMs consider the content of questions more than the question categories when selecting data. (3) Overall, Gemini demonstrates the best data-location capabilities, while Claude performs poorly. As shown in Fig. 6 in Appendix, alluvial plot illustrates the table positing ability of different LLMs.

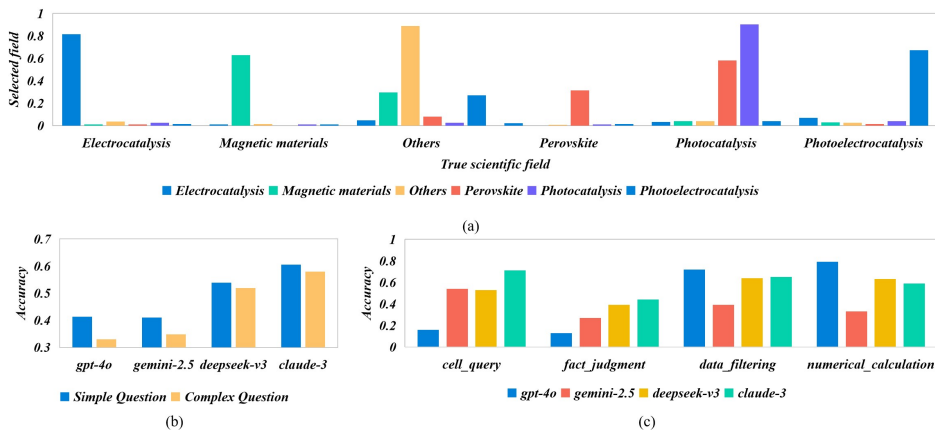


Figure 4: (a) Statistical Analysis of the Table Positioning Ability of LLMs, (b) Comparison of condition analysis accuracy of different LLMs on questions with different difficulty levels, (c) Comparison of condition analysis accuracy of different LLMs on questions with different type. capabilities

5.3 COLUMN SELECTION ABILITY&CONDITION ANALYSIS ABILITY

Table 3 shows that regardless of which LLM it is, its performance in field selection is superior to other capabilities, and among these models, Gemini performs the best. It can also be easily inferred that question difficulty has little impact, as the process relies more on semantic information than complex reasoning.

As shown Fig. 4(b) & (c), in terms of conditional analysis capabilities, it can be observed that the performance of Claude-3 is better than that of other models, which may be attributed to its strong logical reasoning ability. Moreover, it can be seen that the difficulty level of the questions has a certain impact on this ability. This is mainly because as the complexity of the questions increases,

the number of conditions also increases accordingly, posing challenges to LLMs. And the type of question also affects the accuracy of conditional judgment. For most LLMs, in problems involving numerical calculations and data filtering, they have demonstrated relatively excellent results in conditional analysis capabilities. Finally, we also conducted experiments on the scope of tool interaction, and the experimental results are presented in D.2.

5.4 TOOL INVOCATION ABILITY

From the results, it can be seen that in the type of questions related to Data Filtering, all large models can effectively invoke the corresponding tools. However, in other types of questions, the performance of large models is not satisfactory. Fig. 5 shows the confusion matrix of the tool invocation results of large models. All LLMs tend to confuse the Data Filtering tool with other tools. This may be due to the inductive summarization ability of large models, resulting in poor performance in questions related to data calculation, situation judgment, etc. Questions of inductive summarization tend to call the Data Filtering tool rather than other types of tools. Except for Claude-3, other models can effectively use the calculate data tool to answer questions.

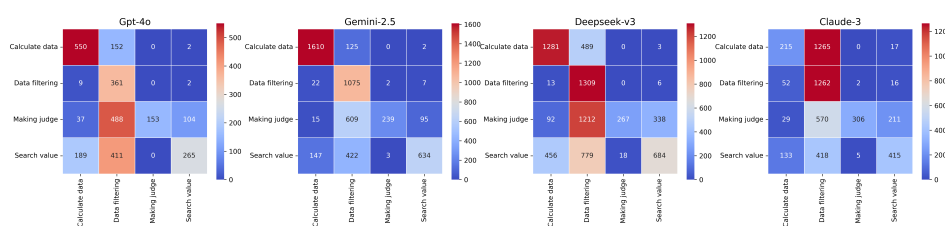


Figure 5: Comparison of Confusion Matrices of Tool Invocation Capabilities from Different LLMs

6 CONCLUSION AND LIMITATIONS

As one of the most intuitive representations of scientific data, tabular data contains a wealth of latent domain knowledge. Making full use of open domain datasets can assist scientists in better exploring the development context of the field and achieving scientific discoveries. Additionally, it facilitates the development of specialized chatbots that can generate answers based on table content. Therefore, we developed CataTQA, designed to evaluate LLMs’ ability to accurately understand background knowledge in scientific domains, precisely locate valid information within vast domain datasets, and invoke external tools to solve complex problems. The development process of CataTQA is highly generalizable, enabling seamless migration to other scientific disciplines. We conducted extensive tests on CataTQA using GPT, DeepSeek, Claude and Gemini, and the results revealed that even the strongest model achieved limited performance on the hard questions of CataTQA.

However, CataTQA has significant limitations. Firstly, it currently doesn’t support cross-table associative retrieval (We have made some attempts in our research in Appendix.D.1.). As a result, the dataset can hardly simulate the complex problems in real scientific research, and it may underestimate the models’ ability to handle associated data. Secondly, CataTQA is mainly focused on the catalyst field at present. To enhance its comprehensiveness, it’s necessary to expand the dataset to cover more disciplines like astronomy, genomics, etc. This tests models’ ability to handle diverse data and scientific problems, enabling a more comprehensive evaluation of generalization. Finally, recent advancements in scientific annotation tools have significantly enhanced both the quality and efficiency of metadata annotation. Representative systems such as Autodive Du et al. (2023) demonstrate this progress by enabling automated annotation of tabular data within PDF-formatted scientific papers, thereby facilitating rapid extraction of scientific table data. Our approach currently lacks integration with established annotation tools, which could better reflect real-world scenarios. In the future, efforts will be made to address these shortcomings and strengthen the role of CataTQA as a benchmark for evaluating large language models in scientific applications.

REFERENCES

- 486
487
488 Bishwaranjan Bhattacharjee, Aashka Trivedi, Masayasu Muraoka, Muthukumaran Ramasubramanian, Takuma Udagawa, Iksha Gurung, Nishan Pantha, Rong Zhang, Bharath Dandala, Rahul
489 Ramachandran, et al. Indus: Effective and efficient language models for scientific applications.
490 *arXiv preprint arXiv:2405.10725*, 2024.
491
- 492 Abeba Birhane, Atoosa Kasirzadeh, David Leslie, and Sandra Wachter. Science in the age of large
493 language models. *Nature Reviews Physics*, 5(5):277–280, 2023.
494
- 495 Ben Bogin, Kejuan Yang, Shashank Gupta, Kyle Richardson, Erin Bransom, Peter Clark, Ashish
496 Sabharwal, and Tushar Khot. Super: Evaluating agents on setting up and executing tasks from
497 research repositories. *arXiv preprint arXiv:2409.07440*, 2024.
- 498 Daniil A Boiko, Robert MacKnight, Ben Kline, and Gabe Gomes. Autonomous chemical research
499 with large language models. *Nature*, 624(7992):570–578, 2023.
- 500 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
501 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
502 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
503
- 504 Sébastien Bubeck, Varun Chadrsekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar,
505 Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence:
506 Early experiments with gpt-4, 2023.
- 507 Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (eds.). *Proceedings*
508 *of the 2022 Conference of the North American Chapter of the Association for Computational*
509 *Linguistics: Human Language Technologies*, Seattle, United States, July 2022. Association for
510 Computational Linguistics. URL [https://aclanthology.org/2022.naacl-main.](https://aclanthology.org/2022.naacl-main.0/)
511 [0/](https://aclanthology.org/2022.naacl-main.0/).
- 512 Zhiyu Chen, Wenhu Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema
513 Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, et al. Finqa: A dataset of numerical
514 reasoning over financial data. *arXiv preprint arXiv:2109.00122*, 2021.
- 515 Subramanian Chidambaram, Li Erran Li, Min Bai, Xiaopeng Li, Kaixiang Lin, Xiong Zhou, and
516 Alex C Williams. Socratic human feedback (sohf): Expert steering strategies for llm code gener-
517 ation. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 15491–
518 15502, 2024.
- 519 DeepSeek-AI. Deepseek-v3 technical report, 2024. URL [https://arxiv.org/abs/2412.](https://arxiv.org/abs/2412.19437)
520 [19437](https://arxiv.org/abs/2412.19437).
- 521
522 José Cassio dos Santos Junior, Rachel Hu, Richard Song, and Yunfei Bai. Domain-driven llm de-
523 velopment: Insights into rag and fine-tuning practices. In *Proceedings of the 30th ACM SIGKDD*
524 *Conference on Knowledge Discovery and Data Mining*, pp. 6416–6417, 2024.
- 525 Yi Du, Ludi Wang, Mengyi Huang, Dongze Song, Wenjuan Cui, and Yuanchun Zhou. Autodive: An
526 integrated onsite scientific literature annotation tool. In *Proceedings of the 61st Annual Meeting*
527 *of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pp. 76–85,
528 2023.
529
- 530 Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and
531 Graham Neubig. Pal: Program-aided language models. In *International Conference on Machine*
532 *Learning*, pp. 10764–10799. PMLR, 2023.
- 533 Jorge Osés Grijalba, L Alfonso Urena Lopez, Eugenio Martínez-Cámara, and Jose Camacho-
534 Collados. Question answering over tabular data with databench: A large-scale empirical eval-
535 uation of llms. In *Proceedings of the 2024 Joint International Conference on Computational Lin-*
536 *guistics, Language Resources and Evaluation (LREC-COLING 2024)*, pp. 13471–13488, 2024.
537
- 538 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
539 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv*
preprint arXiv:2103.03874, 2021.

- 540 Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane
541 Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. Few-shot learning with re-
542 trieval augmented language models. *arXiv preprint arXiv:2208.03299*, 1(2):4, 2022.
- 543
- 544 Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang,
545 Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM*
546 *computing surveys*, 55(12):1–38, 2023.
- 547 Yeonghun Kang and Jihan Kim. Chatmof: an artificial intelligence system for predicting and gen-
548 erating metal-organic frameworks using large language models. *Nature communications*, 15(1):
549 4705, 2024.
- 550
- 551 Julia Koblitz, Sabine Eva Will, S Alexander Riemer, Thomas Ulas, Meina Neumann-Schaal, and
552 Dietmar Schomburg. The metano modeling toolbox mmtb: an intuitive, web-based toolbox intro-
553 duced by two use cases. *Metabolites*, 11(2):113, 2021.
- 554 Agustinus Kristiadi, Felix Strieth-Kalthoff, Marta Skreta, Pascal Poupart, Alán Aspuru-Guzik, and
555 Geoff Pleiss. A sober look at llms for material discovery: Are they actually good for bayesian
556 optimization over molecules? *arXiv preprint arXiv:2402.05015*, 2024.
- 557
- 558 Bernadette Lamb, Jonathan Herskovitz, Marta Jonson, Harlan Sayles, and Faruq Pradhan. S2185
559 use of large language model (llm) chatbots for the generation of specialized gastrointestinal diet
560 meal plans: A pilot study. *Official journal of the American College of Gastroenterology— ACG*,
561 119(10S):S1562, 2024.
- 562 Gyubok Lee, Hyeonji Hwang, Seongsu Bae, Yeonsu Kwon, Woncheol Shin, Seongjun Yang, Min-
563 joon Seo, Jong-Yeup Kim, and Edward Choi. Ehrsql: A practical text-to-sql benchmark for
564 electronic health records. *Advances in Neural Information Processing Systems*, 35:15589–15601,
565 2022.
- 566
- 567 Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ra-
568 masesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative
569 reasoning problems with language models. *Advances in Neural Information Processing Systems*,
570 35:3843–3857, 2022.
- 571 Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scient-
572 ist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*,
573 2024.
- 574
- 575 Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu,
576 and Jianfeng Gao. Chameleon: Plug-and-play compositional reasoning with large language mod-
577 els. *Advances in Neural Information Processing Systems*, 36:43447–43478, 2023.
- 578 Aman Madaan and Amir Yazdanbakhsh. Text and patterns: For effective chain of thought, it takes
579 two to tango. *arXiv preprint arXiv:2209.07686*, 2022.
- 580
- 581 Rodrigo Nogueira, Zhiying Jiang, and Jimmy Lin. Investigating the limitations of transformers with
582 simple arithmetic tasks. *arXiv preprint arXiv:2102.13019*, 2021.
- 583
- 584 Karen O’Leary. Llms get a medical education. *Nature Medicine*, 2023.
- 585 Maciej P Polak and Dane Morgan. Extracting accurate materials data from research papers with
586 conversational language models and prompt engineering. *Nature Communications*, 15(1):1569,
587 2024.
- 588
- 589 Jing Qian, Hong Wang, Zekun Li, Shiyang Li, and Xifeng Yan. Limitations of language models in
590 arithmetic and symbolic induction. *arXiv preprint arXiv:2208.05051*, 2022.
- 591 Bernardino Romera-Paredes, Mohammadamin Barekatin, Alexander Novikov, Matej Balog,
592 M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang,
593 Omar Fawzi, et al. Mathematical discoveries from program search with large language models.
Nature, 625(7995):468–475, 2024.

- 594 Sohini Roychowdhury, Marko Krema, Anvar Mahammad, Brian Moore, Arijit Mukherjee, and Punit
595 Prakashchandra. Eratta: Extreme rag for enterprise-table to answers with large language models.
596 In *2024 IEEE International Conference on Big Data (BigData)*, pp. 4605–4610. IEEE, 2024.
597
- 598 Mobashir Sadat, Zhengyu Zhou, Lukas Lange, Jun Araki, Arsalan Gundroo, Bingqing Wang,
599 Rakesh R Menon, Md Rizwan Parvez, and Zhe Feng. Delucionqa: Detecting hallucinations
600 in domain-specific question answering. *arXiv preprint arXiv:2312.05200*, 2023.
- 601 Abdullahi Saka, Ridwan Taiwo, Nurudeen Saka, Babatunde Abiodun Salami, Saheed Ajayi, Kabiru
602 Akande, and Hadi Kazemi. Gpt models in construction industry: Opportunities, limitations, and
603 a use case validation. *Developments in the Built Environment*, 17:100300, 2024.
- 604 Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro,
605 Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can
606 teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:68539–
607 68551, 2023.
- 608
- 609 Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew
610 Hausknecht. Alfworld: Aligning text and embodied environments for interactive learning. *arXiv
611 preprint arXiv:2010.03768*, 2020.
- 612 Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan
613 Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. Large language models encode
614 clinical knowledge. *Nature*, 620(7972):172–180, 2023.
- 615
- 616 Ethan Waisberg, Joshua Ong, Mouayad Masalkhi, and Andrew G Lee. Large language model (llm)-
617 driven chatbots for neuro-ophthalmic medical education. *Eye*, 38(4):639–641, 2024.
- 618 Xingyao Wang, Sha Li, and Heng Ji. Code4struct: Code generation for few-shot structured predic-
619 tion from natural language. *arXiv preprint arXiv:2210.12810*, 3, 2022.
- 620
- 621 Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang,
622 Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, et al. Ethical and social risks of harm
623 from language models. *arXiv preprint arXiv:2112.04359*, 2021.
- 624 Yiran Wu, Feiran Jia, Shaokun Zhang, Hangyu Li, Erkang Zhu, Yue Wang, Yin Tat Lee, Richard
625 Peng, Qingyun Wu, and Chi Wang. Mathchat: Converse to tackle challenging math problems
626 with llm agents. *arXiv preprint arXiv:2306.01337*, 2023.
- 627 Qiantong Xu, Fenglu Hong, Bo Li, Changran Hu, Zhengyu Chen, and Jian Zhang. On the tool
628 manipulation capability of open-sourced large language models. In *NeurIPS 2023 Foundation
629 Models for Decision Making Workshop*.
- 630
- 631 Fangkai Yang, Pu Zhao, Zezhong Wang, Lu Wang, Jue Zhang, Mohit Garg, Qingwei Lin, Saravan
632 Rajmohan, and Dongmei Zhang. Empower large language model to perform better on industrial
633 domain-specific question answering. *arXiv preprint arXiv:2305.11541*, 2023.
- 634 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
635 React: Synergizing reasoning and acting in language models. In *International Conference on
636 Learning Representations (ICLR)*, 2023.
- 637
- 638 Guanghui Ye, Huan Zhao, Zixing Zhang, and Zhihua Jiang. Unide: A multi-level and low-resource
639 framework for automatic dialogue evaluation via llm-based data augmentation and multitask
640 learning. *Information Processing & Management*, 62(3):104035, 2025.
- 641 Mohd Zaki, NM Anoop Krishnan, et al. Mascqa: investigating materials science knowledge of large
642 language models. *Digital Discovery*, 3(2):313–327, 2024.
- 643
- 644 Di Zhang, Wei Liu, Qian Tan, Jingdan Chen, Hang Yan, Yuliang Yan, Jiatong Li, Weiran Huang,
645 Xiangyu Yue, Wanli Ouyang, et al. Chemllm: A chemical large language model. *arXiv preprint
646 arXiv:2402.06852*, 2024a.
- 647 Huan Zhang, Yu Song, Ziyu Hou, Santiago Miret, and Bang Liu. Honeycomb: A flexible llm-based
agent system for materials science. *arXiv preprint arXiv:2409.00135*, 2024b.

648 Rongzhi Zhang, Yue Yu, Pranav Shetty, Le Song, and Chao Zhang. Prboost: Prompt-based rule dis-
649 covery and boosting for interactive weakly-supervised learning. *arXiv preprint arXiv:2203.09735*,
650 2022.

651 Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng,
652 Tianyue Ou, Yonatan Bisk, Daniel Fried, et al. Webarena: A realistic web environment for build-
653 ing autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023.

654 Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng,
655 and Tat-Seng Chua. Tat-qa: A question answering benchmark on a hybrid of tabular and textual
656 content in finance. *arXiv preprint arXiv:2105.07624*, 2021.

657 Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and Chao Zhang. Toolqa: A dataset for llm
658 question answering with external tools. *Advances in Neural Information Processing Systems*, 36:
659 50117–50143, 2023.

660 Andreas Ziegler, David Joseph, Thomas Gossard, Emil Moldovan, and Andreas Zell. Biasbench: A
661 reproducible benchmark for tuning the biases of event cameras. *arXiv preprint arXiv:2504.18235*,
662 2025.

663 A DATA SOURCE

664 A.1 DATA SOURCE

665 We have retrieved data from multiple open-source databases. The following lists the top 9 categories
666 of data with large volumes. For the complete data sources, please refer to the tabular metadata in
667 our GitHub repository.

- 668 • **Catalyst Acquisition by Data Science(CADS)**¹: An innovative web-based integrated cat-
669 alyst informatics platform, Catalyst Acquisition by Data Science (CADS), is developed for
670 use towards the discovery and design of catalysts.
- 671 • **Catalytic Material Database(CMD)**²: CMD contains material composition, properties,
672 reactions, products and other information.
- 673 • **Catalyst Hub**³: A featured database for surface reactions contains more than 100,000
674 chemisorption and reaction energies obtained from electronic structure calculations, and is
675 continuously being updated with new datasets.
- 676 • **Crystallography Open Database(COD)**⁴: An open-access collection of crystal structures
677 of organic, inorganic, metal-organic compounds and minerals.
- 678 • **Materials Project**⁵: The Materials Project provides computed information on known and
679 predicted materials as well as powerful analysis tools to inspire and design novel materials.
- 680 • **2DMatPedia dataset**⁶: DMatPedia dataset is a collection of 2D materials, contains 6351
681 materials.
- 682 • **Alexandria_DB PBE 3D**⁷: A dataset of 2.5m+ stable and metastable materials calculated
683 with the PBE functional.
- 684 • **OQMD-3D dataset**⁸: The OQMD is a database of DFT calculated thermodynamic and
685 structural properties of 1,226,781 materials, created in Chris Wolverton’s group at North-
686 western University.
- 687 • **Associated Data of Papers**: We have also obtained the data related to the paper. For the
688 specific information of the paper, please refer to the project on GitHub.

695 ¹<https://cads.eng.hokudai.ac.jp/>

696 ²<http://cmd.us.edu.pl/catalog/>

697 ³<http://www.catalysthub.net/>

698 ⁴<http://www.crystallography.net/cod/>

699 ⁵<https://next-gen.materialsproject.org/>

700 ⁶<http://www.2dmatpedia.org/>

701 ⁷<https://alexandria.icams.rub.de/>

⁸<https://www.oqmd.org/download/>

702 A.2 ANNOTATION INSTANCE
703

704 We annotated each dataset as shown below.
705

706
707
708 **tabular name:** table1

709
710 **tabular description:** This table evaluates Mo-doped BaTiO₃ photocatalysts for hydrogen produc-
711 tion, correlating doping levels with bandgap (E_g) and activity (RH₂). It identifies optimal
712 Mo content (2%) for peak efficiency, aiding material design for enhanced solar-driven wa-
713 ter splitting. Contributions include optimizing dopant ratios to balance light absorption
714 and charge separation.

715 **field description:** 1.[Molecular formula] - [elemental composition of compounds]
716 2.[RH₂(μmol h⁻¹ g⁻¹)] - [Hydrogen production rate]
717 3.[E_g(eV)] - [Band gap energy]
718 4.[Preparation method] - [Material synthesis technique]
719 5.[Calcination temperature(K)] - [Heating temperature during synthesis]
720 6.[Calcination time(h)] - [Heating duration during synthesis]
721 7.[Light intensity(W)] - [Light source power used]
722 8.[Reaction solution] - [Chemical solution in reaction]
723 9.[Co-catalyst] - [Catalyst used with photocatalyst]
724 10.[Photocatalyst dose(g L⁻¹)] - [Photocatalyst concentration used]
725 11.[Ref] - [Reference source citation]
726

727
728
729 A.3 DATA INSTANCE
730

731
732
733 **question:** Find the partial pressure of argon (Par) when the methane partial pressure (Pch₄) is 0.3,
734 oxygen partial pressure (Po₂) is 0.06, and temperature is 850°C.
735

736 **refer_dataset:** table28

737
738 **column names:** {Pch₄, Po₂, Temperature, Par}

739
740 **condition_column:** {Pch₄, Po₂, Temperature}

741
742 **answer_column:** {Par}

743
744 **condition:** {Pch₄: 0.3, Po₂: 0.06, Temperature: 850}

745
746 **tool:** search_value

747
748 **answer: Par:** 0.6

749
750 **level:** *complex*

751 **question description:** In a tabular data structure, locate the cells that meet the requirements.

752
753 **refer_template:** Find the partial pressure of argon (Par) when the methane partial pressure (Pch₄)
754 is {}, oxygen partial pressure (Po₂) is {}, and temperature is {}°C.
755

756 B QUESTION TEMPLATES

757

758

759

B.1 CELL QUERY TEMPLATES

760

761

Cell Query Templates

762

763

764

765

766

767

768

769

- Find the Pauling electronegativity of element B (χ_{pb}) in {}.
- How much is the molar mass (M) of the compound named {}?
- What license type applies to the paper authored by {} and last updated on {}?

770

771

B.2 FACT JUDGMENT TEMPLATES

772

773

Fact Judgment Templates

774

775

776

777

778

779

780

781

782

783

B.3 DATA FILTERING TEMPLATES

784

785

786

787

Data Filtering Templates

788

789

790

791

792

793

794

795

796

797

- For support material {}, which M1 elements paired with M2 {} yield H2 over {}% and O2 conversion under {}%?
- Find compounds in sub_family with axial distortion exceeding {} and slope_CSM under {}.
- Show materials with phase separation energy ; {} and atoms matching lattice type {}?

B.4 NUMERICAL CALCULATION TEMPLATES

798

799

800

801

802

803

804

805

806

807

808

809

Numerical Calculation Templates

- What is the minimum e_hull value for materials with bulk modulus exceeding {}?
- Calculate the average PBE0 open-circuit voltage for entries with PBE0 LUMO energy level greater than {}.
- What is the highest First Ionization Energy of B (I1b) observed in materials with EAa greater than {}, EAb less than {}, and ρ_a exceeding {}?

810 C PROMPTS

811 C.1 PROMPTS FOR TABLE ANNOTATION

812 <Table Annotation>Prompt

813 **Annotate the table data, summarize the main problems that this table can solve and its**
 814 **contributions based on the content of the table data.**

815 **Table data:**

816 **{dataset}**

817 **Output requirements:**

- 818 **1. Summarize the role of data and avoid discussing a single column or row of data**
- 819 **2. Output is limited to 50 words or less**

820 C.2 PROMPTS FOR TEMPLATE QUESTIONS GENERATION

821 <Template Questions Generation>Prompt

822 **Annotate the table data, summarize the main problems that this table can solve and its**
 823 **contributions based on the content of the table data.**

824 **Please generate questions according to the following rules:**

825 **1. Requirements to be met:**

- 826 **- template questions type: {question_description}**
- 827 **- number of columns required to obtain answers: at least two columns**
- 828 **- level: The level of the template questions is differentiated according to the number of**
 829 **columns used.**

830 **Including two levels of simple and complex.**

831 **2. Example:**

832 **Input:**

833 **- table description: {example_tabular_description}**

834 **- [column names] - [description]: {example_field_description}**

835 **Output: {example}**

836 **3. output format:**

- 837 **- Mark the level of each question. At least ten questions per level.**
- 838 **- Mark the column names that need to be used to answer this question template.**
- 839 **- Use "{ }" for template variables. The template variable must be one of the columns**
 840 **of the table.**
- 841 **- Use of multiple sentence structures. Questions need to be phrased in a way that is easy**
 842 **to understand.**

843 **Use the information in the table below to generate template questions according to the**
 844 **above rules:**

845 **Input:**

846 **- table description: {tabular_description}**

847 **- [column names] - [description]: {field_description}**

850 C.3 PROMPTS FOR TABLE ORDER

851 <Table Order>Prompt

852 **Please analyze the relevance of the table according to the problem.**

853 **question: {question}**

854 **Available tables and table descriptions:**

855 **{table_desc}**

856 **Please sort the tables by relevance from high to low, give the first five possible tables,**
 857 **and directly return the table name list.**

858 **For example: ["table1", "table2"]**

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

C.4 PROMPTS FOR COLUMN SELECTION

<Column Selection>Prompt

Only provide the column names required to answer the question:
question: {question}
{table} information: [column]-[column name explanation]
{table_field}
directly return to the column name list, for example: ['col1','col2']

C.5 PROMPTS FOR CONDITION EXTRACTION

<Condition Extraction>Prompt

Please extract the query criteria from the question and return the results according to the table structure:
question: {question}
table information:
column name: [{tar}]
Output requirements:
1. Return to dictionary format, with the key being the column name and the value being the query condition
2. The values corresponding to all column names must be output. The output not found in the problem is "".
3. Extract and preserve comparison symbols (<, >, =, etc.)
4. example: "column1": "<50", "column2": "Liming"
Please return the JSON dictionary directly without including any other content

C.6 PROMPTS FOR TOOL INVOCATION

<Tool Invocation>Prompt

Please use the tools needed to answer the questions according to the question analysis. You need to specify the calculated column name when you need to perform calculation, but you do not need to specify it when you are performing table lookup.
Question: {question}
Description of available tools:
{tool_desc}
Please return the tool and column name to be used directly. For example:
{{"tool": "", "colnum name": ""}}

C.7 PROMPTS FOR GENERATING CODE

<Generating Code>Prompt

Generate python code based on the question and given conditions, index CSV file data, and answer the question.

Generate a function named "get_answer"(No parameters required). The function must use the "return" keyword to return the variable "answer", which is the answer to the question.

question: {question['question']}

refer_dataset: {question['refer_dataset']}

column names: {question['column names']}

condition: {question['condition']}

Code must be used in markdown format("python").

Do not return redundant content. The returned results must be saved in the "answer" variable.

D SUPPLEMENTARY EXPERIMENTS

D.1 CROSS-TABLE REASONING

In the field of materials science, there is a need to perform reasoning across table information to address more complex questions and derive answers. In this paper, while constructing a benchmark for single tables, we also attempted to build a multi-table association reasoning task. We migrated the current pipeline to multi-table task reasoning, using metadata and table descriptions as additional inputs to prompt LLMs to generate cross-table task question templates. In total, we generated 426 cross-table questions and invited domain experts to evaluate the rationality and complexity of these questions. However, after expert evaluation, only 37.5% of the questions were satisfactory. Considering the value and accuracy of the benchmark comprehensively, we did not include them in CataTQA, and we also described this part in the "limitation" section of the paper. In the future, we will continue to update the pipeline and further expand the evaluation data according to the characteristics of cross-table reasoning tasks.

Nevertheless, as can be seen from the evaluation results in this paper, even in reasoning for single tables, the performance of mainstream LLMs still has significant room for improvement. Therefore, we believe that this work can serve as a new benchmark for evaluating whether large models can accurately understand information in the catalytic field, locate tables and conditions, and invoke the correct tools.

Cross-Table Reasoning Examples:

- **Question 1: For materials with a molecular formula of {table1.Molecular formula} and a bandgap energy of {table2.Eg(eV)}, what is the co-catalyst used and the first ionization energy of element B?**

Annotation from expert: The two tables differ in their research objects and data structures, so directly correlating molecular formulas with band gap energies may lead to data confusion. In addition, molecular formulas and band gap energies belong to the basic properties of materials, while Cocatalysts and the ionization energy of element B fall under process and elemental attributes. There is no direct scientific or database logical connection between these two categories.

- **Question 2: What is the effect of the light type {table4.Light Type} on the bandgap energy {table2.Eg(eV)} and the density of element A {table2.a} in materials with a primary element A of {table4.A}?**

Annotation from expert: This question has a muddled causal relationship. The type of light (table4.Light Type) is an external experimental condition, while the band gap energy (table2.Eg(eV)) and element density (table2.a) are intrinsic properties of the material. The type of light does not directly affect the band gap or element density of the material, and there is no scientific causal relationship between them.

- **Question 3: What is the correlation between the publication year {table12.Year} of studies on catalysts with surface strain {table12.Surface Strain} and the number of times they were cited {table13.Cited_Time}?**

Annotation from expert: The correlation between the different variables proposed in this question lacks scientific significance. The publication year is a temporal attribute of a paper, while surface strain is a physical property of materials; there is no direct scientific connection between the two. The number of citations a paper receives is mainly influenced by factors such as research impact and field hotspots, and has no direct causal relationship with the surface strain of materials.

D.2 SCOPE OF TOOL INTERACTION

We have evaluated and analyzed the accuracy or failure modes of tool components to conduct a more comprehensive assessment of the capabilities of large language models.

For the tasks of Table Positioning Ability and Column Selection, LLM can output in the required format without discovering any other faults, due to our limited output constraints. LLM has excellent ability to complete the answers in these two tasks. For the tasks of Condition Analysis and Tool Invocation, we require LLM to output a more rigorous JSON format. Our diagnostic analysis revealed that 31% of failures stemmed exclusively from JSON formatting errors in the LLM's output. While this represents an improvement over previous format requirements, we acknowledge this remains a significant error source.

For our primary research focus on assessing LLMs' tabular reasoning capabilities. We intentionally excluded formatting errors from the final accuracy calculations. Only properly formatted outputs were evaluated for task performance. This approach ensures our metrics reflect the LLM's core table reasoning abilities rather than format compliance. This methodological choice aligns with our research objective to isolate and measure the model's table processing competencies, while recognizing formatting as an orthogonal challenge.

In the QA Ability task, we conducted a detailed code check and, based on our task division, LLM deduced the correct preconditions. By calling relevant tools, we were able to almost complete the correct answer retrieval.

We also further analyzed the reasons for the incorrect invocation of Data Filtering tool, and the results of the analysis are as follows:

As can be seen from Fig.5, when facing all types of problems, LLMs often incorrectly call data filtering tools to answer the questions. We further analyzed the possible causes for different types of problems.

For questions that require calling the calculate_data tool for answers, LLMs often incorrectly invoke the data_filtering tool because they fail to understand the calculation instructions in the questions (such as minimum, maximum, average, etc.). For example, regarding the question "Determine the minimum Fermi energy level (efermi) for materials where SCF valence band maximum (scf_vbm) is greater than 5.176." the large model failed to accurately comprehend "minimum efermi" when analyzing the required conditions, and only generated the condition information "llm_condition": "scf_vbm": ">5.176", "efermi": """. Therefore, it chose to call the data_filtering tool to filter the table data.

For questions that require invoking the makin_judge tool for answers, LLMs incorrectly call the data_filtering tool because they fail to fully understand that the instructions demand a judgmental response of "yes or no". This may be due to the fact that data-related information in the question's syntax is positioned later, thereby causing the misunderstanding. For instance, in the question "Is there a material where M06 method reports both HOMO of -0.235, LUMO of -0.084, and Scharber efficiency of 0.00034537252143?", the LLMs focused heavily on the information about "HOMO", "LUMO", and "Scharber efficiency" while ignoring the question's intent "Is there a material", leading to the invocation of the data_filtering tool.

In the analysis of tool invocation confusion in questions that require the search_value tool, we found that the issue of incorrect tool calls by large models may lie in the ambiguity in their judgment of column name selection. For example, in the question "What is the chemical formula for the material

containing the element '-C-H?', the large model mistakenly identified "containing the element '-C-H'" as a data filtering condition, and thus invoked the data_filtering tool to answer the question.

It can be seen from the results that current LLMs still have significant room for improvement in the correctness of tool calls in such specialized fields as catalysis. Therefore, we believe that CataTQA can serve as a benchmark for evaluating the application effects of LLMs in practical scientific research fields.

E OTHER INFORMATION

E.1 VERIFY THE ORIGIN OF THE ANSWER

To verify whether LLMs can answer relevant questions based on their own memorization, we conducted supplementary experiments. We directly input the questions into the large model and asked it to provide answers, then evaluated the accuracy of its responses. As can be seen from the examples below, the LLMs can not answer the questions accurately. Therefore, It can prove to a certain extent that this benchmark is mainly used to test the reasoning ability of LLMs rather than their memorization, in the catalytic field.

We use the model 'gpt-4o-2024-11-20' and the following prompt:

<Verify the origin of the answer >Prompt

Please provide the answer directly to the question without giving any explanation.

question:{question}

The output format is: ``json{{"answer": ""}}``

Table 5: Result of verify the origin of the answer.

Question Type	Level	ACC	AVG ACC
Cell Query	simple	0.004	0.002
	complex	0	
Fact Judgment	simple	0.253	0.236
	complex	0.217	
Data Filtering	simple	0	0
	complex	0	
Numerical Calculation	simple	0.034	0.056
	complex	0.094	

The table(Table.5) shows our experimental results.By randomly selecting 500 questions from each category in the dataset we created for LLM to answer directly. In the end, we calculate an average accuracy of 0.002 for Cell_Query, 0.236 for Fact_Judgment, 0.0 for Data_Filtering, and 0.056 for Numerical Calculation. Although the accuracy rate is slightly higher on Fact Judge questions, this does not necessarily mean that LLM has answered the questions based on their own knowledge, because for Fact Judge questions, LLM only needs to provide 'true' or 'false' answers, and the accuracy rate of random guessing on this question may reach 50%.However, now it is 23.6% for Fact_Judgment.

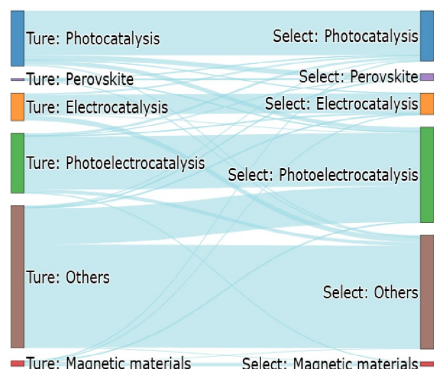
Therefore, this experiment proves that the data we used does not overlap with the data pre-trained by LLM.

CataTQA codebase is hosted and version-tracked via GitHub. It will be permanently available under the link <https://github.com/qcui2025/CataTQA>. The download link of all the datasets can be found in the GitHub repository.

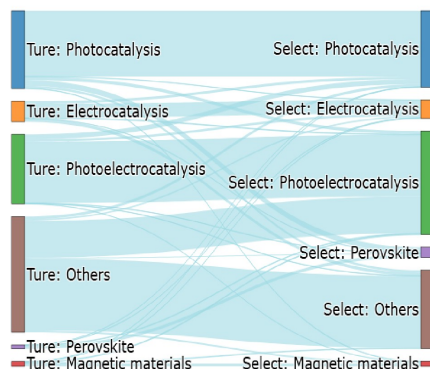
CataTQA stands as a community-driven open-source initiative. We are not only firmly committed but also well-resourced to sustain and vigorously advance it in the times ahead. Envisioning its expansion, we aim to incorporate a more diverse range of tasks, tools, and baseline methodologies. Additionally, we warmly welcome external contributors to partake in this endeavor, as their involvement will be pivotal in shaping the future trajectory of CataTQA.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

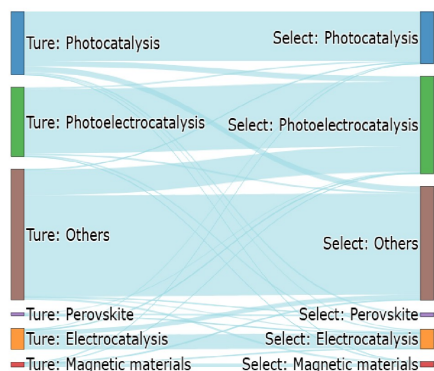
Gpt-4o



Claude-3



Deepseek-v3



Gemini-2.5

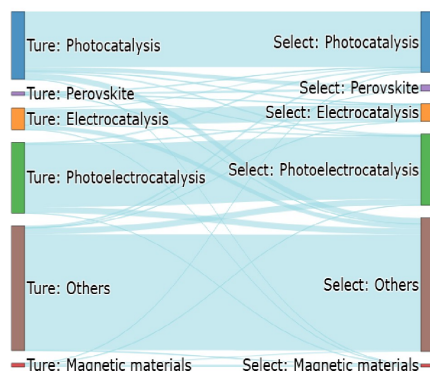


Figure 6: Alluvial plot of the Table Positioning Ability of Large Language Models

We license our work using Apache 2.0. All the datasets will be publicly released through the aforementioned GitHub link.

The authors will bear all responsibility in case of violation of rights.

E.2 FULL DATASET INFORMATION

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166

Table 6: Detailed Information of CataTQA Raw Data

Domain	Dataset or paper name	Access url	Index
photo-catalysis	Machine learning aided design of perovskite oxide materials for photocatalytic water splitting	https://www.sciencedirect.com/science/article/pii/S2095495621000644#s0090	table1 table2 table3
	Data mining in photocatalytic water splitting over perovskites literature for higher hydrogen production	https://www.sciencedirect.com/science/article/pii/S0926337318309470#sec0130	table4 table5
	An insight into tetracycline photocatalytic degradation by MOFs using the artificial intelligence technique	https://www.nature.com/articles/s41598-022-10563-8#Sec10	table6
	Analysis of photocatalytic CO2 reduction over MOFs using machine learning	https://pubs.rsc.org/en/content/articlelanding/2024/ta/d3ta07001h	table7
	Data-driven for accelerated design strategy of photocatalytic degradation activity prediction of doped TiO2 photocatalyst	https://www.sciencedirect.com/science/article/pii/S2214714422005700#s0055	table8
	A generalized predictive model for TiO2-Catalyzed photo-degradation rate constants of water contaminants through artificial neural network	https://www.sciencedirect.com/science/article/pii/S0013935120305909	table9

1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199

Domain	Dataset or paper name	Access url	Index
	Statistical information review of CO2 photocatalytic reduction via bismuth-based photocatalysts using artificial neural network	https://www.sciencedirect.com/science/article/pii/S1110016824008640?via=ihub	table10
	Accelerated Design for Perovskite-Oxide-Based Photocatalysts Using Machine Learning Techniques	https://www.mdpi.com/1996-1944/17/12/3026	table11
electro-catalysis	Building Blocks for High Performance in Electrocatalytic CO2 Reduction: Materials, Optimization Strategies, and Device Engineering	https://acs.figshare.com/articles/dataset/Building_Blocks_for_High_Performance_in_Electrocatalytic_CO_sub_2_sub_Reduction_Materials_Optimization_Strategies_and_Device_Engineering/5293804?file=9064090	table12
	Unlocking New Insights for Electrocatalyst Design: A Unique Data Science Workflow Leveraging Internet-Sourced Big Data	https://github.com/ruiding-uchicago/InCrEDible-MaT-GO	table13
	Perovskite-based electrocatalyst discovery and design using word embeddings from retrained SciBERT language model	https://github.com/aranm917/Perovskite-based-electrocatalyst-design-and-discovery	table14

Domain	Dataset or paper name	Access url	Index
	Exploring the Composition Space of High-Entropy Alloy Nanoparticles for the Electrocatalytic H ₂ /CO Oxidation with Bayesian Optimization	https://github.com/vamints/Scripts_BayesOpt_PtRuPdRhAu_paper	table15
	High Throughput Discovery of Complex Metal Oxide Electrocatalysts for the Oxygen Reduction Reaction	https://data.caltech.edu/records/1km87-52j70	table16
photo-electro-catalysis	High-throughput OCM data	https://cads.eng.hokudai.ac.jp/datamanagement/datasources/21010bbe-0a5c-4d12-a5fa-84eea540e4be/	table17
	CatApp Data	https://cads.eng.hokudai.ac.jp/datamanagement/datasources/20de069b-53cf-4310-9090-1738f53231e2/	table18
	Oxidative Coupling of Methane	https://cads.eng.hokudai.ac.jp/datamanagement/datasources/9436f770-a7e2-4e87-989b-c5a9ce2312bf/	table19
	ChemCatChem	https://cads.eng.hokudai.ac.jp/datamanagement/datasources/224dd7ad-7677-4161-b744-a0c796bf5347/	table20
	HTP OCM data obtained with catalysts designed on the basis of heuristics derived from random catalyst data	https://cads.eng.hokudai.ac.jp/datamanagement/datasources/92200ba4-7644-44ca-9801-ed3cc52fc32f/	table21

1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232

1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265

Domain	Dataset or paper name	Access url	Index
	Perovskite Data	https://cads.eng.hokudai.ac.jp/datamanagement/datasources/f1b42c58-a423-4ec2-8bcf-e66c6470ff7d/	table22
	Random catalyst OCM data by HTE	https://cads.eng.hokudai.ac.jp/datamanagement/datasources/f7e30001-e440-4c1a-be64-ea866b2f77cb/	table23
	Synthesis of Heterogeneous Catalysts in Catalyst Informatics to Bridge Experiment and High-Throughput Calculation	https://cads.eng.hokudai.ac.jp/datamanagement/datasources/f2a6d4f2-91be-48ba-bf13-ffeabd90f6ee/	table24
	Multi-component La2O3-based catalysts in OCM	https://cads.eng.hokudai.ac.jp/datamanagement/datasources/d6347fc1-e4d7-412e-aed5-a8ffa415a703/	table25
	Catalyst Modification in OCM via Manganese Promoter	https://cads.eng.hokudai.ac.jp/datamanagement/datasources/32dbec2c-c3d5-43ec-962a-90dba719bb44/	table26
	Leveraging Machine Learning Engineering to Uncover Insights in Heterogeneous Catalyst Design for Oxidative Coupling of Methane	https://cads.eng.hokudai.ac.jp/datamanagement/datasources/d84c1e22-ceb9-488a-8d45-4c7cf1c603b5/	table27
	Oxidative of Coupling Literature and Highthroughput Data	https://cads.eng.hokudai.ac.jp/datamanagement/datasources/adb27910-d0e5-4a22-9415-580bf597035a/	table28

Domain	Dataset or paper name	Access url	Index
	Catalytic Material Database	http://cmd.us.edu.pl/catalog/	table29 table30
	Catalyst Hub	-	table31
magnetic material	Magnetic Database	https://doi.org/10.15131/shef.data.24008055.v1	table32
	Materials database of Curie and Néel magnetic phase transition temperatures	https://doi.org/10.6084/m9.figshare.5702740.v1	table33
perovskite	Data-driven design of molecular nanomagnets	https://go.uv.es/rosaleny/SIMDAVIS	table34
	Predicting the thermodynamic stability of perovskite oxides using machine learning models	-	table35 table36 table37
	Crystallography Open Database(COD)	http://www.crystallography.net/cod/	table38
	Alloy synthesis and processing by semi-supervised text mining	https://www.nature.com/articles/s41524-023-01138-w	table39
	A Machine Learning Approach to Zeolite Synthesis Enabled by Automatic Literature Data Extraction	https://github.com/olivettigroup/table_extractor	table40
others	ZeoSyn: A Comprehensive Zeolite Synthesis Dataset Enabling Machine-Learning Rationalization of Hydrothermal Parameters	https://github.com/eltonpan/zeosyn_dataset	table41

1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298

1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331

Domain	Dataset or paper name	Access url	Index
	Unveiling the Potential of AI for Nanomaterial Morphology Prediction	https://github.com/acid-design-lab/Nanomaterial_Morphology_Prediction	table42
	AFLOW-2 CFID dataset	https://doi.org/10.1016/j.commatsci.2012.02.005	table43
	Alexandria_DB PBE 3D all	https://alexandria.icams.rub.de/	table44
	arXiv dataset	https://www.kaggle.com/Cornell-University/arxiv	table45
	CCCBDB dataset	https://cccbdb.nist.gov/	table46
	3D dataset	https://www.nature.com/articles/s41524-020-00440-1	table47
	2D dataset	https://www.nature.com/articles/s41524-020-00440-1	table48
	halide perovskite dataset	https://doi.org/10.1039/D1EE02971A	table49
	hMOF dataset	https://doi.org/10.1021/acs.jpcc.6b08729	table50
	HOPV15 dataset	https://www.nature.com/articles/sdata201686	table51
	Surface property dataset	https://doi.org/10.1039/D4DD00031E	table52
	JARVIS-FF	https://www.nature.com/articles/s41524-020-00440-1	table53
	MEGNET-3D CFID dataset	-	table54
	Materials Project-3D CFID dataset	https://next-gen.materialsproject.org/	table55
	Materials Project-3D CFID dataset 84k	-	table56
	OQMD-3D dataset	https://www.oqmd.org/download/	table57

1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364

Domain	Dataset or paper name	Access url	Index
	Polymer genome	https://datadryad.org/dataset/doi:10.5061/dryad.5ht3n	table58
	QETB dataset	https://arxiv.org/abs/2112.11585	table59
	QM9 dataset	https://www.nature.com/articles/sdata201422	table60
	QM9 standardized dataset 130k	-	table61
	QMOF dataset	https://www.cell.com/matter/fulltext/S2590-2385(21)00070-9	table62
	SNUMAT Hybrid functional dataset	https://www.nature.com/articles/s41597-020-00723-8	table63
	SSUB dataset	https://github.com/wolverton-research-group/qmpy	table64
	chem dataset	https://www.nature.com/articles/s41524-018-0085-8	table65
	InterMat dataset	https://doi.org/10.1039/D4DD00031E	table66
	2DMatPedia dataset	http://www.2dmatpedia.org/	table67
	vacancy dataset	https://doi.org/10.1063/5.0135382	table68