

TOOLBRIDGE: AN OPEN-SOURCE DATASET TO EQUIP LLMs WITH EXTERNAL TOOL CAPABILITIES

Anonymous authors

Paper under double-blind review

ABSTRACT

Through the integration of external tools, large language models (LLMs) such as GPT-4o and Llama 3.1 significantly expand their functional capabilities, evolving from elementary conversational agents to general-purpose assistants. We contend that the primary drivers of these advancements are the quality and diversity of the training data. However, the existing LLMs with external tool integration provide only limited transparency regarding their datasets and data collection approaches, which has led to the initiation of this study. Specifically, in this work, we endeavor to present a detailed exposition of the methodology for constructing datasets that facilitate LLMs in effectively learning how to utilize external tools and make this process available to the public through the introduction of ToolBridge. ToolBridge proposes to leverage a collection of general open-access datasets as its raw dataset pool and incorporates a series of strategies to identify the appropriate data entries for external tool API insertions. By supervised fine-tuning (SFT) on these curated data entries, LLMs can invoke external tools in appropriate contexts to boost their predictive accuracy, particularly for essential functions including factual retrieval, data processing and numerical computation. Our experiments meticulously isolate model architectures and training configurations, zeroing in exclusively on the role of data. The experimental results indicate that LLMs trained on ToolBridge exhibit consistent performance gains on both standard benchmarks and custom evaluation datasets. All associated code and data will be released as open source, promoting transparency and facilitating the broader community to explore methodologies for equipping LLMs with external tools capabilities.

1 INTRODUCTION

Large language models (LLMs) have revolutionized natural language processing, excelling in tasks including question answering, summarization, and text generation Jiang et al. (2023); Achiam et al. (2023); Dubey et al. (2024); Gunter et al. (2024); Team et al. (2024). Despite the impressive achievements of LLMs, they persistently underperform in fundamental areas, such as arithmetic and factual lookup, where external tools can effectively provide solutions Schick et al. (2024).

Consequently, there has been a rise in research efforts committed to equipping LLMs with the ability to utilize external tools. These efforts can be broadly categorized into two distinct areas: (1) function calling, such as Toolllm Qin et al. (2023) and API-Bank Li et al. (2023), which offer curated datasets and standardized evaluation protocols to systematically assess and compare the tool-use capabilities of LLMs, focusing on enabling LLMs to function as agents that call APIs without requiring explicit training on their usage. (2) tool former, exemplified by GPT-4o and Toolformer Schick et al. (2024), which seek to empower LLMs with the capability to autonomously invoke and utilize external tools during inference, emphasizing advanced reasoning by enabling LLMs to discern which external tool to call, the appropriate timing for the call, and the systematic parsing and integration of the execution results into the generated text. This study primarily concentrates on the latter area.

Within the latter context, Llama 3.1 Dubey et al. (2024) is trained to utilize Brave Search, Wolfram Alpha API and Python interpreter by incorporating synthetic data combined with human-annotated data. Toolformer Schick et al. (2024) adopts a self-supervised training strategy to equip LLMs with the capability to leverage diverse tools effectively. However, these existing works tend to either only open-source their trained models and code or restrict access to the models through usage interfaces

054	User: "Calculate the answer to: 5^2 ."
055	Assistant: "The answer to 5^2 is <code><python>answer = 5**2\nprint(answer)</python><result>25</result></code> 25."
056	
057	User: "Standardize the following date to the ISO 8601 format.\nJuly 26th, 2021."
058	Assistant: "The standardized date in ISO 8601 format is: <code><python>from datetime import datetime\ndate_str = 'July 26th, 2021'\ndate_obj = datetime.strptime(date_str, '%B %dth, %Y')\niso_date = date_obj.date().isoformat()\nprint(iso_date)</python><result>2021-07-26</result></code> 2021-07-26."
059	
060	
061	User: "Question: What's the answer to that question: what currency do the ukraine use?\nAnswer:"
062	Assistant: "The currency used in Ukraine is the <code><python>import requests\n\nresponse = requests.get('https://restcountries.com/v3.1/name/Ukraine')\ncurrency = response.json()[0]['currencies']\ncurrency_code = list(currency.keys())[0]\ncurrency_name = currency[currency_code]['name']\nprint(currency_name)</python><result>Ukrainian hryvnia</result></code> Ukrainian hryvnia (UAH)."
063	
064	
065	
066	
067	User: "Solve $5z - 180 = -165$ for z . What is the answer?"
068	Assistant: "To solve for z in the equation $5z - 180 = -165$, we need to isolate z on one side of the equation. \n\nFirst, we can add 180 to both sides of the equation to get rid of the constant term on the left side: <code><python>5*z - 180 + 180 = -165 + 180</code> \n\nSimplifying the left side, we get: <code><python>5*z = 15</code> \n\nNext, we can divide both sides of the equation by 5 to isolate z : <code><python>5*z/5 = 15/5</code> \n\nSimplifying, we get: <code><python>from sympy import symbols, Eq, solve\nz = symbols('z')\nequation = Eq(5*z - 180, -165)\nsolution = solve(equation, z)\nprint(solution[0])</python><result>3</result></code> $z = 3$ \n\nTherefore, the solution to the equation $5z - 180 = -165$ is $z = 3$."
069	
070	
071	
072	
073	

Figure 1: **Examples of data entries from ToolBridge.** The proposed pipeline is capable of generating a large volume of data entries to facilitate model learning to utilize external tools.

Achiam et al. (2023); Dubey et al. (2024); Gunter et al. (2024). The data underlying the training of these models is commonly withheld from public disclosure, enveloped in secrecy.

The widespread availability of pre-trained models stands in stark contrast to the limited transparency surrounding their training datasets, which impedes the community’s progress in advancing discoveries in this area. This emphasizes the critical need to explore ways to create high-quality, publicly available datasets for training LLMs to utilize external tools to facilitate their reasoning process.

In response to the challenges discussed, this paper proposes a pipeline for the large-scale creation of datasets tailored to equip LLMs with the capability to effectively utilize external tools. Specifically, we begin by aggregating a substantial collection of open-source datasets used for LLMs supervised fine-tuning (SFT) from the community, which circumvents proprietary concerns including copyright issues. Upon establishing the dataset pool, we propose a systematic strategy to assist in identification of valuable data entries and convert them into a standardized format. Finally, consistency validation is conducted to further boost the quality of the converted data entries and ToolBridge is constructed. As demonstrated in Figure 1, we showcase the data entries from ToolBridge, which function to guide LLMs in understanding how to incorporate external tools in appropriate contexts, thereby improving the accuracy and reliability of their outputs. For instance, the third case in Figure 1 serves to instruct LLMs in leveraging the *requests* module to gather factual information from web sources. Moreover, LLMs can employ the fourth scenario to understand how *sympy* library can be applied to solve linear equation of one variable.

In summary, the contributions of this paper are as follows,

- We propose a pipeline capable of producing large volumes of entries for training LLMs to incorporate various external tools. The collection of over 178K yielded data entries, named ToolBridge, will be open-sourced to the community, marking a significant advancement in the transparency and accessibility of the data for training LLMs to leverage external tools.
- Our experimental results demonstrates that LLMs supervised fine-tuned on ToolBridge can achieve consistent performance improvements on several standard benchmarks.
- We also propose to curate a set of data entries to examine the abilities of LLMs in numerical calculation, data processing and factual retrieval before and after supervised fine-tuning on ToolBridge. Experimental outcomes reaffirm the effectiveness of ToolBridge

This is, to our knowledge, the first work in the domain of enabling LLMs to learn to utilize external tools that open-sources the training data. We anticipate that ToolBridge will facilitate the community

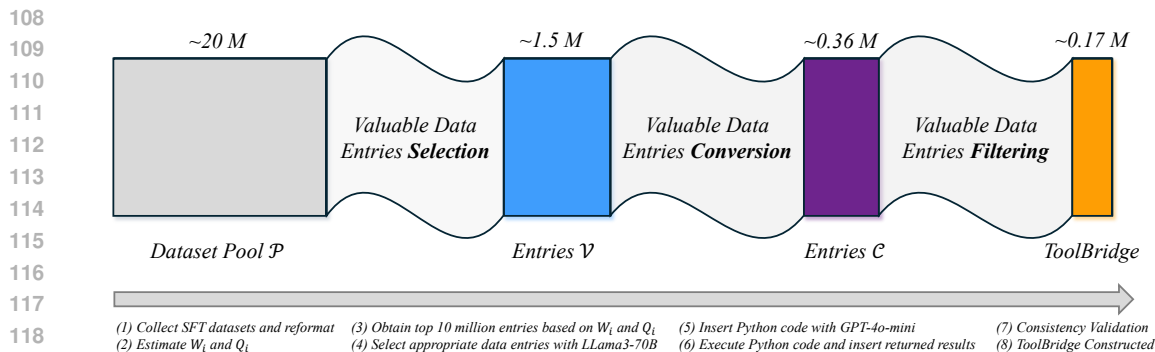


Figure 2: **Overview of the construction pipeline for ToolBridge.** The construction of ToolBridge follows three main steps: identifying valuable data entries in the dataset pool, converting these data entries through the integration of external tool calls and finally conducting a thorough cleanup of the converted data entries by consistency validation.

in further investigating the ability of LLMs to use external tools, thereby advancing LLMs from basic conversational models to versatile general-purpose assistants.

2 RELATED WORK

Tool Use for LLMs. Enabling LLMs to use external tools like search engines and code interpreters significantly broadens the range of tasks LLMs can address and strengthens their predictive accuracy. The methodologies for equipping LLMs with the capability to employ external tools can be broadly categorized into two paradigms, *i.e.*, function calling and tool former. In particular, function calling emphasizes allowing LLMs to act as agents that invoke APIs using predefined functions and prompts without necessitating explicit training on their usage Li et al. (2023); Shen et al. (2024). Conversely, tool former prioritize empowering LLMs with the ability to autonomously identify appropriate tools, determine the optimal invocation timings, and incorporate the outputs into their reasoning processes. For example, Komeili (2021) proposed to enable LLMs to adopt a search engine by learning to yield an internet search query based on the context, and then condition its generated response on the search results. Cobbe et al. (2021) facilitated LLMs’ utilization of a calculator during inference by training the models with calculation annotations injected into the datasets. Thoppilan et al. (2022) proposed to assist LLMs in invoking external tools from a toolset, comprising an information retrieval system, a calculator and a translator, by training it to produce a special string TS . Gao et al. (2023) suggested adopting LLMs to interpret natural language problems and yield programs as intermediate reasoning, while delegating the solution process to a runtime environment like a Python interpreter. Toolformer Schick et al. (2024) allowed LLMs to learn how to adopt the external tools through a self-supervised learning approach. Of late, the works like GPT-4o, Llama 3.1 Dubey et al. (2024) and Apple LLMs further strengthened LLMs’ ability to leverage external tools through improvements in training data, model architectures, *etc.* This study falls into the latter category, namely tool former.

Although previous research in the domain of tool former are highly praiseworthy, they seldom make the data required for training their models publicly available, which is crucial for the community to advance research and build upon their contributions. This paper presents a pipeline aimed at yielding data entries for training models in external tool utilization, along with open-sourcing all data entries produced using this methodology This open access facilitates the development of more effective and efficient algorithms for the next generation of LLMs integrated with external tool functionalities.

Training Datasets for Tool Use. Previous datasets designed to train LLMs to utilize external tools primarily fall within the function calling paradigm. For instance, Qin et al. (2023) collected a high-quality instruction-tuning dataset ToolBench, which is constructed automatically adopting ChatGPT. Li et al. (2023) introduced API-Bank, which encompasses 1,888 tool-use dialogues from 2,138 APIs spanning 1,000 distinct domains. However, to the best of our knowledge, there is a notable absence of research efforts that have open-sourced training datasets within the tool former domain.

To address this significant gap, this paper presents ToolBridge - a dataset of more than 178,000 data entries to support LLMs in effectively learning to utilize external tools within tool former paradigm.

3 TOOLBRIDGE

Previous LLMs like GPT-4o and Llama 3.1 only provide limited information on how they curate the data entries to empower themselves to employ external tools. To address the lack of transparency in training data, we propose a generic pipeline for constructing large-scale datasets from public sources to enable LLMs to use external tools. As indicated in Figure 1, the whole pipeline follows three main steps: valuable data entries **selection**, **conversion** and **filtering**.

Source	# of Entries	Source	# of Entries
School Math 0.25M	248,481	LIMA	1,330
code_instructions_120k_alpaca	121,959	TigerBot	1,199,030
Platypus	24,926	TSI-v0	5,607,620
ShareGPT90K	90,665	LaMini-Instruction	2,585,615
WizardLM_Orca	54,974	Bactrian-X	67,017
WizardLM_evol_instruct_70k	70,000	Baize	210,311
tiny-codes	1,632,309	COIG	178,246
WizardLM_evolve_instruct_V2	143,000	MOSS SFT	1,074,551
No Robots	10,000	AlpacaDataCleaned	51,760
ign_clean_instruct_dataset_500k	508,620	GPT-4all	808,812
GPT-4-LLM	113,003	Alpaca	52,002
ChatAlpaca	20,000	self-instruct	82,439
OpenOrca	4,233,923		

Table 1: The composition of our dataset pool \mathcal{P} to construct ToolBridge.

3.1 DATASET POOL CONSTRUCTION

Our work starts with a review of the data accessible in the community for the purpose of supervised fine-tuning (SFT). Table 1 summarizes the results. Owing to the diverse range of teams contributing the SFT datasets, there is significant heterogeneity in their formats, which introduces difficulties on effective model training. Hence, we first reformat all candidate datasets into a standardized ChatML format for further processing:

```
data entry = [{"role": "user", "content": "..."},
              {"role": "assistant", "content": "..."},
              ...].
```

After reformatting all datasets, we construct the dataset pool as $\mathcal{P} = \{(\mathcal{D}_i, W_i, Q_i) \mid i \geq 0\}$, where \mathcal{D}_i denotes one candidate dataset, W_i measures the proportion of valuable entries for each dataset, and Q_i serves as a metric for assessing the quality of each dataset.

Practically, to obtain W_i , we first perform random sampling on \mathcal{D}_i to produce a subset \mathcal{S}_i , containing 1% data entries of \mathcal{D}_i . Llama3-70B is then applied to judge the appropriateness of each entry in \mathcal{S}_i for external tool invocation to enhance reasoning, where the prompt employed is shown in Appendix A.1. In generally, if an entry is deemed suitable for invoking external tools to help LLMs' reasoning process, we label it as a valuable entry. At last, we determine W_i as the ratio between the number of valuable data entries and the total number of data entries in \mathcal{S}_i .

Additionally, we observe that some candidate datasets within the dataset pool \mathcal{P} are partially sourced from the Internet via web scraping, resulting in the inclusion of certain meaningless HTML tags and other irrelevant content. Thus, we incorporate Q_i as an additional metric to evaluate the data quality of each candidate dataset. To compute Q_i , we randomly sample N data entries from \mathcal{D}_i and conduct a manual review to identify any presence of irrelevant characters or content, where we configure N as 100 by default. Q_i is then derived as the fraction of data entries devoid of irrelevant contents over the total number of sampled entries N .

3.2 VALUABLE DATA ENTRIES SELECTION

Upon constructing the dataset pool \mathcal{P} , due to the large scale of candidate data entries, we propose to select 10 million data entries from \mathcal{P} for further processing tailored to the dataset attributes W_i and

Q_i . Particularly, we first arrange \mathcal{D}_i in descending order guided by the value of $Q_i \times W_i$. Then, the data entries are selected from the top-ranked datasets sequentially until the overall volume amounts to 10 million data entries.

Subsequently, Llama3-70B with the prompt detailed in Appendix A.1 is applied to ascertain whether each entry within the 10 million samples is appropriate for LLMs to enhance reasoning via utilizing external tools. And we represent the collection of these appropriate data entries with \mathcal{V} , namely, the valuable data entries. In Table 2, we present the distribution of \mathcal{V} across the respective datasets from which they are derived. It is observed that leveraging the capabilities of Llama3-70B, we can refine the 10 million data entries down to 1,527,153 valuable entries.

Source	# of Entries	Source	# of Entries
School Math 0.25M	205,996	ChatAlpaca	2,643
Platypus	7,776	ShareGPT90K	24,348
WizardLM_Orca	8,659	WizardLM_evol_instruct_70k	28,293
WizardLM_evol_instruct_V2	5,399	MOSS SFT	136,603
TigerBot	182,249	GPT-4all	47,627
COIG	15,181	LIMA	309
AlpacaDataCleaned	13,805	GPT-4-LLM	9,978
Bactrian-X	3,313	OpenOrca	834,974

Table 2: The composition of the selected valuable data entries \mathcal{V} .

3.3 VALUABLE DATA ENTRIES CONVERSION

Following valuable data entries selection, we further convert the selected entries, allowing LLMs to learn how to invoke the external tools effectively within the proper context to support their reasoning process. In particular, we draw on previous methodologies Schick et al. (2024); Dubey et al. (2024); Thoppilan et al. (2022) by embedding special characters in each selected entry to enable the external tool invocation, and LLMs are used to pinpoint the appropriate context for calling external tools.

Algorithm 1 Process Special Tokens During Inference

Require: Python interpreter and *inputs*.
Ensure: Execute the code enclosed by `<python>` and `</python>`, and use the captured output as a condition for the subsequent text generation.

- 1: Initialize an empty list *outputs*
- 2: Initialize *condition* \leftarrow None
- 3: **while** *outputs* is empty **or** *outputs*[-1] \neq `<|end.of.text|>` **do**
- 4: *output* \leftarrow LLM(*inputs*)
- 5: **if** *output* = `<python>` **then**
- 6: *start* \leftarrow length of *outputs* + length of `<python>`
- 7: **else if** *output* = `</python>` **then**
- 8: Extract substring *code* \leftarrow *outputs*[*start* :]
- 9: *condition* \leftarrow ExecutePython(*code*)
- 10: **end if**
- 11: Append *output* to both *inputs* and *outputs*
- 12: **if** *condition* is not None **then**
- 13: Append *condition* to both *inputs* and *outputs*
- 14: Reset *condition* \leftarrow None
- 15: **end if**
- 16: **end while**
- 17: Post-process and return *outputs*

As illustrated in the examples in Figure 1, `<python>` and `</python>` are represented as a pair of special tokens. The content enclosed by the special tokens specifies the Python invocation for using external tools. During the construction of ToolBridge, we predominantly use GPT-4o-mini to insert the special tokens in the appropriate context within each data entry identified in Section 3.2, as well as to create the associated code for invoking external tools. To facilitate the return of tool execution

270 results, we examine the code generated by GPT-4o when calling Python API as part of its reasoning
 271 process. It is observed that the final results are always printed at the end of the code. In line with the
 272 strategy of GPT-4o, we include a directive in the prompt for GPT-4o-mini to print the final result of
 273 the tool execution as the concluding line of the code (refer to Appendix A.2).

274 Following GPT-4o-mini’s processing of \mathcal{V} , we retrieve the code segments between `<python>` and
 275 `</python>`, execute them, and finally insert the captured output after corresponding `</python>`,
 276 where the output will be wrapped within another pair of special tokens, denoted as `<result>` and
 277 `</result>`. To summarize, the tool invocation in ToolBridge can be formatted as,
 278

279 `<python>generated code</python><result>captured output</result>`.
 280

281 During the reasoning process of the models, it is sufficient to verify the existence of the pre-defined
 282 special tokens `<python>` and `</python>`. If identified, the Python interpreter will be adopted to
 283 execute the enclosed code between the special tokens, and the final result is wrapped in `<result>`
 284 and `</result>` to condition further text generation. Algorithm 1 describes the primary inference
 285 process of LLMs post-SFT on the ToolBridge dataset.

286 Table 3 summarizes the composition of the data entries after converting \mathcal{V} , denoted as \mathcal{C} .
 287

Source	# of Entries	Source	# of Entries
School Math 0.25M	150,104	ChatAlpaca	116
Platypus	82	ShareGPT90K	3
WizardLM_Orca	88	WizardLM_evol_instruct_70k	3,716
WizardLM_evol_instruct_V2	8	OpenOrca	101,715
TigerBot	66,793	GPT-4all	8,936
COIG	7,877	LIMA	107
AlpacaDataCleaned	4,510	GPT-4-LLM	2,577
Bactrian-X	84		

288
289
290
291
292
293
294
295
296
297 Table 3: The composition of the converted data entries \mathcal{C} .

298
299 It is worth noting that the data entries in \mathcal{C} , totaling 364,605, are significantly fewer than those in \mathcal{V}
 300 which amount to 1,527,153. We attribute this to the following factors,

- 301 • The returned data entries that lack the Python code inserted by GPT-4o-mini are eliminated.
 302 Approximately 19.2% of the data entries in \mathcal{V} fall under this category, which indicates that
 303 GPT-4o-mini regards these entries as not requiring external tools to aid in LLMs’ reasoning.
 304
- 305 • The returned data format did not match the expected structure, resulting in parsing failure.
 306 Instances include an unequal count of `<python>` and `</python>` tokens, modifications
 307 to the original content alongside Python code insertion, to name a few. Such entries account
 308 for approximately 27.2% of the data entries in \mathcal{V} .
- 309 • The request to GPT-4o-mini fails. These data entries constitute approximately 2.1% of \mathcal{V} .
- 310 • With Python’s `ast` library, we filter out the returned data entries where the inserted code is
 311 only an assignment followed by a print statement. Such entries represent about 4.8% of \mathcal{V} .
 312 We provide the detailed algorithm we use to filter such data entries in Appendix A.3.
- 313 • We filter around 22.8% data entries of \mathcal{V} to remove tool calls that could not be executed or
 314 the execution time of the tool exceeded 30 seconds (refer to Appendix A.4 for details).
 315

316 In summary, 23.9% of the entries from \mathcal{V} remain in the converted data entries \mathcal{C} .
 317

318 3.4 DATA ENTRIES FILTERING BY CONSISTENCY VALIDATION

319
320 In practice, we observe that LLMs trained on \mathcal{C} do not always base their subsequent contents on the
 321 results produced by the yielded Python code during inference. So, we conduct a reassessment of the
 322 data entries within \mathcal{C} and observe that the execution results from the code generated by GPT-4o-mini
 323 also does not always align with the ensuing text, which can explain LLMs’ sporadic inconsistencies
 between tool execution results and further contents during inference.

To alleviate the issues above, we propose to filter out the entries in \mathcal{C} where the tool execution results are inconsistent with the following text, which is accomplished by validating if the execution results are included in the subsequent content in our approach. Upon the conclusion of the filtering process, the open-source dataset ToolBridge is constructed. In Appendix A.5, we compare the generated text of Llama3-8B after SFT on \mathcal{C} and ToolBridge, which demonstrates the necessity for the data entries filtering by consistency validation.

The data sources that comprise ToolBridge, totaling 178,023 entries, are outlined in Table 4, which represents 48.8% of the total data entries in \mathcal{C} .

Source	# of Entries	Source	# of Entries
School Math 0.25M	100,836	ChatAlpaca	17
Platypus	35	ShareGPT90K	3
WizardLM_Orca	29	WizardLM_evol_instruct_70k	794
WizardLM_evol_instruct_V2	1	OpenOrca	46,449
TigerBot	22,306	GPT-4all	2,616
COIG	2,706	LIMA	27
AlpacaDataCleaned	1,129	GPT-4-LLM	1,043
Bactrian-X	32		

Table 4: The final composition of our proposed ToolBridge.

4 EXPERIMENTS

In this section, we first present the statistics of ToolBridge and then investigate whether LLMs with SFT on ToolBridge could leverage external tools to facilitate their reasoning process. The evaluation consists of two components: (1) by comparing LLMs’ performance on standard benchmarks before and after SFT on ToolBridge (Section 4.3); (2) by evaluating the accuracy of the models on custom datasets RandomQA and FACT (Section 4.4).

Source	# of Tool Use	# of Libraries	Source	# of Tool Use	# of Libraries
School Math 0.25M	104,983	8	ChatAlpaca	73	1
Platypus	36	3	ShareGPT90K	3	0
WizardLM_Orca	33	4	WizardLM_evol_instruct_70k	836	22
WizardLM_evol_instruct_V2	1	0	OpenOrca	46,832	28
TigerBot	22,507	40	GPT-4all	2,870	42
COIG	2,719	28	LIMA	27	2
AlpacaDataCleaned	1,139	23	GPT-4-LLM	1,052	14
Bactrian-X	36	3			

Table 5: Statistics of the usage of external tools in ToolBridge.

4.1 DATASET STATISTICS OF TOOLBRIDGE

In Table 4, we present 15 source datasets involved in the data entries of ToolBridge, along with their respective composition ratios. To prevent any confusion for the datasets (*e.g.*, other datasets with the same name and the same datasets with different versions), we also provide download links for these datasets in Appendix A.6.

Besides, we provide a summary of the frequency of external tool calls in ToolBridge and the variety of Python packages used for these calls, as presented in Table 5. By comparing Table 4 and Table 5, we can observe that the majority of the data entries in ToolBridge involve only a single external tool call. Furthermore, Table 5 also suggests that the data entries in ToolBridge originating from datasets with narrower topics (*e.g.*, School Math 0.25M) generally leverage fewer kinds of Python packages compared to those from datasets with broader topics (*e.g.*, TigerBot), which is consistent with logic. To summarize, there are 183,147 external tool calls in our ToolBridge dataset, utilizing a total of 60 Python packages, including *requests*, *math*, *datetime*, *sklearn*, to name a few.

In Appendix A.12, we also quantify the usage frequency of the 60 Python packages incorporated in ToolBridge, and their distribution is illustrated in Table 13.

4.2 EXPERIMENTAL SETUP

Baseline Models. Our experiments incorporate four baseline models: the base model of Mistral-7B, Llama2-7B, Llama3-8B and Llama3-70B. Also, we remove all the external tool invocation sections in each entry in ToolBridge (denote as ToolBridge[§]) and report the accuracy of four baseline models SFT on ToolBridge[§] as four additional baseline models.

Benchmark Datasets. The standard benchmark datasets leveraged in our experiments include GSM 8K Cobbe et al. (2021), GSM Plus Li et al. (2024), MathBench Liu et al. (2024), Stanford WebQA Berant et al. (2013) and TruthfulQA Lin et al. (2021). For GSM 8k and GSM Plus, the performance is evaluated on their respective test sets under few-shot setting, where we leverage a fixed CoT-n-shot prompt template, as outlined in Li et al. (2024). For MathBench, we report results on MathBench-A, where we transform the multiple-choice questions in the College, High and Middle categories into a question-and-answer format for CoT-n-shot evaluation. To differentiate from standard MathBench, we refer to this adjusted dataset as MathBench*.

We also develop two custom datasets, termed RandomQA and FACT, to evaluate the capabilities of LLMs in data processing, numerical computation and factual retrieval. Section 4.4 elaborates on the specifics of both datasets.

SFT Settings. All the models in our experiments are trained with the open-source TRL library from Hugging Face von Werra et al. (2020). The LoRA module Hu et al. (2021) is employed to carry out SFT on the base model of Mistral-7B, Llama2-7B, Llama3-8B and Llama3-70B using ToolBridge[§] or ToolBridge, with a LoRA rank of 16. Model training is conducted on $64 \times$ MI200 64GB GPUs, with each processing a batch size of 2 (*i.e.*, total batch size is 128). AdamW is employed to optimize the parameters of LoRA, with a cosine learning rate scheduler, configuring the initial lr at $3e-5$ and the total training epoch at 3.

Inference Settings. The primary process of model inference is described in Algorithm 1. Moreover, to handle potential tool call failures during inference, we propose to eliminate failed tool calls from the current output before conditioning the generation of further text. All trained models are evaluated on $16 \times$ MI200 64 GB GPUs, with the max new tokens set to 512 and the temperature set to zero.

Models	SFT data	GSM 8k	GSM Plus	MathBench*
Llama2-7B	-	13.6	8.9	18.0
Llama2-7B-Lora	ToolBridge [§]	16.9	9.9	19.1
Llama2-7B-Lora	ToolBridge	18.1	11.0	21.4
Llama3-8B	-	52.3	36.9	33.0
Llama3-8B-Lora	ToolBridge [§]	53.4	37.8	35.2
Llama3-8B-Lora	ToolBridge	55.8	40.0	37.4
Mistral-7B	-	38.1	25.1	27.8
Mistral-7B-Lora	ToolBridge [§]	42.8	27.6	28.9
Mistral-7B-Lora	ToolBridge	45.0	29.8	31.0
Llama3-70B	-	75.3	54.4	42.1
Llama3-70B-Lora	ToolBridge [§]	78.5	57.6	44.1
Llama3-70B-Lora	ToolBridge	80.1	59.8	46.9

Table 6: Ablation studies on GSM 8k, GSM Plus and MathBench* with (8 shots, CoT) setting.

4.3 RESULTS ON STANDARD BENCHMARKS

In this section, we conduct ablation studies on standard benchmark datasets, including TruthfulQA, GSM 8k, GSM Plus, MathBench and Stanford WebQA, where GSM 8k, GSM Plus and MathBench are primarily responsible for evaluating the capability of LLMs in numerical reasoning and computation, and TruthfulQA and Stanford WebQA are primarily adopted to assess the ability of LLMs in factual retrieval. Table 6, 7 and 8 demonstrates the evaluation results.

It is observed that the models SFT on ToolBridge significantly outperform the baseline models. For example, Llama3-8B SFT on ToolBridge brings 3.1% and 2.2% accuracy improvements to the base model of Llama3-8B and the Llama3-8B model SFT on ToolBridge[§] when evaluating on GSM Plus, respectively. These results indicate that SFT on ToolBridge can help strengthen LLMs’ capabilities in handling numerical computations. Furthermore, ToolBridge is able to facilitate abilities of LLMs

432
433
434
435
436
437
438

Models	SFT data	Stanford WebQA
Llama3-8B	-	21.2
Llama3-8B-Lora	ToolBridge [§]	37.7
Llama3-8B-Lora	ToolBridge	39.9
Mistral-7B	-	34.4
Mistral-7B-Lora	ToolBridge [§]	35.8
Mistral-7B-Lora	ToolBridge	39.1

439
440

Table 7: Ablation studies on Stanford WebQA under zero-shot setting.

441
442
443
444
445
446

Models	SFT data	ROUGE1	BLEURT
Llama3-8B	-	41.2	34.6
Llama3-8B-Lora	ToolBridge [§]	47.0	42.8
Llama3-8B-Lora	ToolBridge	48.7	44.4
Mistral-7B	-	43.5	39.4
Mistral-7B-Lora	ToolBridge [§]	44.9	42.3
Mistral-7B-Lora	ToolBridge	47.7	44.9

447

Table 8: Ablation studies on TruthfulQA under zero-shot setting.

448
449
450
451

in factual retrieval. As shown in Table 7 and 8, SFT on ToolBridge enables Llama3-8B and Mistral-7B to achieve notable gains on Stanford WebQA and TruthfulQA. Specifically, ToolBridge increases the accuracy of Llama3-8B on Stanford WebQA from 21.2% to 39.9%, and on TruthfulQA, it boosts ROUGE1 from 41.2% to 48.7% and BLEURT from 34.6% to 44.4%.

452
453
454
455
456
457

The results above demonstrate that LLMs can effectively learn how to use external tools to enhance their capabilities in basic functions after SFT on ToolBridge. Moreover, it is worthy noting that there is considerable room for improvements in these results, as our emphasis is on the training data, with minimal adjustments made to the model architectures and training strategies, which may help LLMs better learn how to employ external tools through ToolBridge.

458
459
460
461
462
463
464
465
466
467

Models	SFT data	RandomQA-DP-B1	RandomQA-DP-B2	RandomQA-NC-B1	RandomQA-NC-B2
Llama2-7B	-	10.0	9.0	3.3	3.2
Llama2-7B-Lora	ToolBridge [§]	19.2	16.6	7.7	8.6
Llama2-7B-Lora	ToolBridge	53.2	54.0	63.4	60.7
Llama3-8B	-	9.6	9.2	5.8	7.0
Llama3-8B-Lora	ToolBridge [§]	30.3	29.0	15.8	13.9
Llama3-8B-Lora	ToolBridge	62.1	60.0	82.1	80.1
Mistral-7B	-	10.8	9.0	13.8	13.6
Mistral-7B-Lora	ToolBridge [§]	24.7	23.2	16.8	16.5
Mistral-7B-Lora	ToolBridge	61.8	60.5	83.3	82.5
Llama3-70B	-	20.0	17.1	9.6	8.9
Llama3-70B-Lora	ToolBridge [§]	32.1	31.7	22.0	20.3
Llama3-70B-Lora	ToolBridge	74.2	69.9	89.7	89.1

468
469

Table 9: Experimental results on RandomQA under zero-shot setting, where DP denotes data processing and NC means numerical computation.

470
471

4.4 RESULTS ON CUSTOM BENCHMARKS

472
473
474

To further assess whether SFT on the ToolBridge dataset can equip LLMs with the ability to leverage external tools for aiding its reasoning process, we propose to design two custom datasets to evaluate LLMs’ performance before and after SFT on the ToolBridge dataset.

475
476
477
478
479

RandomQA. To assess LLMs’ accuracy in data processing and numerical computation capabilities after SFT on ToolBridge, we propose to design 30 templates capable of generating question-answer pairs to validate the abilities of LLMs in data processing and numerical computations, respectively. Here is one example,

480
481
482
483
484
485

```

481 1 # Templatel: Reverse the order of elements in a list
482 2 array = [random.randint(1, 10000) for _ in range(random.randint(5, 15))]
483 3 question = f"Reverse the order of the elements in the list {array} and
484   ↪ then plus 3 for each element."
485 4 answer = array[::-1]
485 5 answer = [a + 3 for a in answer]
```

The complete list can be found in Appendix A.7. With these pre-defined templates, we first generate four RandomQA datasets, each consisting of 1,000 data entries, and focusing on data processing or numerical computation, *i.e.*, RandomQA-DP-B1/2, RandomQA-NC-B1/2, where DP signifies using template related to data processing to yield the dataset and NC means numerical computation. Then, we evaluate the accuracy of LLMs on the four datasets before and after SFT on ToolBridge. Table 9 demonstrates the results. It is observed that after SFT on ToolBridge, the models shows a significant increase in accuracy on RandomQA.

Models	SFT data	FACT-200-Batch1	FACT-200-Batch2	FACT-200-Batch3
Llama2-7B	-	69.5	55.0	49.0
Llama2-7B-Lora	ToolBridge [§]	86.0	67.5	65.7
Llama2-7B-Lora	ToolBridge	88.5	72.5	73.2
Llama3-8B	-	79.0	60.5	62.2
Llama3-8B-Lora	ToolBridge [§]	89.0	73.0	73.7
Llama3-8B-Lora	ToolBridge	90.0	73.5	80.2
Mistral-7B	-	85.0	67.5	65.9
Mistral-7B-Lora	ToolBridge [§]	86.5	70.0	66.2
Mistral-7B-Lora	ToolBridge	90.5	72.0	77.3
Llama3-70B	-	76.0	53.5	54.0
Llama3-70B-Lora	ToolBridge [§]	88.3	72.4	70.7
Llama3-70B-Lora	ToolBridge	91.2	74.6	82.6

Table 10: Experimental results on FACT under zero-shot setting.

FACT. To determine if the factual retrieval skills of LLMs can be improved by SFT on ToolBridge, we construct the FACT datasets. Specifically, we begin by prompting GPT-4o to produce thousands of question-answer pairs focused on factual retrieval. One example prompt is as following,

1. Generate 100 Q&A pairs for LLM factual retrieval testing. The
 → question topic should be related with Geography. Return them
 → as a Python dictionary, with concise answers (3–5 words).

Appendix A.8 contains the entire set of the adopted prompts for constructing FACT. Upon obtaining the candidate question-answer pairs, we continuously draw random entries from them and manually check their correctness until 200 correct data entries are verified. Through iterating the above process three times with five different prompts each time, we construct three FACT datasets, each with a size of 200 entries, termed FACT-200-Batch1/2/3. Table 10 presents a comparison of the performance of LLMs on the three datasets before and after SFT on the ToolBridge dataset, which shows that LLMs demonstrates enhanced abilities in retrieving factual content.

Besides, one should be aware that relying on external tools to improve the fact retrieval performance of LLMs can be inconsistent. During our experiments, we uncover the following deficiencies

- Using external tools like the *requests* module to obtain information from external websites can lead to potentially unreliable information and LLMs thus yield incorrect predictions. A solution to address this issue is to construct a curated list of trustworthy domains, restricting LLMs from accessing websites beyond this approved whitelist. Besides, we should further refine the training entries to exclude unreliable access links.
- The content scraped from web pages may be overly lengthy, hindering the ability of LLMs to efficiently extract useful information from the scraped contents.
- When LLMs access some external websites, they may attempt to download files from those sites, which poses potential security risks in real-world applications.

Therefore, we recommend that ToolBridge is only used for academic research, with controlled measures in place when it attempts to access external websites.

5 CONCLUSION

This paper proposes ToolBridge, which is the first open-source dataset for training LLMs in the use of external tools in the area of tool former. The principal steps to build ToolBridge involves valuable data entries selection, conversion and filtering. Ablation studies on standard and custom benchmark datasets demonstrate the effectiveness of ToolBridge.

REFERENCES

- 540
541
542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
543 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
544 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 545 Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. Semantic parsing on Freebase from
546 question-answer pairs. In *Proceedings of the 2013 Conference on Empirical Methods in Natural*
547 *Language Processing*, pp. 1533–1544, Seattle, Washington, USA, October 2013. Association for
548 Computational Linguistics. URL <https://www.aclweb.org/anthology/D13-1160>.
- 549 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
550 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to
551 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 552 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
553 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
554 *arXiv preprint arXiv:2407.21783*, 2024.
- 555 Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and
556 Graham Neubig. Pal: Program-aided language models. In *International Conference on Machine*
557 *Learning*, pp. 10764–10799. PMLR, 2023.
- 558 Tom Gunter, Zirui Wang, Chong Wang, Ruoming Pang, Andy Narayanan, Aonan Zhang, Bowen
559 Zhang, Chen Chen, Chung-Cheng Chiu, David Qiu, et al. Apple intelligence foundation language
560 models. *arXiv preprint arXiv:2407.21075*, 2024.
- 561 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,
562 and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint*
563 *arXiv:2106.09685*, 2021.
- 564 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
565 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.
566 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 567 M Komeili. Internet-augmented dialogue generation. *arXiv preprint arXiv:2107.07566*, 2021.
- 568 Minghao Li, Yingxiu Zhao, Bowen Yu, Feifan Song, Hangyu Li, Haiyang Yu, Zhoujun Li, Fei
569 Huang, and Yongbin Li. Api-bank: A comprehensive benchmark for tool-augmented llms. *arXiv*
570 *preprint arXiv:2304.08244*, 2023.
- 571 Qintong Li, Leyang Cui, Xueliang Zhao, Lingpeng Kong, and Wei Bi. Gsm-plus: A comprehensive
572 benchmark for evaluating the robustness of llms as mathematical problem solvers. *arXiv preprint*
573 *arXiv:2402.19255*, 2024.
- 574 Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human
575 falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.
- 576 Hongwei Liu, Zilong Zheng, Yuxuan Qiao, Haodong Duan, Zhiwei Fei, Fengzhe Zhou, Wen-
577 wei Zhang, Songyang Zhang, Dahua Lin, and Kai Chen. Mathbench: Evaluating the theory
578 and application proficiency of llms with a hierarchical mathematics benchmark. *arXiv preprint*
579 *arXiv:2405.12209*, 2024.
- 580 Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru
581 Tang, Bill Qian, et al. Toolllm: Facilitating large language models to master 16000+ real-world
582 apis. *arXiv preprint arXiv:2307.16789*, 2023.
- 583 Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Eric Hambro,
584 Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can
585 teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36, 2024.
- 586 Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugging-
587 gpt: Solving ai tasks with chatgpt and its friends in hugging face. *Advances in Neural Information*
588 *Processing Systems*, 36, 2024.

594 Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya
595 Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open
596 models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
597
598 Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze
599 Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. Lamda: Language models for dialog
600 applications. *arXiv preprint arXiv:2201.08239*, 2022.
601
602 Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan
603 Lambert, Shengyi Huang, Kashif Rasul, and Quentin Gallouédec. Trl: Transformer reinforcement
604 learning. <https://github.com/huggingface/trl>, 2020.
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

648 A APPENDIX

649

650 A.1 LLAMA3-70B PROMPTS

651

652 The prompt applied while using Llama3-70B to identify valuable data is as follow,

653

654 Your task is to determine whether you can add calls to a Python API to a piece
 655 ↪ of text. The calls should help you get information required to complete the
 656 ↪ text. You only need to respond with "Yes" or "No", "Yes" means you can and
 657 ↪ "No" means you can't. Here are some examples:

657

658 Input:

```
659 {"messages": [{"role": "user", "content": "Sort the numbers in the list arr =
660 ↪ [1, 10, 2, 5, -2, 11, 12] in descending order."}, {"role": "assistant",
661 ↪ "content": "The sorted list in descending order is: [12, 11, 10, 5, 2, 1,
662 ↪ -2]."}]}
```

663

664 Output:

665

666 Yes

667

668 Input:

```
669 {"messages": [{"role": "user", "content": "Can you tell me a little bit about
670 ↪ what LaTeX is?"}, {"role": "assistant", "content": "LaTeX is a high-quality
671 ↪ typesetting system; it includes features designed for the production of
672 ↪ technical and scientific documentation. LaTeX is the de facto standard for
673 ↪ the communication and publication of scientific documents. It is widely used
674 ↪ by mathematicians, scientists, engineers, philosophers, linguists,
675 ↪ economists, and other scholars in academia and the professional world."}]}
```

674

675 Output:

676

677 No

678

679 Input:

```
680 {"messages": [{"role": "user", "content": "What is the value of sin 40 degrees
681 ↪ plus cos 31 degrees?"}, {"role": "assistant", "content": "The value is
682 ↪ approximately sin 40 + cos 31 = 0.6428 + 0.8572 = 1.500."}]}
```

682

683 Output:

684

685 Yes

686

687 Input:

```
688 {"messages": [{"role": "user", "content": "Write a Python script that reads an
689 ↪ image and recognizes the text on it."}, {"role": "assistant", "content": "To
690 ↪ read an image and recognize text on it in Python, you can use the
691 ↪ pytesseract library along with Pillow for image processing. Here's a simple
692 ↪ example: \nimport pytesseract; from PIL import Image;
693 ↪ print(pytesseract.image_to_string(Image.open('path_to_image.jpg'))). \nMake
694 ↪ sure you have Tesseract OCR installed on your machine and the required
695 ↪ libraries (pytesseract and Pillow) installed in your Python environment. You
696 ↪ can install them using: \npip install pytesseract pillow"}]}
```

694

695 Output:

696

697 No

698

699 Input:

700

701 PLACEHOLDER

702

703 Output:

702 A.2 GPT-4O-MINI PROMPTS

703

704 The prompt applied while using GPT-4o-mini to convert data entries is as follow,

705

706 Your task is to add calls to a Python API to a piece of text. The calls should help you
 707 ↪ get information required to complete the text. You can call the API by writing
 708 ↪ "<python>code</python>" where "code" is the code to be executed. The last line of
 709 ↪ all code should print the variable that stores the final result. Here are some
 710 ↪ examples of API calls:

711 Input:

```
712 {"messages": [{"role": "user", "content": "Which number is greater, 13.11 or 13.8?"},
  713 ↪ {"role": "assistant", "content": "13.8 is greater than 13.11."}]}
```

714 Output:

```
715 {"messages": [{"role": "user", "content": "Which number is greater, 13.11 or 13.8?"},
  716 ↪ {"role": "assistant", "content": "<python>greater_number = max(13.11,
  717 ↪ 13.8)\nprint(greater_number)</python> 13.8 is greater than 13.11."}]}
```

718 Input:

```
719 {"messages": [{"role": "user", "content": "How many unique words are there in the
  720 ↪ sentence 'The quick brown fox jumps over the lazy dog?'}, {"role": "assistant",
  721 ↪ "content": "There are eight unique words in the sentence 'The quick brown fox jumps
  722 ↪ over the lazy dog.'}]}
```

723 Output:

```
724 {"messages": [{"role": "user", "content": "How many unique words are there in the
  725 ↪ sentence 'The quick brown fox jumps over the lazy dog?'}, {"role": "assistant",
  726 ↪ "content": "There are <python>unique_words = len(set('The quick brown fox jumps over
  727 ↪ the lazy dog'.lower().split()))\nprint(unique_words)</python> eight unique words in
  728 ↪ the sentence 'The quick brown fox jumps over the lazy dog.'}]}
```

729 Input:

```
730 {"messages": [{"role": "user", "content": "What is the area of a circle with a radius of
  731 ↪ 5?"}, {"role": "assistant", "content": "The area of a circle with radius 5 is
  732 ↪ 78.54."}]}
```

733 Output:

```
734 {"messages": [{"role": "user", "content": "What is the area of a circle with a radius of
  735 ↪ 5?"}, {"role": "assistant", "content": "The area of a circle with radius 5 is
  736 ↪ <python>import math\narea = math.pi * 5**2\nprint(area)</python> 78.54."}]}
```

737 Input:

```
738 {"messages": [{"role": "user", "content": "Sort the numbers [5, 3, 8, 1, 2] in ascending
  739 ↪ order."}, {"role": "assistant", "content": "The sorted list is [1, 2, 3, 5, 8]."}]}
```

740 Output:

```
741 {"messages": [{"role": "user", "content": "Sort the numbers [5, 3, 8, 1, 2] in ascending
  742 ↪ order."}, {"role": "assistant", "content": "The sorted list is <python>lst =
  743 ↪ sorted([5, 3, 8, 1, 2])\nprint(lst)</python> [1, 2, 3, 5, 8]."}]}
```

744 Input:

```
745 {"messages": [{"role": "user", "content": "Extract the domain from the email
  746 ↪ 'example@test.com'. cant be empty"}, {"role": "assistant", "content": "The domain of the email
  747 ↪ 'example@test.com' is 'test.com'. cant be empty"}]}
```

748 Output:

```
749 {"messages": [{"role": "user", "content": "Extract the domain from the email
  750 ↪ 'example@test.com'. cant be empty"}, {"role": "assistant", "content": "The domain of the email
  751 ↪ 'example@test.com' is <python>domain =
  752 ↪ 'example@test.com'.split('@')[1]\nprint(domain)</python> 'test.com'. cant be empty"}]}
```

753 Input:

754 PLACEHOLDER

755 Output:

A.3 DETECTING SIMPLE CODE STRUCTURES: ASSIGNMENT FOLLOWED BY PRINT STATEMENT

We use the following code to determine if the inserted code in each entry is simply an assignment followed by printing the variable.

Algorithm 2 Identifying Code Patterns: Distinguishing Assignments Followed by Print Statements

```

756 1 import ast
757 2
758 3 ...
759 4 Function:
760 5     Identifying Code Patterns: Distinguishing Assignments Followed
761 6     by Print Statements
762 7 Argument:
763 8     node: node = ast.parse(code)
764 9 ...
765 10 def isuselesscode(node):
766 11     # Check if the node is an AST Module
767 12     if isinstance(node, ast.Module):
768 13         # Ensure the module has exactly two statements:
769 14         # an assignment and an expression.
770 15         if len(node.body) == 2 and isinstance(node.body[0], ast.
771 16             Assign) and isinstance(node.body[1], ast.Expr):
772 17             assign_node = node.body[0]
773 18             expr_node = node.body[1]
774 19             # Check if the assignment targets a variable and
775 20             # the value is a constant.
776 21             if isinstance(assign_node.targets[0], ast.Name) and
777 22                 isinstance(assign_node.value, ast.Constant):
778 23                 # Check if the expression is a function call
779 24                 # to 'print'.
780 25                 if isinstance(expr_node.value, ast.Call) and
781 26                     isinstance(expr_node.value.func, ast.Name) and
782 27                     expr_node.value.func.id == 'print':
783 28                     # Ensure 'print' has exactly one argument.
784 29                     if len(expr_node.value.args) == 1:
785 30                         arg = expr_node.value.args[0]
786 31                         # Check if the argument to 'print' is
787 32                         # the same variable assigned earlier.
788 33                         if isinstance(arg, ast.Name) and arg.id ==
789 34                             assign_node.targets[0].id:
790 35                             return True
791 36                         # Alternatively, check if 'print' uses an
792 37                         # f-string format with the variable.
793 38                         elif isinstance(arg, ast.JoinedStr):
794 39                             for value in arg.values:
795 40                                 if isinstance(value, ast.
796 41                                     FormattedValue) and isinstance(
797 42                                         value.value, ast.Name):
798 43                                     # Confirm the formatted
799 44                                     # variable is the same as
800 45                                     # the assigned variable.
801 46                                     if value.value.id ==
802 47                                         assign_node.targets[0].id:
803 48                                         return True
804 49
805 50 return False

```

A.4 OPTIMIZING DATASET QUALITY BY REMOVING NON-EXECUTABLE AND SLOW TOOL CALLS

The core code for removing data entries that either fail to execute or exceed the execution time limit is as follows.

Algorithm 3 Efficient Dataset Curation: Filtering Non-Executable and Time-Consuming Tool Calls

```

817 1 import re
818 2 import io
819 3 import contextlib
820 4 import multiprocessing
821 5
822 6 # Execute Python code
823 7 def executecode(code):
824 8     with io.StringIO() as buf, contextlib.redirect_stdout(buf):
825 9         try:
826 10             exec(code)
827 11             return buf.getvalue().strip()
828 12         except Exception:
829 13             return None
830 14
831 15 # Execute Python code with timeout = 30 second
832 16 def safeexecutecode(code, timeout=30):
833 17     result_queue = multiprocessing.Queue()
834 18     def target():
835 19         result = executecode(code)
836 20         result_queue.put(result)
837 21     process = multiprocessing.Process(target=target)
838 22     process.start()
839 23     process.join(timeout)
840 24     if process.is_alive():
841 25         process.terminate()
842 26         process.join()
843 27         return None
844 28     else:
845 29         return result_queue.get() if not result_queue.empty() else
846 30         None
847 31
848 32 # Filtering
849 33 DATASET_FILTERED = []
850 34 for item in tqdm(DATASET):
851 35     success_calling = False
852 36     for message in item['messages']:
853 37         if message['role'] == 'assistant':
854 38             code_matches = re.findall(r'<python>(.*?)</python>',
855 39                                     message['content'], re.DOTALL)
856 40             if code_matches:
857 41                 for code in code_matches:
858 42                     python_str = f'<python>{code}</python>'
859 43                     result = safeexecutecode(code=code)
860 44                     if result is None:
861 45                         message['content'] = message['content'].
862 46                             replace(python_str, '')
863 47                     else:
864 48                         result_str = f'<result>{result}</result>'
865 49                         message['content'] = message['content'].
866 50                             replace(python_str, f'{python_str}{
867 51                                 result_str}')
868 52                 success_calling = True
869 53     if success_calling:
870 54         DATASET_FILTERED.append(item)

```

A.5 COMPARE OUTPUTS FROM LLAMA3 TRAINED ON \mathcal{C} AND TOOLBRIDGE

We notice in actual applications that the models train on \mathcal{C} may not always use the Python output as a basis for generating further content. Here are some examples.

<i>Inconsistent outputs from Llama3 trained on \mathcal{C}</i>
<i>Example 1:</i>
User: "Count the number of vowels in the sentence 'This is a simple sentence'."
Assistant: "There are <code><python>sentence = 'This is a simple sentence'\nvowels = 'aeiou'\ncount = sum(1 for char in sentence if char.lower() in vowels)\nprint(count)</python><result>8</result></code> 5 vowels in the sentence 'This is a simple sentence'."
<i>Example 2:</i>
User: "What is the date 30 days after January 31, 2022?"
Assistant: "The date 30 days after January 31, 2022 is <code><python>from datetime import datetime, timedelta\n\ndate = datetime.strptime('2022-01-31', '%Y-%m-%d')\nnew_date = date + timedelta(days=30)\nresult = new_date.strftime('%B %d, %Y')\nprint(result)</python><result>March 02, 2022</result></code> March 1, 2022. Note that February only has 28 days in non-leap years, so the 30th day would fall on March 1, 2022 in a leap year."

Figure 3: Inconsistency between LLMs' outputs and the returned results from Python interpreter.

By filtering out the entries in \mathcal{C} that show inconsistencies between the Python execution results and the following content, we create the ToolBridge dataset. This refined dataset is then used to train the Llama3 model, yielding outputs consistent with the Python results. Here are some examples.

<i>Consistent outputs from Llama3 trained on ToolBridge</i>
<i>Example 1:</i>
User: "Count the number of vowels in the sentence 'This is a simple sentence'."
Assistant: "There are <code><python>sentence = 'This is a simple sentence'\nvowels = ['a', 'e', 'i', 'o', 'u']\ncount = sum(1 for char in sentence.lower() if char in vowels)\nprint(count)</python><result>8</result></code> 8 vowels in the sentence 'This is a simple sentence'."
<i>Example 2:</i>
User: "What is the date 30 days after January 31, 2022?"
Assistant: "The date 30 days after January 31, 2022 is <code><python>from datetime import datetime, timedelta\n\ndate = datetime.strptime('2022-01-31', '%Y-%m-%d')\nnew_date = date + timedelta(days=30)\nresult = new_date.strftime('%B %d, %Y')\nprint(result)</python><result>March 02, 2022</result></code> February 30, 2022. However, this date is not valid as February only has 28 days in a non-leap year. So, the date would be March 02, 2022."

Figure 4: Consistency between LLMs' outputs and the returned results from Python interpreter.

A.6 LINKS TO DATA SOURCES USED IN THE DATASET POOL

To avoid any confusion regarding the datasets, the download links for each dataset are also provided in Table 11. Note that, links to all datasets from the dataset pool \mathcal{P} are presented in the table, even those not incorporated into ToolBridge.

Dataset	Source Link
School Math 0.25M	https://huggingface.co/datasets/BelleGroup/school_math.0.25M
ChatAlpaca	https://github.com/cascip/ChatAlpaca
Platypus	https://huggingface.co/datasets/garage-bAInd/Open-Platypus
ShareGPT90K	https://huggingface.co/datasets/liyucheng/ShareGPT90K
WizardLM_Orca	https://huggingface.co/datasets/pankajmathur/WizardLM_Orca
WizardLM_evol_instruct_70k	https://huggingface.co/datasets/WizardLMTeam/WizardLM_evol_instruct_70k
WizardLM_evol_instruct_V2	https://huggingface.co/datasets/WizardLMTeam/WizardLM_evol_instruct_V2.196k
OpenOrca	https://huggingface.co/datasets/Open-Orca/OpenOrca
TigerBot	https://huggingface.co/datasets/TigerResearch/sft_en , TigerResearch/sft_zh
GPT-4all	https://huggingface.co/datasets/nomic-ai/gpt4all-j-prompt-generations
COIG	https://huggingface.co/datasets/BAAI/COIG
LIMA	https://huggingface.co/datasets/GAIR/lima
AlpacaDataCleaned	https://huggingface.co/datasets/yahma/alpaca-cleaned
GPT-4-LLM	https://github.com/Instruction-Tuning-with-GPT-4/GPT-4-LLM
Bactrian-X	https://huggingface.co/datasets/MBZUAI/Bactrian-X
code_instructions_120k_alpaca	https://huggingface.co/datasets/iamtarun/code_instructions_120k_alpaca
TSI-v0	https://huggingface.co/datasets/tasksource/tasksource-instruct-v0
Alpaca	https://github.com/tatsu-lab/stanford_alpaca
No Robots	https://huggingface.co/datasets/HuggingFaceH4/no_robots
Baize	https://github.com/project-baize/baize-chatbot
LaMini-Instruction	https://huggingface.co/datasets/MBZUAI/LaMini-instruction
tiny-codes	https://huggingface.co/datasets/nampdn-ai/tiny-codes
self-instruct	https://github.com/yizhongw/self-instruct
ign_clean_instruct_dataset_500k	https://huggingface.co/datasets/ignmilton/ign_clean_instruct_dataset_500k
MOSS SFT	https://github.com/OpenMOSS/MOSS

Table 11: Source links for the datasets utilized in the dataset pool.

Here is a brief description of each dataset,

- School Math 0.25M: It includes approximately 250,000 Chinese math problems generated by the BELLE project, along with their solution processes.
- ChatAlpaca: ChatAlpaca is a comprehensive dataset created to assist researchers in building models for instruction-following across multi-turn conversations. It expands upon the Stanford Alpaca dataset by incorporating a broader range of multi-turn instructions and their corresponding responses.
- Platypus: This dataset is designed to enhance the logical reasoning capabilities of LLMs and was instrumental in training the Platypus2 models. It combines eleven public datasets, carefully curated through keyword filtering and further refined using Sentence Transformers to exclude questions with more than 80% similarity.
- ShareGPT90K: A high quality dataset generated by using GPT-4.
- WizardLM_Orca: Enhanced WizardLM dataset, generated using the Orca methodology.
- WizardLM_evol_instruct_70k: This is the training data of WizardLM.
- WizardLM_evol_instruct_V2: The dataset contains 143K mixed evolved data derived from Alpaca and ShareGPT. It represents the latest optimized version of Evol-Instruct training data for the WizardLM model.
- OpenOrca: A collection of augmented FLAN data, generated using the methodology described in the Orca paper.
- TigerBot: Datasets used to train TigerBot include pretraining data, STF data, and domain-specific datasets such as financial research reports.
- GPT-4all: A curated mix of subsets from OIG, P3, and StackOverflow, focusing on topics such as general question-answering and customized creative queries.
- COIG: A Chinese-focused dataset encompassing domains such as general-purpose QA, Chinese exams, and coding. Its quality has been verified by human annotators.
- LIMA: High quality SFT dataset used by LIMA.

- 972 • AlpacaDataCleaned: An improved and cleaned iteration of the Alpaca, GPT_LLM, and
973 GPTeacher datasets.
- 974 • GPT-4-LLM: It is generated using GPT-4 and other LLMs to produce improved pairs and
975 data for RLHF.
- 976 • Bactrian-X: A multilingual adaptation of the Alpaca and Dolly-15K datasets.
- 977 • code_instructions_120k_alpaca: Code instruction data formatted for instruction fine-tuning.
- 978 • TSI-v0: A multi-task instruction-tuning dataset derived from 475 Tasksource datasets, de-
979 signed in a manner similar to the Flan and Natural Instructions datasets.
- 980 • Alpaca: It consists of 52K instruction-following examples, specifically designed for fine-
981 tuning the Alpaca model.
- 982 • No Robots: High-quality, human-generated STF data in a single-turn format.
- 983 • Baize: A dialogue dataset generated by GPT-4 through self-talking, with questions and
984 topics sourced from Quora, StackOverflow, and various medical knowledge bases.
- 985 • LaMini-Instruction: A dataset distilled from the FLAN collection, P3, and Self-Instruct.
- 986 • tiny-codes: This synthetic dataset comprises 1.6 million concise and clear code snippets,
987 designed to help LLM models develop reasoning skills in both natural and programming
988 languages. The dataset spans a wide range of programming languages, including Python,
989 TypeScript, JavaScript, Ruby, Julia, Rust, C++, Bash, Java, C#, and Go.
- 990 • self-instruct: This dataset is generated using the methodology outlined in Self-Instruct:
991 Aligning Language Models with Self-Generated Instructions.
- 992 • ign_clean_instruct_dataset_500k: A large-scale SFT dataset synthetically generated from a
993 subset of Ultrachat prompts.
- 994 • MOSS SFT: A conversational dataset curated and developed by the MOSS team, with each
995 entry annotated with labels for usefulness, loyalty, and harmlessness.
- 996
- 997
- 998
- 999
- 1000
- 1001
- 1002
- 1003
- 1004
- 1005
- 1006
- 1007
- 1008
- 1009
- 1010
- 1011
- 1012
- 1013
- 1014
- 1015
- 1016
- 1017
- 1018
- 1019
- 1020
- 1021
- 1022
- 1023
- 1024
- 1025

1026 A.7 TEMPLATES FOR YIELDING RANDOMQA DATASET
 1027

1028 The complete list of templates used to generate the RandomQA dataset is as follows,
 1029

```

1030 1 '''RandomQAGenerator'''
1031 2 class RandomQAGenerator():
1032 3     question_types_data_processing = [
1033 4         "Sort an array in ascending order",
1034 5         "Transpose a 2D matrix",
1035 6         "Reverse the string",
1036 7         "Extract first N elements in a list",
1037 8         "Reverse the order of elements in a list",
1038 9         "Count the frequency of one character in a string",
103910        "Find the intersection of two strings",
103911        "Find the length of the longest word in a string",
104012        "Count the number of vowels in a string",
104113        "Convert a list of Celsius temperatures to Fahrenheit",
104214        "Calculate time difference between two time zones",
104315        "Find the leap year after a year",
104416        "Find the most common word in a paragraph",
104417        "Find the first recurring word in a string",
104518        "Extract all the numbers in a string",
104619        "Convert a decimal number to its binary equivalent",
104720        "Calculate the difference between two lists",
104821        "Find out all the numbers that are not unique",
104922        "Flatten a 2D list into a 1D list",
105023        "Remove duplicates from a list",
105124        "Filter elements in a list based on a condition",
105225        "Merge two dictionaries into one",
105326        "Extract all words of a specific length from a text",
105427        "Extract email addresses from a text",
105528        "Sort a list of strings by their length",
105629        "Check if two strings are anagrams",
105730        "Extract hashtags from a social media post",
105831        "Capitalize each word in a string",
105932        "Find the index of a substring in a string",
106033        "Replace all vowels in a string with a specific character",
106134    ]
106235    question_types_numerical_computation = [
106336        "Calculate the average of an array",
106437        "Find the maximum and minimum values of an array",
106538        "Calculate the dot product of two arrays",
106639        "Generate a set of random integers and find their sum",
106740        "Generate the smallest prime number greater than x",
106841        "Calculate the standard deviation of a list of floating-point
106942        ↪ numbers",
107043        "Generate a random matrix and find its inverse",
107144        "Find the median of an array",
107245        "Generate Fibonacci sequence up to n-th term",
107346        "Find the GCD (Greatest Common Divisor) of two numbers",
107447        "Calculate the factorial of a number",
107548        "Find the mode of a list of numbers",
107649        "Calculate the sum of even numbers in a list",
107750        "Calculate the cumulative sum of an array",
107851        "Calculate cosine value",
107952        "Square every number in a list",
108053        "Calculate the sum of squares of numbers in an array",
108154        "Find the n-th smallest number in an array",
108255        "Calculate the Euclidean distance between two points in a plane",
108356        "Calculate the compound interest given principal, rate, and
108457        ↪ time",
108558        "Calculate the perimeter of a rectangle given its length and
108659        ↪ width",
108760        "Sum all the digits of a given number",
  
```

```

1080     "Calculate the area of a triangle given its base and height",
1081     "Find the real roots of a quadratic equation",
1082     "Calculate the sum of the cubes of a list",
1083     "Round all elements in a list to two decimal places",
1084     "Calculate the hypotenuse of a right triangle given the other two
1085     ↪ sides",
1086     "Sum all odd numbers in a list",
1087     "Generate the smallest N primes",
1088     "Find the sum of all elements above the main diagonal of a
1089     ↪ matrix"
1090 ]
1091 def __init__(self, num_gen_qa=1000):
1092     self.num_gen_qa = num_gen_qa
1093     '''generate'''
1094 def generate(self):
1095     qa_pairs = []
1096     for _ in range(self.num_gen_qa):
1097         ↪ qa_pairs.append(self.randomgenone(self.question_types_data_processing))
1098     pickle.dump(qa_pairs,
1099     ↪ open(f'random_qa_dp_{int(time.time())}.pkl', 'wb'))
1100     time.sleep(1)
1101     qa_pairs = []
1102     for _ in range(self.num_gen_qa):
1103         ↪ qa_pairs.append(self.randomgenone(self.question_types_data_processing))
1104     pickle.dump(qa_pairs,
1105     ↪ open(f'random_qa_dp_{int(time.time())}.pkl', 'wb'))
1106     time.sleep(1)
1107     qa_pairs = []
1108     for _ in range(self.num_gen_qa):
1109         ↪ qa_pairs.append(self.randomgenone(self.question_types_numerical_computation))
1110     pickle.dump(qa_pairs,
1111     ↪ open(f'random_qa_nc_{int(time.time())}.pkl', 'wb'))
1112     time.sleep(1)
1113     qa_pairs = []
1114     for _ in range(self.num_gen_qa):
1115         ↪ qa_pairs.append(self.randomgenone(self.question_types_numerical_computation))
1116     pickle.dump(qa_pairs,
1117     ↪ open(f'random_qa_nc_{int(time.time())}.pkl', 'wb'))
1118     '''randomgenone'''
1119 def randomgenone(self, question_types):
1120     # randomly choose a question type
1121     question_type = random.choice(question_types)
1122     # generate question and answer based on type
1123     # 1. Calculate the average of an array
1124     if question_type == "Calculate the average of an array":
1125         array = [round(random.uniform(-10000, 10000)) for _ in
1126         ↪ range(random.randint(5, 15))]
1127         question = f"Calculate the average of the array {array} and
1128         ↪ round the result to two decimal places."
1129         answer = round(sum(array) / len(array), 2)
1130     # 2. Find the maximum and minimum values of an array
1131     elif question_type == "Find the maximum and minimum values of an
1132     ↪ array":
1133         array = [round(random.uniform(-10000, 10000)) for _ in
1134         ↪ range(random.randint(5, 15))]
1135         max_or_min = random.choice(['maximum', 'minimum'])
1136         question = f"Find the {max_or_min} value of the array
1137         ↪ {array}, give the result of multiplying it by 7."
1138         answer = max(array) if max_or_min == 'maximum' else
1139         ↪ min(array)
1140         answer = answer * 7

```

```

1134
1135 # 3. Calculate the dot product of two arrays
1136 elif question_type == "Calculate the dot product of two arrays":
1137     length = random.randint(5, 15)
1138     array1 = [random.randint(20, 1000) for _ in range(length)]
1139     array2 = [random.randint(20, 1000) for _ in range(length)]
1140     question = f"Calculate the dot product of the arrays {array1}
1141     ↪ and {array2}."
1142     answer = sum(x * y for x, y in zip(array1, array2))
1143
1144 # 4. Sort an array in ascending order
1145 elif question_type == "Sort an array in ascending order":
1146     array = [random.randint(-10000, 10000) for _ in
1147     ↪ range(random.randint(5, 15))]
1148     question = f"Sort the array {array} in ascending order."
1149     answer = sorted(array)
1150
1151 # 5. Generate a set of random integers and find their sum
1152 elif question_type == "Generate a set of random integers and find
1153 ↪ their sum":
1154     array = [random.randint(1000, 100000) for _ in
1155     ↪ range(random.randint(5, 15))]
1156     question = f"Here is a set of random integers {array}, please
1157     ↪ find their sum."
1158     answer = sum(array)
1159
1160 # 6. Generate the smallest prime number greater than x
1161 elif question_type == "Generate the smallest prime number greater
1162 ↪ than x":
1163     num = random.randint(2000, 100000)
1164     question = f"Generate the smallest prime number greater than
1165     ↪ {num}."
1166     answer = nextprime(num)
1167
1168 # 7. Calculate the standard deviation of a list of floating-point
1169 ↪ numbers
1170 elif question_type == "Calculate the standard deviation of a list
1171 ↪ of floating-point numbers":
1172     array = [round(random.uniform(10, 1000), 2) for _ in
1173     ↪ range(random.randint(5, 15))]
1174     mean = sum(array) / len(array)
1175     variance = sum((x - mean) ** 2 for x in array) / len(array)
1176     question = f"Calculate the standard deviation of the array
1177     ↪ {array} and round the result to two decimal places."
1178     answer = round(variance ** 0.5, 2)
1179
1180 # 8. Generate a random matrix and find its inverse
1181 elif question_type == "Generate a random matrix and find its
1182 ↪ inverse":
1183     matrix_len = random.randint(2, 10)
1184     matrix = [[random.randint(1, 1000) for _ in
1185     ↪ range(matrix_len)] for _ in range(matrix_len)]
1186     question = f"Here is a random matrix {matrix}, please find
1187     ↪ its inverse, you can answer with 'not invertible' if its
1188     ↪ inverse does not exist."
1189     det = np.linalg.det(matrix)
1190     if int(det) != 0:
1191         inv_matrix = np.linalg.inv(matrix).tolist()
1192     else:
1193         inv_matrix = "not invertible"
1194     answer = inv_matrix
1195
1196 # 9. Count the frequency of one character in a string
1197 elif question_type == "Count the frequency of one character in a
1198 ↪ string":
1199     char = random.choice('abcdefghijklmnopqrstuvwxyz')
1200     string = ''.join(random.choices('abcdefghijklmnopqrstuvwxyz',
1201     ↪ k=random.randint(50, 100))) + char * 101
1202     question = f"Count the frequency of character {char} in the
1203     ↪ string '{string}'."
1204     answer = string.count(char)
1205
1206 # 10. Square every number in a list

```

```

1188 elif question_type == "Square every number in a list":
1189     array = [random.randint(1, 10000) for _ in
1190             range(random.randint(5, 15))]
1191     question = f"Square every number in the list {array}."
1192     answer = [x ** 2 for x in array]
1193     # 11. Find the median of an array
1194     elif question_type == "Find the median of an array":
1195         array = [random.randint(200000, 10000000) for _ in
1196                 range(random.randint(5, 15))]
1197         sorted_array = sorted(array)
1198         question = f"Find the median of the array {array}, give the
1199                 result of multiplying it by 9."
1200         answer = sorted_array[len(sorted_array) // 2]
1201         answer = answer * 9
1202     # 12. Generate Fibonacci sequence up to n-th term
1203     elif question_type == "Generate Fibonacci sequence up to n-th
1204             term":
1205         n = random.randint(5, 20)
1206         question = f"Generate the Fibonacci sequence up to the {n}-th
1207                 term."
1208         fib = [0, 1]
1209         for i in range(2, n):
1210             fib.append(fib[-1] + fib[-2])
1211         answer = fib
1212     # 13. Transpose a 2D matrix
1213     elif question_type == "Transpose a 2D matrix":
1214         matrix_len = random.randint(2, 10)
1215         matrix = [[random.randint(-1000, 1000) for _ in
1216                  range(matrix_len)] for _ in range(matrix_len)]
1217         question = f"Transpose the matrix {matrix}."
1218         answer = [list(row) for row in zip(*matrix)]
1219     # 14. Reverse the string
1220     elif question_type == "Reverse the string":
1221         string = ''.join(random.choices('abcdefghijklmnopqrstuvwxyz',
1222                                     k=random.randint(10, 20)))
1223         question = f"Reverse the string {string}, and splice it
1224                 behind the string 'appleiphone'."
1225         answer = 'appleiphone' + string[::-1]
1226     # 15. Find the GCD (Greatest Common Divisor) of two numbers
1227     elif question_type == "Find the GCD (Greatest Common Divisor) of
1228             two numbers":
1229         answer = 0
1230         while answer <= 100:
1231             a, b = random.randint(200, 1000000), random.randint(200,
1232                     1000000)
1233             question = f"Find the GCD of the numbers {a} and {b}."
1234             answer = math.gcd(a, b)
1235     # 16. Calculate the factorial of a number
1236     elif question_type == "Calculate the factorial of a number":
1237         num = random.randint(10, 100)
1238         question = f"Calculate the factorial of {num}."
1239         answer = math.factorial(num)
1240     # 17. Find the mode of a list of numbers
1241     elif question_type == "Find the mode of a list of numbers":
1242         array = [random.randint(113333, 113343) for _ in range(15)]
1243         question = f"Find the mode of the array {array}, give the
1244                 result of multiplying it by 3."
1245         answer = max(set(array), key=array.count)
1246         answer = answer * 3
1247     # 18. Calculate the sum of even numbers in a list
1248     elif question_type == "Calculate the sum of even numbers in a
1249             list":
1250         array = [random.randint(1000, 1000000) for _ in
1251                 range(random.randint(10, 25))]

```

```

1242
205      question = f"Calculate the sum of even numbers in the list
1243      ↪ {array}."
1244
206      answer = sum(x for x in array if x % 2 == 0)
1245
207      # 19. Calculate the cumulative sum of an array
1246      elif question_type == "Calculate the cumulative sum of an array":
1247
209          array = [random.randint(1, 10000) for _ in
1248          ↪ range(random.randint(5, 15))]
210
211          question = f"Calculate the cumulative sum of the array
1249          ↪ {array}."
1250
212          answer = [sum(array[:i+1]) for i in range(len(array))]
1251
213      # 20. Extract first N elements in a list
1252      elif question_type == "Extract first N elements in a list":
214
215          N = random.randint(5, 10)
1253
216          array = [random.randint(1, 10000) for _ in
1254          ↪ range(random.randint(15, 35))]
1255
217          question = f"Extract first {N} elements in the list {array}
1256          ↪ and then plus 7 for each element in the sub-list."
1257
218          answer = array[:N]
1258
219          answer = [a + 7 for a in answer]
1259
220      # 21. Calculate cosine value
1260      elif question_type == "Calculate cosine value":
1261
221          degree = random.randint(0, 360) + 0.5
1262
222          question = f"Calculate cosine value for {degree} degree and
1263          ↪ round the result to two decimal places."
1264
225          answer = round(math.cos(math.radians(degree)), 2)
1265
226      # 22. Reverse the order of elements in a list
1266      elif question_type == "Reverse the order of elements in a list":
1267
228          array = [random.randint(1, 10000) for _ in
1268          ↪ range(random.randint(5, 15))]
1269
229          question = f"Reverse the order of the elements in the list
1270          ↪ {array} and then plus 3 for each element."
1271
230          answer = array[::-1]
1272
231          answer = [a + 3 for a in answer]
1273
232      # 23. Calculate the sum of squares of numbers in an array
1274      elif question_type == "Calculate the sum of squares of numbers in
1275      ↪ an array":
1276
233          array = [random.randint(10, 10000) for _ in
1277          ↪ range(random.randint(5, 15))]
1278
234          question = f"Calculate the sum of squares of the numbers in
1279          ↪ the array {array}."
1280
235          answer = sum(x ** 2 for x in array)
1281
236      # 24. Find the n-th smallest number in an array
1282      elif question_type == "Find the n-th smallest number in an
1283      ↪ array":
1284
237          array = [random.randint(1000, 10000000) for _ in
1285          ↪ range(random.randint(5, 15))]
1286
238          n = random.randint(1, len(array))
1287
239          question = f"Find the {n}-th smallest number in the array
1288          ↪ {array}, give the result of multiplying it by 3."
1289
240          answer = sorted(array)[n - 1] * 3
1290
241      # 25. Calculate the Euclidean distance between two points in a
1291      ↪ plane
1292      elif question_type == "Calculate the Euclidean distance between
1293      ↪ two points in a plane":
1294
243          x1, y1 = round(random.uniform(-100, 100), 2),
1295          ↪ round(random.uniform(-100, 100), 2)
1296
244          x2, y2 = round(random.uniform(-100, 100), 2),
1297          ↪ round(random.uniform(-100, 100), 2)
1298
245          question = f"Calculate the Euclidean distance between points
1299          ↪ ({x1}, {y1}) and ({x2}, {y2}), round the result to two
1300          ↪ decimal places."
1301
246          answer = round(math.sqrt((x2 - x1)**2 + (y2 - y1)**2), 2)
1302
247      # 26. Find the intersection of two strings
1303      elif question_type == "Find the intersection of two strings":

```



```

1296
1297         str1 = ''.join(random.choices('abcdefghijklmnopqrstuvwxyz',
1298                                     ↪ k=random.randint(50, 100)))
1299         str2 = ''.join(random.choices('abcdefghijklmnopqrstuvwxyz',
1300                                     ↪ k=random.randint(50, 100)))
1301         question = f"Find the intersection of string '{str1}' and
1302                                     ↪ string '{str2}'."
1303         answer = ''.join(set(str1) & set(str2))
1304     # 27. Calculate the compound interest given principal, rate, and
1305     ↪ time
1306     elif question_type == "Calculate the compound interest given
1307     ↪ principal, rate, and time":
1308         principal = random.randint(1000, 10000)
1309         rate = round(random.uniform(1, 10), 2)
1310         time = random.randint(1, 5)
1311         question = f"Calculate the compound interest for principal
1312                 ↪ {principal}, rate {rate}%, and time {time} years, round
1313                 ↪ the result to two decimal places."
1314         answer = round(principal * (1 + rate/100)**time, 2)
1315     # 28. Find the length of the longest word in a string
1316     elif question_type == "Find the length of the longest word in a
1317     ↪ string":
1318         words = [''.join(random.choices('abcdefghijklmnopqrstuvwxyz',
1319                                         ↪ k=random.randint(101, 200))) for _ in
1320                 ↪ range(random.randint(5, 15))]
1321         string = ' '.join(words)
1322         question = f"Find the length of the longest word in the
1323                 ↪ string '{string}'."
1324         answer = max(len(word) for word in words)
1325     # 29. Count the number of vowels in a string
1326     elif question_type == "Count the number of vowels in a string":
1327         string = ''.join(random.choices('abcdefghijklmnopqrstuvwxyz',
1328                                         ↪ k=random.randint(20, 50))) + 'a' * 101
1329         question = f"Count the number of vowels in the string
1330                 ↪ '{string}'."
1331         answer = sum(1 for char in string if char in 'aeiou')
1332     # 30. Convert a list of Celsius temperatures to Fahrenheit
1333     elif question_type == "Convert a list of Celsius temperatures to
1334     ↪ Fahrenheit":
1335         celsius_list = [random.randint(-20, 40) for _ in range(5)]
1336         question = f"Convert the list of Celsius temperatures
1337                 ↪ {celsius_list} to Fahrenheit, round the result to two
1338                 ↪ decimal places."
1339         answer = [round(c * 9/5 + 32, 2) for c in celsius_list]
1340     # 31. Calculate time difference between two time zones
1341     elif question_type == "Calculate time difference between two time
1342     ↪ zones":
1343         tz1, tz2 = random.sample(pytz.all_timezones, 2)
1344         now = datetime.datetime.now()
1345         time1 = pytz.timezone(tz1).localize(now)
1346         time2 = pytz.timezone(tz2).localize(now)
1347         time_difference = abs((time1 - time2).total_seconds())
1348         question = f"Calculate time difference between {tz1} and {tz2}
1349                 ↪ in seconds."
1350         answer = time_difference
1351     # 32. Find the leap year after a year
1352     elif question_type == "Find the leap year after a year":
1353         year = random.randint(1900, 2100)
1354         while calendar.isleap(year):
1355             year = random.randint(1900, 2100)
1356         question = f"Find the leap year after year {year}."
1357         answer = next(y for y in range(year + 1, year + 10000) if
1358                 ↪ calendar.isleap(y))
1359     # 33. Find the most common word in a paragraph
1360     elif question_type == "Find the most common word in a paragraph":

```

```

1350 words = ['apple', 'banana', 'orange', 'grape', 'pear',
1351 ↪ 'hello', 'iphone', 'newspaper']
1352 paragraph = ' '.join(random.choices(words, k=30))
1353 question = f"Find the most common word in the paragraph
1354 ↪ '{paragraph}', concatenate it with the second common word
1355 ↪ in this paragraph."
1356 answer =
1357 ↪ Counter(paragraph.lower().split()).most_common(2)[0][0] +
1358 ↪ Counter(paragraph.lower().split()).most_common(2)[1][0]
1359 # 34. Calculate the perimeter of a rectangle given its length and
1360 ↪ width
1361 elif question_type == "Calculate the perimeter of a rectangle
1362 ↪ given its length and width":
1363 length, width = random.randint(100, 10000),
1364 ↪ random.randint(100, 10000)
1365 question = f"Calculate the perimeter of a rectangle with
1366 ↪ length {length} and width {width}."
1367 answer = 2 * (length + width)
1368 # 35. Sum all the digits of a given number
1369 elif question_type == "Sum all the digits of a given number":
1370 num = int(str(random.randint(100, 99999)) +
1371 ↪ '9999999999999999')
1372 question = f"Sum all the digits of the number {num}."
1373 answer = sum(int(digit) for digit in str(num))
1374 # 36. Calculate the area of a triangle given its base and height
1375 elif question_type == "Calculate the area of a triangle given its
1376 ↪ base and height":
1377 base = round(random.uniform(100, 500), 2)
1378 height = round(random.uniform(100, 500), 2)
1379 question = f"Calculate the area of a triangle with base
1380 ↪ {base} and height {height}, round the result to two
1381 ↪ decimal places."
1382 answer = round(0.5 * base * height, 2)
1383 # 37. Find the real roots of a quadratic equation
1384 elif question_type == "Find the real roots of a quadratic
1385 ↪ equation":
1386 a = round(random.uniform(10, 200), 2)
1387 b = round(random.uniform(10, 200), 2)
1388 c = round(random.uniform(10, 200), 2)
1389 question = f"Find the real roots of the quadratic equation
1390 ↪ {a}x^2 + {b}x + {c} = 0, round the result to two decimal
1391 ↪ places."
1392 discriminant = b**2 - 4*a*c
1393 if discriminant > 0:
1394 root1 = (-b + math.sqrt(discriminant)) / (2*a)
1395 root2 = (-b - math.sqrt(discriminant)) / (2*a)
1396 answer = (round(root1, 2), round(root2, 2))
1397 elif discriminant == 0:
1398 root = -b / (2*a)
1399 answer = round(root, 2)
1400 else:
1401 answer = "no real roots"
1402 # 38. Calculate the sum of the cubes of a list
1403 elif question_type == "Calculate the sum of the cubes of a list":
1404 sequence = [random.randint(100, 10000) for _ in
1405 ↪ range(random.randint(5, 15))]
1406 question = f"Calculate the sum of the cubes of the list
1407 ↪ {sequence}."
1408 answer = sum([n**3 for n in sequence])
1409 # 39. Round all elements in a list to two decimal places
1410 elif question_type == "Round all elements in a list to two
1411 ↪ decimal places":
1412 array = [random.uniform(100, 10000) for _ in
1413 ↪ range(random.randint(5, 15))]

```

```

1404 question = f"Round all elements in the list {array} to two
1405 ↪ decimal places."
1406 answer = [round(num, 2) for num in array]
1407 # 40. Find the first recurring word in a string
1408 elif question_type == "Find the first recurring word in a
1409 ↪ string":
1410 words = [''.join(random.choices('abcdefghijklmnopqrstuvwxy',
1411 ↪ k=random.randint(5, 15))) for _ in
1412 ↪ range(random.randint(5, 10))]
1413 words = words * 3
1414 random.shuffle(words)
1415 paragraph = ' '.join(words)
1416 question = f"Find the first recurring word in the paragraph
1417 ↪ '{paragraph}', concatenate it with the second recurring
1418 ↪ word in this paragraph."
1419 def _find_recurring_words(paragraph):
1420 words = paragraph.lower().split()
1421 seen = set()
1422 first, second = None, None
1423 for word in words:
1424 if word in seen:
1425 if first is None:
1426 first = word
1427 elif second is None and word != first:
1428 second = word
1429 break
1430 seen.add(word)
1431 return first + second
1432 answer = _find_recurring_words(paragraph)
1433 # 41. Calculate the hypotenuse of a right triangle given the
1434 ↪ other two sides
1435 elif question_type == "Calculate the hypotenuse of a right
1436 ↪ triangle given the other two sides":
1437 side1 = random.randint(100, 20000)
1438 side2 = random.randint(100, 20000)
1439 question = f"Calculate the hypotenuse of a right triangle
1440 ↪ with sides {side1} and {side2}, round the result to two
1441 ↪ decimal places."
1442 answer = round(math.sqrt(side1**2 + side2**2), 2)
1443 # 42. Extract all the numbers in a string
1444 elif question_type == "Extract all the numbers in a string":
1445 string1 = random.choices('abcdefghijklmnopqrstuvwxy',
1446 ↪ k=random.randint(20, 50))
1447 string2 = random.choices('0123456789', k=random.randint(20,
1448 ↪ 50))
1449 string = string1 + string2
1450 random.shuffle(string)
1451 string = ''.join(string)
1452 question = f"Extract all the numbers in the string '{string}'
1453 ↪ in order and concatenate them."
1454 answer = ''.join(re.findall(r'\d+', string))
1455 # 43. Convert a decimal number to its binary equivalent
1456 elif question_type == "Convert a decimal number to its binary
1457 ↪ equivalent":
1458 num = random.randint(1000, 1000000)
1459 question = f"Convert the decimal number {num} to its binary
1460 ↪ equivalent."
1461 answer = bin(num)[2:]
1462 # 44. Calculate the difference between two lists
1463 elif question_type == "Calculate the difference between two
1464 ↪ lists":
1465 list1 = [random.randint(1, 50) for _ in range(10)]
1466 list2 = [random.randint(1, 50) for _ in range(10)]
1467 question = f"Calculate the difference between the lists
1468 ↪ {list1} and {list2}."

```

```

1458     answer = list(set(list1) - set(list2))
1459 # 45. Sum all odd numbers in a list
1460 elif question_type == "Sum all odd numbers in a list":
1461     array = [random.randint(1000, 1000000) for _ in
1462             ↪ range(random.randint(5, 15))]
1463     question = f"Sum all the odd numbers in the list {array}."
1464     answer = sum(x for x in array if x % 2 != 0)
1465 # 46. Find out all the numbers that are not unique
1466 elif question_type == "Find out all the numbers that are not
1467     ↪ unique":
1468     array = [random.randint(20, 35) for _ in range(20)]
1469     question = f"Find out all the numbers that are not unique in
1470     ↪ the array {array}."
1471     answer = [num for num, count in Counter(array).items() if
1472             ↪ count > 1]
1473 # 47. Flatten a 2D list into a 1D list
1474 elif question_type == "Flatten a 2D list into a 1D list":
1475     array_len = random.randint(2, 10)
1476     array = [[random.randint(1, 1000) for _ in range(array_len)]
1477             ↪ for _ in range(array_len)]
1478     question = f"Flatten the 2D list {array} into a 1D list."
1479     answer = [item for sublist in array for item in sublist]
1480 # 48. Remove duplicates from a list
1481 elif question_type == "Remove duplicates from a list":
1482     array = [random.randint(1, 20) for _ in range(15)]
1483     while len(array) == len(set(array)):
1484         array = [random.randint(1, 20) for _ in range(15)]
1485     question = f"Remove duplicates from the list {array}."
1486     answer = list(set(array))
1487 # 49. Generate the smallest N primes
1488 elif question_type == "Generate the smallest N primes":
1489     n = random.randint(5, 20)
1490     primes = []
1491     candidate = 2
1492     while len(primes) < n:
1493         if all(candidate % i != 0 for i in range(2, int(candidate
1494             ↪ ** 0.5) + 1)):
1495             primes.append(candidate)
1496             candidate += 1
1497     question = f"Generate the smallest {n} prime numbers."
1498     answer = primes
1499 # 50. Find the sum of all elements above the main diagonal of a
1500     ↪ matrix
1501 elif question_type == "Find the sum of all elements above the
1502     ↪ main diagonal of a matrix":
1503     matrix_len = random.randint(2, 10)
1504     matrix = [[random.randint(1000, 1000000) for _ in
1505             ↪ range(matrix_len)] for _ in range(matrix_len)]
1506     question = f"Find the sum of all elements above the main
1507     ↪ diagonal of the matrix {matrix}."
1508     answer = sum(matrix[i][j] for i in range(matrix_len) for j in
1509             ↪ range(i + 1, matrix_len))
1510 # 51. Filter elements in a list based on a condition
1511 elif question_type == "Filter elements in a list based on a
1512     ↪ condition":
1513     array = [random.randint(-100, 100) for _ in
1514             ↪ range(random.randint(10, 20))]
1515     condition = random.randint(-50, 50)
1516     question = f"Filter all elements in the array {array} that
1517     ↪ are greater than {condition}."
1518     answer = [x for x in array if x > condition]
1519 # 52. Merge two dictionaries into one
1520 elif question_type == "Merge two dictionaries into one":
1521     dict1 = {chr(65 + i): random.randint(1, 100) for i in
1522             ↪ range(random.randint(10, 20))}

```

```

1512 dict2 = {chr(67 + i): random.randint(1, 100) for i in
1513         ↪ range(random.randint(10, 20))}
1514 question = f"Merge the dictionaries {dict1} and {dict2},
1515         ↪ summing values for duplicate keys."
1516 answer = {k: dict1.get(k, 0) + dict2.get(k, 0) for k in
1517         ↪ set(dict1) | set(dict2)}
1518 # 53. Extract all words of a specific length from a text
1519 elif question_type == "Extract all words of a specific length
1520     ↪ from a text":
1521     text = '
1522         ↪ ''.join([''.join(random.choices('abcdefghijklmnopqrstuvwxyz',
1523         ↪ k=random.randint(5, 10)) for _ in
1524         ↪ range(random.randint(10, 20)))
1525         ↪ length = random.randint(5, 10)
1526         ↪ question = f"Find all words in the text '{text}' that have
1527         ↪ exactly {length} characters."
1528         ↪ answer = [word for word in text.split() if len(word) ==
1529         ↪ length]
1530 # 54. Extract email addresses from a text
1531 elif question_type == "Extract email addresses from a text":
1532     answer = [Faker().email() for _ in range(random.randint(2,
1533     ↪ 4))]
1534     text = answer +
1535         ↪ [''.join(random.choices('abcdefghijklmnopqrstuvwxyz',
1536         ↪ k=random.randint(5, 10)) for _ in
1537         ↪ range(random.randint(10, 20))]
1538     random.shuffle(text)
1539     text = ' '.join(text)
1540     question = f"Find all email addresses in the text: '{text}'"
1541 # 55. Sort a list of strings by their length
1542 elif question_type == "Sort a list of strings by their length":
1543     strings =
1544         ↪ [''.join(random.choices('abcdefghijklmnopqrstuvwxyz',
1545         ↪ k=random.randint(5, 20)) for _ in
1546         ↪ range(random.randint(10, 20))]
1547     question = f"Sort the list {strings} by the length of each
1548     ↪ string."
1549     answer = sorted(strings, key=len)
1550 # 56. Check if two strings are anagrams
1551 elif question_type == "Check if two strings are anagrams":
1552     string1 = random.choices('abcdefghijklmnopqrstuvwxyz',
1553     ↪ k=random.randint(10, 20))
1554     string2 = random.choices('abcdefghijklmnopqrstuvwxyz',
1555     ↪ k=random.randint(10, 20)) if random.random() > 0.5 else
1556     ↪ string1
1557     random.shuffle(string2)
1558     string1 = ''.join(string1)
1559     string2 = ''.join(string2)
1560     question = f"Check if '{string1}' and '{string2}' are
1561     ↪ anagrams."
1562     answer = sorted(string1) == sorted(string2)
1563 # 57. Extract hashtags from a social media post
1564 elif question_type == "Extract hashtags from a social media
1565     ↪ post":
1566     topic = [''.join(random.choices('abcdefghijklmnopqrstuvwxyz',
1567     ↪ k=random.randint(5, 10)) for _ in
1568     ↪ range(random.randint(10, 20))]
1569     hashtags = ['#' +
1570         ↪ ''.join(random.choices('abcdefghijklmnopqrstuvwxyz',
1571         ↪ k=random.randint(5, 10)) for _ in
1572         ↪ range(random.randint(2, 5))]
1573     text = topic + hashtags
1574     random.shuffle(text)
1575     text = ' '.join(text)
1576     question = f"Extract all hashtags from the post: '{text}'"

```

```

1566         answer = [word for word in text.split() if
1567                    ↪ word.startswith("#")]
1568     # 58. Capitalize each word in a string
1569     elif question_type == "Capitalize each word in a string":
1570         text = '
1571             ↪ '.join([''.join(random.choices('abcdefghijklmnopqrstuvwxyz',
1572             ↪ k=random.randint(5, 10))) for _ in
1573             ↪ range(random.randint(10, 20))])
1574         question = f"Capitalize each word in the string '{text}'."
1575         answer = text.title()
1576     # 59. Find the index of a substring in a string
1577     elif question_type == "Find the index of a substring in a
1578     ↪ string":
1579         string =
1580             ↪ [''.join(random.choices('abcdefghijklmnopqrstuvwxyz',
1581             ↪ k=random.randint(5, 10))) for _ in
1582             ↪ range(random.randint(10, 20))]
1583         substring = random.choice(string)
1584         string = ' '.join(string)
1585         question = f"Find the index of the substring '{substring}' in
1586         ↪ the string '{string}'."
1587         answer = string.find(substring)
1588     # 60. Replace all vowels in a string with a specific character
1589     elif question_type == "Replace all vowels in a string with a
1590     ↪ specific character":
1591         string = '
1592             ↪ '.join([''.join(random.choices('abcdefghijklmnopqrstuvwxyz',
1593             ↪ k=random.randint(5, 10))) for _ in
1594             ↪ range(random.randint(10, 20))])
1595         replacement = random.choice(["*", "$", "%", "&", "#", "@"])
1596         question = f"Replace all vowels in the string '{string}' with
1597         ↪ '{replacement}'."
1598         answer = ''.join([replacement if char.lower() in "aeiou" else
1599         ↪ char for char in string])
1600     # not defined question
1601     else:
1602         raise ValueError(f'{question_type} is not defined')
1603     # format and return
1604     random_qa = {'question': question, 'answer': answer}
1605     return random_qa

```

1600

1601

1602

1603

1604

1605

1606

1607

1608

1609

1610

1611

1612

1613

1614

1615

1616

1617

1618

1619

1620 A.8 PROMPTS FOR CONSTRUCTING FACT
16211622 We construct the FACT datasets by prompting GPT-4o with,
1623

- 1624 1. Generate 100 Q&A pairs for LLM factual retrieval testing. The
1625 ↪ question topic should be related with Geography. Return them
1626 ↪ as a Python dictionary, with concise answers (3-5 words).
- 1627 2. Generate 100 Q&A pairs for LLM factual retrieval testing. The
1628 ↪ question topic should be related with History. Return them as
1629 ↪ a Python dictionary, with concise answers (3-5 words).
- 1630 3. Generate 100 Q&A pairs for LLM factual retrieval testing. The
1631 ↪ question topic should be related with Science. Return them as
1632 ↪ a Python dictionary, with concise answers (3-5 words).
- 1633 4. Generate 100 Q&A pairs for LLM factual retrieval testing. The
1634 ↪ question topic should be related with Technology. Return them
1635 ↪ as a Python dictionary, with concise answers (3-5 words).
- 1636 5. Generate 100 Q&A pairs for LLM factual retrieval testing. The
1637 ↪ question topic should be related with Mathematics. Return them
1638 ↪ as a Python dictionary, with concise answers (3-5 words).
- 1639 6. Generate 100 Q&A pairs for LLM factual retrieval testing. The
1640 ↪ question topic should be related with Culture and Arts. Return
1641 ↪ them as a Python dictionary, with concise answers (3-5 words).
- 1642 7. Generate 100 Q&A pairs for LLM factual retrieval testing. The
1643 ↪ question topic should be related with Sports. Return them as a
1644 ↪ Python dictionary, with concise answers (3-5 words).
- 1645 8. Generate 100 Q&A pairs for LLM factual retrieval testing. The
1646 ↪ question topic should be related with Politics. Return them as
1647 ↪ a Python dictionary, with concise answers (3-5 words).
- 1648 9. Generate 100 Q&A pairs for LLM factual retrieval testing. The
1649 ↪ question topic should be related with Language and Grammar.
1650 ↪ Return them as a Python dictionary, with concise answers (3-5
1651 ↪ words).
- 1652 10. Generate 100 Q&A pairs for LLM factual retrieval testing. The
1653 ↪ question topic should be related with Current Affairs. Return
1654 ↪ them as a Python dictionary, with concise answers (3-5 words).
- 1655 11. Generate 100 Q&A pairs for LLM factual retrieval testing. The
1656 ↪ question topic should be related with Entertainment. Return
1657 ↪ them as a Python dictionary, with concise answers (3-5 words).
- 1658 12. Generate 100 Q&A pairs for LLM factual retrieval testing. The
1659 ↪ question topic should be related with Medicine and Health.
1660 ↪ Return them as a Python dictionary, with concise answers (3-5
1661 ↪ words).
- 1662 13. Generate 100 Q&A pairs for LLM factual retrieval testing. The
1663 ↪ question topic should be related with Economics and Business.
1664 ↪ Return them as a Python dictionary, with concise answers (3-5
1665 ↪ words).
- 1666 14. Generate 100 Q&A pairs for LLM factual retrieval testing. The
1667 ↪ question topic should be related with Religion and Mythology.
1668 ↪ Return them as a Python dictionary, with concise answers (3-5
1669 ↪ words).
- 1670 15. Generate 100 Q&A pairs for LLM factual retrieval testing. The
1671 ↪ question topic should be related with General Knowledge.
1672 ↪ Return them as a Python dictionary, with concise answers (3-5
1673 ↪ words).

A.9 SOME PROMISING RESULTS

Here, we show some promising results and compare them with other LLMs.



Figure 5: Some promising results from LLMs trained on ToolBridge.

1728 A.10 EXTEND TO NATURAL WEB DATA
1729

1730 Our algorithm is not strictly reliant on the existing SFT datasets as it is equally applicable to natural
1731 data sourced from the web. Because, by using LLMs like GPT-4o, raw web data can be transformed
1732 into QA pairs, which can then be processed leveraging the proposed pipeline outlined in Section 3,
1733 including the selection, conversion, and filtering of valuable data entries.

1734 For instance, Google’s C4 dataset can be systematically transformed into QA pairs by using GPT-4o
1735 with the following example prompt:

1736
1737 You are an expert assistant tasked with creating concise and
1738 ↪ relevant Question-Answer (QA) pairs based on the content of a
1739 ↪ given web page. Follow these guidelines:

- 1740
1741 1. Content Focus: Read the provided web page content carefully and
1742 ↪ base all questions and answers solely on the information within
1743 ↪ the page. Do not include external knowledge.
1744 2. Question Types: Include a mix of factual, explanatory, and
1745 ↪ inferential questions to test comprehension and understanding.
1746 3. Question Structure: Ensure questions are clear and specific. Use
1747 ↪ diverse formats such as: What/Why/How questions.
1748 4. Answer Structure: Provide direct, accurate, and concise answers.
1749 ↪ Avoid ambiguous or overly lengthy responses.

1750 Example Web Content:

1751
1752 Coffee is one of the most popular beverages in the world. It is made
1753 ↪ from roasted coffee beans, which are seeds of the Coffea plant.
1754 ↪ A standard cup of coffee contains approximately 95 milligrams of
1755 ↪ caffeine. Studies suggest that consuming 400 milligrams of
1756 ↪ caffeine per day is generally safe for most adults. Many people
1757 ↪ drink coffee daily as part of their morning routine, with some
1758 ↪ consuming 2 to 4 cups per day depending on their preference.

1759 Example Output:

- 1760
1761 - Question 1: How much caffeine is in a standard cup of coffee?
1762 - Answer: A standard cup of coffee contains approximately 95
1763 ↪ milligrams of caffeine.
1764 - Question 2: If a person drinks 3 cups of coffee, how much caffeine
1765 ↪ do they consume?
1766 - Answer: They consume 285 milligrams of caffeine ($95 * 3 = 285$).
1767 - Question 3: If a person drinks 2 cups of coffee in the morning and
1768 ↪ 1 in the evening, how much caffeine do they consume in total?
1769 - Answer: They consume 285 milligrams of caffeine ($95 * 2 + 95 =$
1770 ↪ 285).
1771 - Question 4: How many cups of coffee would a person need to drink
1772 ↪ to consume exactly 400 milligrams of caffeine?
1773 - Answer: They would need to drink approximately 4.2 cups of coffee
1774 ↪ ($400 \div 95 = 4.2$).
1775 - Question 5: If a person reduces their daily coffee intake from 4
1776 ↪ cups to 2 cups, how much less caffeine do they consume in a day?
1777 - Answer: They consume 190 milligrams less caffeine ($95 * 4 - 95 * 2$
1778 ↪ $= 380 - 190 = 190$).

1778 Web Content:

1779
1780 PLACEHOLDER

1781 Output:

Some examples of the output results are presented below,

Source data:

Schwarz has been the president of the Greater Cleveland Film Commission for over ten years. He has over thirty years of experience in the film and television industries, working predominantly as a location manager and producer. His work includes "Wild At Heart," "My Name Is Earl," "Entourage," "The Comeback," "From the Earth to the Moon" and Tom Hanks and Steven Spielberg's "Band of Brothers," for which Ivan won a Golden Globe. The Film Commission is the only nonprofit organization dedicated to bringing jobs and economic impact to Northeast Ohio through the growth of a sustainable media production industry. Since passing the Ohio Motion Picture Tax Credit in 2009, more than 110 productions have shot in Cleveland, including "The Avengers," "Captain America: The Winter Soldier," and most recently, Matthew McConaughey's "White Boy Rick" and Helen Hunt's "I See You".

Generated QA:

- Question1: How long has Schwarz been the president of the Greater Cleveland Film Commission?
- Answer: Schwarz has been the president of the Greater Cleveland Film Commission for over ten years.
- Question2: What is the primary goal of the Greater Cleveland Film Commission?
- Answer: The primary goal of the Greater Cleveland Film Commission is to bring jobs and economic impact to Northeast Ohio through the growth of a sustainable media production industry.
- Question3: What significant legislation related to film production was passed in Ohio in 2009?
- Answer: The Ohio Motion Picture Tax Credit was passed in 2009.

Source data:

Earn monthly interest on our Citibank Time Deposits (also known as Fixed Deposits). What's more, you get to enjoy the flexibility of making partial withdrawals before maturity date of your Time Deposit. Partial withdrawals in multiples of RM5,000 before the maturity date. Option to pledge your Time Deposit for overdraft facility. Continue to earn the original interest rate with your remaining balance. View your Citibank® Time Deposit accounts across participating countries with a single login via Online Global View. Online account inquiries and banking transactions via Citibank® or Citigold® Online. Make banking inquiries anytime, anywhere via live e-Chat, 'Call Me' function and 20 seconds1 call pickup. Eligible for protection by PIDM**. * Monthly interest is only available for minimum deposit of RM50,000 and for a minimum tenure of 6 months. ** Click here for more details on PIDM protection. 1From the moment you press 0 to speak to a CitiPhone officer. For the Citibank Account Terms and Conditions, click here.

Generated QA:

- Question1: If a customer deposits RM100,000 into a Citibank Time Deposit account, how many partial withdrawals of RM5,000 can they make before the maturity date?
- Answer: They can make 20 partial withdrawals of RM5,000 each ($RM100,000 \div RM5,000 = 20$).
- Question2: A customer earns monthly interest on a Time Deposit of RM200,000 with an interest rate of 2% per annum. How much interest will they earn in one month?
- Answer: The monthly interest is RM333.33 ($RM200,000 \times 2\% \div 12 = RM333.33$).
- Question3: If a customer deposits RM75,000 for 6 months with an interest rate of 3% per annum, what is the total interest earned at the end of the tenure?
- Answer: The total interest earned is RM1,125 ($RM75,000 \times 3\% \div 12 \times 6 = RM1,125$).
- Question4: A customer deposits RM120,000 for 12 months with a 4% annual interest rate. How much interest will they earn in total?
- Answer: The total interest earned is RM4,800 ($RM120,000 \times 4\% = RM4,800$).
- Question5: If the monthly interest on a Time Deposit of RM60,000 is RM100, what is the annual interest rate?
- Answer: The annual interest rate is 2% ($RM100 \times 12 \div RM60,000 = 2\%$).

Figure 6: Sample outcomes of converting natural web data from Google's C4 dataset into QA format with the help of GPT-4o.

1836 A.11 CONSTRUCT FACT WITH GEMINI

1837
1838 In Section 4.4, we employ GPT-4o, together with the prompts outlined in Appendix A.8, to construct
1839 the FACT dataset, which is utilized to evaluate the extent to which ToolBridge enhances the LLMs’
1840 capability for factual retrieval. Here, we also use Google Gemini to yield three batches of the FACT
1841 datasets termed Gemini-FACT-B(n), with each batch corresponding sequentially to the five prompts
1842 outlined in Appendix A.8. Table 12 demonstrates the evaluation results.

1843 Models	SFT data	Gemini-FACT-B1	Gemini-FACT-B2	Gemini-FACT-B3
1844 Llama3-8B	-	75.8	52.5	60.3
1845 Llama3-8B-Lora	ToolBridge [§]	83.4	61.7	66.2
1846 Llama3-8B-Lora	ToolBridge	89.2	63.3	71.2
1847 Mistral-7B	-	77.5	59.2	67.8
1848 Mistral-7B-Lora	ToolBridge [§]	85.8	61.5	70.4
1849 Mistral-7B-Lora	ToolBridge	90.8	64.7	74.7

1850 Table 12: Experimental results on Gemini-FACT under zero-shot setting.

1851
1852 We can observe that the models trained on ToolBridge consistently achieve superior performance.
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889

A.12 DISTRIBUTION OF PYTHON PACKAGES ADOPTED IN TOOLBRIDGE

Here, we provide the distribution of all Python packages used in ToolBridge.

Python Package	Frequency	Python Package	Frequency	Python Package	Frequency
math	2669	re	2234	sympy	1838
nlTK	1616	datetime	1512	langdetect	489
numpy	271	collections	248	string	246
fractions	209	pandas	205	itertools	121
statistics	118	sklearn	92	io	65
requests	54	difflib	35	calendar	34
bs4	32	os	31	random	23
functools	23	urllib	20	json	17
matplotlib	12	scipy	10	operator	8
xml	8	base64	7	codecs	6
sys	5	PIL	5	bisect	4
csv	3	subprocess	3	cmath	3
time	2	pytz	2	ipaddress	2
decimal	2	unicodedata	2	dateutil	2
pytest	2	enum	2	heapq	2
keyword	1	typing	1	inspect	1
configparser	1	struct	1	ctypes	1
networkx	1	lxml	1	pylab	1
cycler	1	torch	1	html	1
textwrap	1	ast	1	locale	1

Table 13: The usage frequency of all Python packages involved in ToolBridge.

1944 A.13 REVIEW PROCESS FOR SFT DATASETS IN TABLE 1
1945

1946 The integration of Supervised Fine-tuning (SFT) datasets presented in Table 1 was primarily carried
1947 out by conducting keyword searches (*e.g.*, "supervised fine-tuning dataset", "SFT dataset", "instruc-
1948 tion tuning dataset", and "LLM dataset") across platforms including Google Scholar, Hugging Face
1949 and GitHub. Representative search results include the following resources:

- 1950 • <https://github.com/Zjh-819/LLMDataHub>,
- 1951 • <https://github.com/RenzeLou/awesome-instruction-learning>,
- 1952 • <https://github.com/raunak-agarwal/instruction-datasets>,
- 1953 • <https://github.com/zhilizju/Awesome-instruction-tuning>,
- 1954 • <https://arxiv.org/abs/2402.18041>,
- 1955 • <https://arxiv.org/abs/2402.06196>.

1958 Based on these resources, we conducted a manual review of all referenced datasets, including verify-
1959 ing whether the dataset qualified as an SFT dataset, assessing its open-source availability, identifying
1960 potential overlaps with existing collected datasets, and examining other potential concerns, such as
1961 copyright issues. At last, we derived Table 1.
1962

1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

A.14 THE PROMPT SELECTION STRATEGY FOR LLAMA3-70B

To select an appropriate prompt for LLama3-70B to perform valuable data entries selection, we first designed the following candidate prompts:

Prompt1: Your task is to determine whether you can add calls to a
 ↪ Python API to a piece of text. The calls should help you get
 ↪ information required to complete the text. You only need to
 ↪ respond with "Yes" or "No", "Yes" means you can and "No" means
 ↪ you can't.

Input:

PLACEHOLDER

Output:

Prompt2: Determine if you can add Python API calls to the text to
 ↪ complete it. Respond with "Yes" or "No".

Input:

PLACEHOLDER

Output:

Prompt3: Prompt1 with example inputs outputs in Appendix A.1.

Prompt4: Prompt2 with example inputs outputs in Appendix A.1.

Then, we randomly sample 50 data entries from each dataset involved in Table 1 to serve as the test set for evaluating the effectiveness of these prompts. For each entry, five people label it to determine whether inserting Python code at the appropriate places can assist in yielding its subsequent contents. The final label for each data entry is determined based on the majority rule.

The evaluation results for each prompt are shown in the following table:

Prompt ID	True Positive	False Positive	False Negative	True Negative	Recall	FPR
1	781	452	0	17	100.0%	96.4%
2	781	432	0	37	100.0%	92.1%
3	732	10	49	459	93.7%	2.1%
4	679	25	102	444	86.9%	5.3%

Table 14: Ablation studies on candidate prompts for LLama3-70B.

Given the role of LLama3-70B as outlined in Section 3.2, the central aim at this stage is to maximize recall — safeguarding against the premature discarding of valuable data entries — while maintaining a sufficiently low false positive rate (FPR) to prevent excessive computational overhead for GPT-4o-mini in the subsequent stage.

In Section 3.3, we elaborate that GPT-4o-mini performs a secondary screening of the entries, refining and further filtering data entries deemed valueless. Its higher accuracy and generalizability compared to LLama3-70B make it the principal decision-making LLMs in our pipeline. However, owing to its significantly higher computational cost, it is essential to minimize the volume of data entries passed to GPT-4o-mini.

Based on the evaluation results in Table 14, Prompt 3 demonstrated the most balanced performance for our requirements, achieving high recall (93.7%) and a notably low FPR (2.1%). While Prompt 1 and Prompt 2 achieved perfect recall, their exceedingly high FPRs (96.4% and 92.1%, respectively)

2052 made them unsuitable for the preliminary filtering task, as they would result in an excessive volume
2053 of data entries being forwarded to GPT-4o-mini. Conversely, Prompt 4's lower recall (86.9%) made
2054 it less effective at retaining valuable entries.

2055
2056 Considering the performance of Prompt 3, we determined it was sufficient to meet the requirements
2057 for LLama3-70B's role in the ToolBridge pipeline. Consequently, we selected Prompt 3 as the final
2058 prompt and did not conduct further ablation studies for LLama3-70B prompt design.

2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105