

TOOLFLOW: Boosting LLM Tool-Calling Through Natural and Coherent Dialogue Synthesis

Zezhong Wang^{1*}, Xingshan Zeng^{2†}, Weiwen Liu², Liangyou Li²,
Yasheng Wang², Lifeng Shang², Xin Jiang², Qun Liu², Kam-Fai Wong¹

¹The Chinese University of Hong Kong, ²Huawei Noah's Ark Lab

{zzwang, kfwong}@se.cuhk.edu.hk

{zeng.xingshan, liuweiw8, liliangyou}@huawei.com

{wangyasheng, Shang.Lifeng, Jiang.Xin, qun.liu}@huawei.com

Abstract

Supervised fine-tuning (SFT) is a common method to enhance the tool calling capabilities of Large Language Models (LLMs), with the training data often being synthesized. The current data synthesis process generally involves sampling a set of tools, formulating a requirement based on these tools, and generating the call statements. However, tools sampled randomly lack relevance, making them difficult to combine and thus reducing the diversity of the data. Additionally, current work overlooks the coherence between turns of dialogues, leading to a gap between the synthesized data and real-world scenarios. To address these issues, we propose a Graph-based Sampling strategy to sample more relevant tool combinations, and a Planned-generation strategy to create plans that guide the synthesis of coherent dialogues. We integrate these two strategies and enable multiple agents to synthesize the dialogue data interactively, resulting in our tool-calling data synthesis pipeline TOOLFLOW. Data quality assessments demonstrate improvements in the naturalness and coherence of our synthesized dialogues. Finally, we apply SFT on LLaMA-3.1-8B using 8,000 synthetic dialogues generated with TOOLFLOW. Results show that the model achieves tool-calling performance comparable to or even surpassing GPT-4, while maintaining strong general capabilities.

1 Introduction

Enabling Large Language Models (LLMs) to perform tool calling significantly enhances their capabilities and practical applications. This requires the models to possess strong understanding, reasoning, and instruction-following abilities. Customized fine-tuning is a widely used method to improve the tool-calling capabilities of LLMs (Abdelaziz et al., 2024; Patil et al., 2023; Schick et al., 2023; Qin

et al., 2023). However, access to fine-tuning data can be limited. One viable solution is to utilize LLMs for data synthesis (Basu et al., 2024; Wang et al., 2023; Xu et al., 2023; Yu et al., 2024).

A typical tool-calling data synthesis process involves three steps: (1) selecting candidate tool(s), (2) generating requirements based on those tools, and (3) creating the call statements (Tang et al., 2023; Liu et al., 2024b). However, the data synthesized through this method often lacks realism and naturalness. Randomly sampled tools frequently fail to interconnect, making it difficult to combine them for complex tasks. Consequently, the requirements for subsequent synthesis tend to be simplistic, which reduces the diversity and complexity of the data. Furthermore, much of the existing research focuses solely on generating single-turn tool-calling instructions, neglecting the coherence between dialogue turns (Qin et al., 2023; Yang et al., 2023). In real-world interactions, LLMs typically engage with users through dialogues rather than single-round Q&A sessions. This creates a gap between Q&A-type training data and its practical application, ultimately diminishing the naturalness of the synthesized data.

To address these two challenges, we propose TOOLFLOW, a tool-calling data synthesis pipeline that employs a graph-based sampling algorithm to improve the correlation among the selected tools and a planned-generation strategy to enhance the naturalness and coherence of the synthesized tool call dialogues.

Specifically, we consider tools with similar parameters or return values to be related. For instance, both "book_flight" and "get_weather" require parameters related to location. In practical scenarios, these two tools are indeed interconnected, as they often occur together in travel contexts. Based on this assumption, we construct a tool graph that represents the similarity between parameters and return values of the tools. Each node in the graph

*Work done during internship at Huawei Noah's Ark Lab.

†Corresponding author

represents a tool, while the edges indicate the relevance between pairs of tools. When sampling tools, we randomly select a subgraph from this tool graph, ensuring that the sampled tools are more likely to interact effectively, thereby facilitating the generation of complex requirements.

On the other hand, before synthesizing dialogues, we first have the LLM create a plan based on the selected subset of tools. This plan outlines the requests that users need to make in each turn of the dialogue. While constructing the plan, the model focuses on establishing the dialogue framework without worrying about phrasing and details. This approach allows the model to concentrate on the logical relationships and interactions between requirements, resulting in more coherent demands. Additionally, we enable the LLM to incorporate non-tool-call requests into the plan. This not only enhances the diversity of the conversation content but also facilitates seamless transitions between topics, naturally leading to new requirements.

We generate dialogues using three agents: *User*, *Assistant*, and *Tool*. Based on the selected tool subset and the established plan, these agents interact to complete the dialogue. By iterating through the "sampling-planning-generation" process, we synthesized a total of 8,000 dialogues. To evaluate the effectiveness of our proposed method, we conduct a comprehensive ablation study on the graph-based sampling and planning strategy by generating dialogues of the same size selectively without these modules. We perform a thorough evaluation of the data quality, which demonstrates that TOOLFLOW can effectively enhance the naturalness, coherence, and diversity of the generated dialogues. Finally, we apply supervised fine-tuning to LLaMA-3.1-8B-Instruct (Dubey et al., 2024) using the synthesized data and validate improvements in the model's tool-calling capabilities while preserving its general abilities, with TOOLFLOW.

We summarize our contributions into the following three key points:

- We propose a Graph-based Sampling strategy to select related tools, aiming to enhance the diversity and complexity of synthetic tool calling requirements.
- We introduce a Planned-Generation strategy to improve the naturalness and coherence of synthetic dialogues.
- We integrate these two strategies and propose TOOLFLOW, a tool calling dialogue data syn-

thesis pipeline, with extensive experiments and analyses showing its effectiveness.

2 Related Works

Integrating external tools with large language models (LLMs) significantly broadens their functional scope, allowing for more specialized, precise, and reliable solutions to complex problems (Qin et al., 2023). There are generally two main strategies for embedding tool-use capabilities into LLMs: prompt-based methods and tool-augmented SFT. Prompt-based methods enable LLMs to utilize tools by providing descriptions and examples of the tools in the prompt, without any incremental training (Ruan et al., 2023; Hsieh et al., 2023; Milon et al., 2023). ReAct (Yao et al., 2023) is one of the notable methods within this category. It allows LLMs to switch between reasoning and executing actions to tackle challenging tasks. However, the effectiveness of these approaches can be limited by the inherent capabilities of the model. On the other hand, tool-augmented tuning is attracting increasing interest due to its direct enhancement of LLM's tool usage capabilities (Abdelaziz et al., 2024; Patil et al., 2023; Schick et al., 2023; Qin et al., 2023). As the limited availability of tool calling datasets, Basu et al. (2024) adapted data from various other domains for application in tool calling studies. Others (Liu et al., 2024b; Tang et al., 2023; Qin et al., 2023) mainly synthesized single-turn instructions that involve basic tool calling requirements. However, LLMs typically interact with users through dialogue rather than single-turn Q&A. This mismatch means the data is unnatural, creating a gap with real-world scenarios. Our TOOLFLOW focuses on enhancing the coherence and naturalness of dialogues in data synthesis, making it more aligned with actual applications.

3 Methodology

To generate realistic and coherent dialogues, we propose a three-step data synthesis process: (1) Selecting a tool subset using graph-based sampling; (2) Generating a dialogue plan based on the selected tool subset; (3) Synthesizing dialogues guided by the tool subset and dialogue plan.

3.1 Graph-based Sampling for Tool Selection

Tool calling data synthesis generally starts by selecting one or more tools from the available toolset. While much of the previous work overlooks the

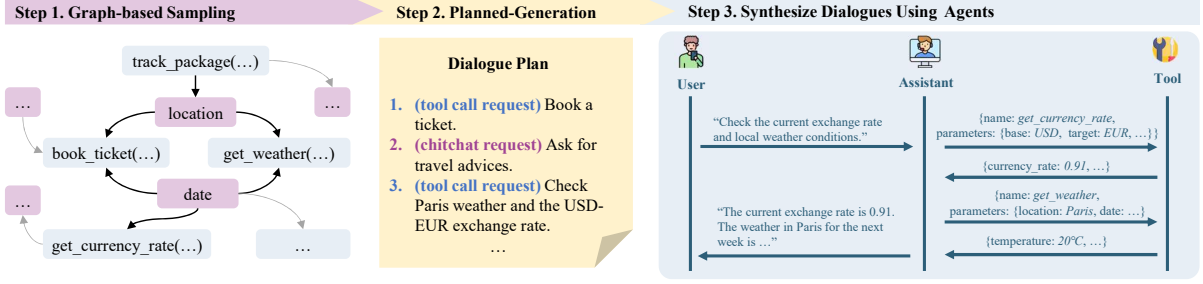


Figure 1: The pipeline of dialogue synthesis. The left side shows the Tool Graph with blue boxes representing tools and purple boxes representing parameters or return values. In the middle is the dialogue synthesis plan generated according to sampled tools. On the right is an example of data synthesis by the *User*, *Assistant*, and *Tool* agents.

significance of tool selection (Qin et al., 2023; Patil et al., 2023), relying solely on random sampling, the chosen tools play a crucial role in shaping the quality of the synthesized dialogues. In real-world scenarios, user requirements often necessitate the combined use of multiple tools to achieve a solution. To synthesize more complex user needs, it is vital to ensure that the selected tools can work together. To address this, we propose a Graph Sampling strategy to identify relevant and compatible tool combinations.

3.1.1 Tool Graph Construction

We first construct a graph, $G = (V, E)$, where node $v_i \in V$ represents the tool i , and the edge $e_{i,j} \in E$ represents whether tool i and tool j are related. The left side of Figure 1 shows an example of the tool graph. We consider tools with similar parameters or return values to be related to each other:

P-P Similarity: When two tools share similar parameters, there is a high probability that the tools are related. For instance, based on "location" and "destination", two semantically similar parameters, we can identify tools like "get_weather" and "book_flight", which are frequently used together in travel-related contexts.

P-R Similarity: If the return value of one tool is similar to the input parameters of another, there is also a high likelihood that the two tools are related. For example, the "check_calendar" tool typically returns the location of events, while the "navigate" tool requires a location as input. When a user requests to "navigate to the location of this afternoon's meeting," both tools would be called.

To derive similarity between parameters or return values, we first concatenate the name and description of a parameter or a return value using the template "{Name}: {Description}". For example, the parameter "Date" of one specific tool is

represented as "Date: Departure date, format as dd/mm/yyyy." Then, we encode these strings using Sentence-BERT (Reimers and Gurevych, 2019) to obtain the corresponding embeddings. We use p_k^i and r_l^i to denote the k -th parameters and l -th return values of the tool v_i , respectively. And we use \mathbf{p}_k^i and \mathbf{r}_l^i to represent the embeddings of p_k^i and r_l^i , respectively. The similarity between v_i 's parameter p_k^i and v_j 's parameter p_l^j is defined as:

$$\cos(\mathbf{p}_k^i, \mathbf{p}_l^j) = \frac{\mathbf{p}_k^i \cdot \mathbf{p}_l^j}{\|\mathbf{p}_k^i\| \|\mathbf{p}_l^j\|} \quad (1)$$

Similarly, the similarity between v_i 's return value r_k^i and v_j 's parameter p_l^j is defined as $\cos(\mathbf{r}_k^i, \mathbf{p}_l^j)$. If the similarity is greater than a predefined threshold τ , we consider the two tools to be correlated. We set τ to be 0.82 according to our preliminary study. This means that when the similarity between any pair of parameters from two tools exceeds τ , or when the similarity between a return value and an input parameter of two tools exceeds τ , we assign an edge between the two tools:

$$e_{i,j} = \begin{cases} 1, & \exists p_k^i \subseteq v_i, p_l^j \subseteq v_j : \cos(\mathbf{p}_k^i, \mathbf{p}_l^j) > \tau \\ & \text{or} \\ & \exists r_k^i \subseteq v_i, p_l^j \subseteq v_j : \cos(\mathbf{r}_k^i, \mathbf{p}_l^j) > \tau \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

where $i, j = 1 \dots N$ and $i \neq j$. We use \subseteq to represent a parameter or a return value is included in the tool.

3.1.2 Graph-based Sampling

With the constructed tool graph, we are able to sample a subset consisting of n tools that might be correlated. Generally, we randomly select a node as the starting point on the graph, and then perform a random walk along the edges of the graph. We

Algorithm 1 Graph-based Sampling

Require: $G = (V, E)$ with $|V| = N$, integer $n \leq N$ (desired sample size)

Ensure: Subset $\hat{V} \subseteq V$ with $|\hat{V}| = n$

```
1: // Randomly choose a node from V and add it to V-hat
2:  $\hat{V} \leftarrow \{\text{Uniform}(V)\}$ 
3: while  $|\hat{V}| < n$  do
4:    $v_i \leftarrow$  last element of  $\hat{V}$ 
5:   // Find all neighbors of  $v_i$ 
6:    $N(v_i) \leftarrow \{v_j \in V \mid e_{i,j} = 1, \forall e_{i,j} \in E\}$ 
7:   // Randomly choose a neighbor of  $v_i$ 
8:    $v_j \leftarrow \text{Uniform}(N(v_i))$ 
9:   if  $v_j \notin \hat{V}$  then
10:    // Add  $v_j$  to  $\hat{V}$  if not already included
11:     $\hat{V} \leftarrow \hat{V} \cup \{v_j\}$ 
12:   end if
13: end while
14: return  $\hat{V}$ 
```

stop when the path length reaches n , and the nodes included in the path constitute the sampled subset of tools. Details are shown in Algorithm 1.

3.2 Dialogue Plan Generation

Coherent dialogues usually involve complex tool-calling scenarios, such as cases where the current tool calling relies on the return value of a previous one or where parameters for the current tool calling are already present in the dialogue history. Moreover, realistic dialogues between humans and AI assistants often don't always require tool use; they frequently involve non-tool-related exchanges, such as chitchat, interspersed with tool-based tasks. To enhance an LLM's performance in such realistic scenarios, training examples that reflect this balance are essential.

To address the need and enhance the coherence of synthesized dialogues, we propose a Planned-Generation strategy, which indicates planning before generating. We have the LLM first formulate a set of user requirements based on the tool subset to create a dialogue plan. These requirements can involve tool call requests – tasks that necessitate the use of these tools – or non-tool interactions, such as chitchat. The middle part of Figure 1 shows an example of a plan. At this stage, the LLM focuses solely on the logic, coherence, and natural flow of the requirements, without delving into the phrasing of interactions or other nuances. Compared to directly synthesized dialogues, those generated based on a dialogue plan show markedly improved coherence. We provide a detailed assessment of this coherence in the subsequent sections. Please refer to Table 14 for the plan synthetic prompt.

3.3 Multi-Agent Dialogue Synthesis

We set up three agents, *user*, *assistant*, and *tool*, with LLMs to collaboratively synthesize dialogues. The right side of Figure 1 illustrates the synthesis process for one dialogue turn.

The *user* agent is responsible for initiating requests based on the dialogue plan. It first checks whether the current request in the plan has been completed, which determines whether to continue with the current task or move on to the next one. This ensures that the dialogue stays aligned with the plan's flow and sequence.

The *assistant* agent evaluates the user's request to determine if a tool is required. In cases where no tool is needed, such as chitchat, the assistant agent responds directly. If a tool call is necessary, the assistant verifies whether all required parameters are present based on the tool documentation. If any parameters are missing, the assistant requests clarification from the user; otherwise, it generates the tool call statement.

The *tool* agent simulates the return values of the requested tool based on the tool documentation and the assistant's call statement.

The interaction among the three agents continues for each dialogue turn until all the requests in plan are addressed or the preset turn limit is reached. Afterward, all dialogue turns are collected and a rule-based data filtering module is applied to remove low-quality data (Liu et al., 2024b,a). The filtering rules primarily check the format of tool call statements, as well as other issues such as incomplete dialogues or missing tool call turns.

3.4 Implementation Details

In this work, for tool selection, we directly utilize the tools from ToolBench (Qin et al., 2023) (including over 16,000 RESTful APIs) as our available tools. To standardize tool descriptions, we follow the setting of OpenAI Function Calling¹ and prompt an LLM (Llama-3.1-8B) to convert all these tools into JSON format. For cases where information is incomplete, such as missing parameter descriptions, we rely on the LLM to infer and fill in the missing details during the conversion process. A demonstration tool is presented in Figure 2.

We use GPT-4 (OpenAI et al., 2024) for all generative tasks, including dialogue plan generation and agent simulation, unless otherwise specified.

¹<https://platform.openai.com/docs/guides/function-calling>

Setting		Statistics		
Graph	Plan	# tokens	# call	# call turn
✓	✓	8,054,298	21,069	17,112
✓	✗	8,145,545	25,158	21,504
✗	✓	7,956,087	18,117	15,931
✗	✗	8,069,304	23,301	20,804

Table 1: Basic information about the dialogue datasets synthesized by TOOLFLOW and its ablation settings. # token, # call, and # call turn represent the number of tokens, tool calls, and turns containing tool calls in the dataset, respectively.

Multiple versions of GPT-4 are randomly selected for each dialogue, potentially enhancing diversity.

4 Data Quality Assessment

4.1 Basic Data Information

In this section, we evaluate the quality of the synthetic data. To assess the effectiveness of the Graph-based Sampling (referred to as Graph) and the Planned-Generation strategy (referred to as Plan), we synthesized three additional sets of comparative data under different conditions: removing Graph, removing Plan, and removing both Graph and Plan. Each dataset contains 8,000 dialogues. Table 1 presents the total number of tokens, the number of tool calls, and the number of dialogue turns containing tool calls for these datasets.

As shown in Table 1, the total number of tokens generated by the different strategies is similar, at approximately 8 million. Dialogues synthesized using the Planned-Generation strategy include more non-tool interactions, resulting in a lower proportion of tool call requests. In contrast, the Graph-based Sampling strategy increases the number of tool calls. This can be attributed to the connections among tools, where relevant information for subsequent tool calls is generally contained in the dialogue history, thereby reducing the need for additional turns to ask for missing information.

4.2 Quality Evaluation

To further assess the quality of the synthetic dialogues, we implemented both an automatic evaluation and a model-based evaluation.

The automatic evaluation primarily assesses the coherence and diversity of the dialogues. Following Dziri et al. (2019), we assess the coherence of the dialogue as a Natural Language Inference (NLI) task. We treat two consecutive turns in the dialogue as the premise and hypothesis, respectively, and calculate the ratio of entailment relation (*EnR*) as well

Automatic Evaluation					
Setting		Coherence		Diversity	
Graph	Plan	SS	EnR	H	D-3
✓	✓	63.36	47.3	10.36	0.4865
✓	✗	62.03	32.1	10.14	0.4364
✗	✓	61.72	48.1	9.82	0.3393
✗	✗	58.65	35.4	9.75	0.3078
GPT-4 Evaluation					
Graph	Plan	NAT	COH	HELP	ACC
✓	✓	3.72	3.91	4.71	4.92
✓	✗	3.00	3.72	4.51	4.84
✗	✓	3.51	3.88	4.39	4.87
✗	✗	2.93	3.66	4.18	4.90

Table 2: Results of automatic evaluation and GPT-4 evaluation on the data synthesized by TOOLFLOW and its ablation settings.

as the semantic similarity (*SS*) between them. A higher *EnR* or *SS* between turns indicates that the dialogue is more coherent.

Regarding diversity, we calculate the text’s Shannon entropy (*H*) based on the word frequency (Shannon, 1948). We also compute the Distinct-N Score (Li et al., 2016) for the dataset, with $N = 3$ (*D-3*). Higher entropy or Distinct-N Score indicates that the dataset contains more information and has greater diversity.

In addition, we randomly sampled 200 dialogues in each dataset for the model-based evaluation. We used GPT-4 (OpenAI et al., 2024) to carefully evaluate each dialogue based on four dimensions: naturalness (*NAT*), coherence (*COH*), helpfulness (*HELP*), and accuracy (*ACC*). The prompt for GPT-4 evaluation is shown in Table 15.

The evaluation results are shown in Table 2. There are two key observations:

- *H* and *D-3* Score demonstrate that Graph Sampling enhances the diversity of the data.
- Both evaluations (*SS*, *EnR*, and *COH*) show that Planned-generation improves the coherence of the dialogue.

For more detailed settings, analysis and explanations, please refer to Appendix A.1.

4.3 Comparison with Natural Dialogue Dataset

To better understand how our synthetic dialogues compare with human-created ones, we conducted a comparative study with an established dataset. We chose the MultiWOZ dataset (Budzianowski et al., 2018) as a natural dialogue dataset for comparison. MultiWOZ (Multi-Domain Wizard-of-Oz) is a well-known task-oriented dialogue dataset, and

Dataset	NAT	COH	HELP	ACC
TOOLFLOW	3.72	3.91	4.71	4.92
MultiWOZ	3.98	4.03	4.41	4.95

Table 3: Evaluation Scores Comparison between TOOLFLOW and MultiWOZ.

we believe comparisons with this dataset would be convincing. We repeated the GPT-4 evaluation experiment by first randomly sampling 200 dialogues from MultiWOZ. Then, using the same scoring prompt, we had GPT-4 evaluate these dialogues across the four dimensions. The results are shown in Table 3.

MultiWOZ scores slightly higher on naturalness, coherence, and accuracy compared to the ToolFlow dataset, though the differences are minimal (average score differences between 0.1-0.2). Regarding helpfulness scores, ToolFlow outperformed MultiWOZ by 0.3 points. These results suggest that our synthetic dialogues in TOOLFLOW achieve comparable quality to human-created dialogues in MultiWOZ, with particularly strong performance in task-oriented aspects such as helpfulness.

5 Experiments

5.1 Settings

5.1.1 Datasets

We conducted experiments on the following three tool-calling datasets to validate the tool call capability of the model trained with TOOLFLOW.

- **BFCL-v2** (Patil et al., 2023) primarily consists of Python-style tool call data, divided into four categories *Simple*, *Multiple*, *Parallel*, and *Parallel Multiple*. Version 2 adds more data from dynamic, real world scenarios. We selected the categories that can be evaluated with the Abstract Syntax Tree (AST), which are statistically stable and easy to evaluate. The accuracy is reported.
- **API-Bank** (Li et al., 2023) is a dialogue-style tool call dataset, including two settings: *Call* and *Retrieve + Call*. The model is required to call predefined local Python tools based on user requirements in the dialogue. Accuracy is measured by evaluating whether the tool return values match the ground truth.
- **ToolAlpaca** (Tang et al., 2023) establishes a multi-agent simulator. It utilizes GPT-4 (OpenAI et al., 2024) to simulate the return values

of tools. The model can make modifications and re-call the tool based on the return values (e.g., error messages). Finally, GPT-4 evaluates the accuracy of *Process* and *Response*.

Additionally, to examine changes in general performance, we evaluated the model’s reasoning and conversational abilities using MMLU (Hendrycks et al., 2021), BBH (Suzgun et al., 2022), and MT-Bench (Zheng et al., 2024).

5.1.2 Models

In our main experiments, we use LLaMA-3.1-8B-Instruct (Dubey et al., 2024) as base model to examine the effectiveness of the synthetic dialogues generated with TOOLFLOW. For simplicity, we use TOOLFLOW to refer to the fine-tuned model throughout the remainder of this paper. The models we compared include GPT-3.5, GPT-4, GPT-4o (OpenAI et al., 2024), Claude (Bai et al., 2022), LLaMA-3.1 (Dubey et al., 2024), etc, as well as baselines from the paper of the evaluated datasets, such as Lynx-7B (Li et al., 2023) and ToolAlpaca-7B (Tang et al., 2023), etc. For specific checkpoint information, please refer to the experimental result tables.

5.2 Results

5.2.1 TOOLFLOW achieves tool-calling capability comparable to GPT-4o.

We evaluated TOOLFLOW’s tool-calling ability on the BFCL. This dataset contains questions from four categories. In the **Simple** category, each question contains one tool, which the LLM must correctly call based on requirements. The **Multiple** question includes 2-4 tools, requiring the model to choose and call the most suitable one. In the **Parallel** category, several tools should be called in one turn. **Multiple Parallel** adds distracting candidate tools to the Parallel setup.

The results are shown in Tabel 4. Overall, TOOLFLOW achieves performance comparable to GPT-4o. On the Non-Live subset, TOOLFLOW outperformed GPT-4 and GPT-4o, but was slightly weaker than Claude-3.5-Sonnet. On the Live subset, TOOLFLOW still lags behind these leading closed-source LLMs. This is because the Live subset added more user-contributed test cases from the real world, thus making it more challenging for the model. We attribute this gap primarily to differences in model size, given that TOOLFLOW only has 8B parameters. The ablation experiment shows that the model trained on data synthesized

Baselines	Non-Live				Live				Overall				
	Simple	Multiple	Parallel	Parallel Multiple	Simple	Multiple	Parallel	Parallel Multiple	Non-Live	Live	Overall		
Claude-3.5-Sonnet	88.55	95.00	91.50	92.50	86.82	80.06	81.25	45.83	91.22	80.75	85.98		
GPT-4-turbo-0409	87.45	96.50	91.00	89.00	87.98	84.14	100.00	79.17	90.28	85.02	87.65		
GPT-4o-0513	80.55	91.00	90.00	83.00	81.78	77.24	87.50	75.00	85.02	78.20	81.61		
LLaMA-3.1 8B	90.36	89.50	73.50	73.50	74.03	73.31	56.25	54.17	83.40	72.88	78.14		
Graph Plan													
	✓	✓	92.18	92.50	<u>90.00</u>	<u>85.00</u>	<u>73.64</u>	<u>75.22</u>	<u>75.00</u>	<u>70.83</u>	<u>90.30</u>	<u>74.83</u>	<u>82.57</u>
Ours	✓	✗	90.73	<u>93.00</u>	89.50	84.50	71.71	74.35	<u>75.00</u>	66.67	89.60	73.71	81.66
	✗	✓	91.45	<u>93.00</u>	88.00	84.00	70.93	74.73	68.75	66.67	89.50	73.78	81.64
	✗	✗	92.36	91.50	86.50	82.50	72.48	73.58	<u>75.00</u>	62.50	89.00	73.18	81.09

Table 4: Results on BFCL-v2 leaderboard (updated on 08/16/2024). "Non-Live" and "Live" indicate the results on v1 and v2 subsets respectively. Table values are shown as a percentage. The best results in each category are marked in **bold**. The best results from our model are underlined.

by strategies including both Graph-based Sampling and Generated Plan performs the best. This is especially evident in Parallel Multiple type questions.

While different models or training strategies exhibit some variance in certain categories, the differences in average performance are not significant. Therefore, we conducted further comparative experiments on additional tool call datasets.

5.2.2 TOOLFLOW achieves SOTA on dialogue data.

BFCL tests the model’s tool calling capability in the form of Q&A. However, we believe that a conversational format is closer to real-world application scenarios. On the other hand, our synthesized training data is also in the form of dialogue. Therefore, in the BFCL test, the advantage of TOOLFLOW cannot be fully demonstrated.

API-Bank is a dialogue dataset. During evaluation, the model needs to make tool call requests after receiving user demands and provides a response based on the tool’s return value. This process may occur multiple times within a single dialogue. It includes two test settings: *Call* and *Retrieve + Call*. In the *Call* setting, the assistant selects tools from a candidate tool set to fulfill user requests. In the *Retrieve + Call* setting, the assistant only has access to a search tool. The assistant needs to search for the relevant tools first, and then call them.

From the results in Table 5, TOOLFLOW achieves state-of-the-art average accuracy. Under the *Call* setting, TOOLFLOW outperforms all baselines. In the *retrieve + call* setting, TOOLFLOW is inferior to GPT-3.5-turbo but superior to other baselines, including GPT-4. The ablation experi-

Baselines	L1	L2	Avg.		
GPT-4-turbo-0409	72.43	39.26	55.85		
GPT-4o-0513	76.19	42.96	59.58		
GPT-3.5-turbo-0125	70.43	52.59	61.51		
Lynx-7B [†]	49.87	30.37	40.12		
Llama-3.1-8B-Instruct	71.18	37.04	54.11		
Graph Plan					
	✓	✓	<u>77.52</u>	<u>46.68</u>	62.10
Ours	✓	✗	76.26	42.23	59.25
	✗	✓	79.30	38.53	58.92
	✗	✗	76.76	39.27	58.02

Table 5: Results on API-Bank Dataset. L1 and L2 refer to the setting of *Call* and *Retrieve + Call*, respectively. Table values are shown as a percentage. [†] indicates results derived from the original paper.

ments show that Graph-based Sampling strategy can improve the model’s accuracy in tool call under this setting. This is because tools obtained through Graph sampling often have sequential correlations. As a result, the training data includes more examples of sequentially calling tools, aligning better with the requirements of the *retrieve + call* setting.

5.2.3 TOOLFLOW can correct mistakes based on error messages.

Correcting errors is a key capability of LLM tool calls (Wang et al., 2024). We conducted the tests in the *Simulated* setting of ToolAlpaca dataset. This dataset established a simulation environment that utilizes GPT-4 to mimic the return values of tools, including the error messages when calls fail. The model is allowed to self-correct based on these error messages and then retry the call. We assess the tool call and correction capabilities of TOOLFLOW

Baselines		Proc.	Resp.	Overall
GPT-3.5 [†]		77	85	75
ToolAlpaca-13B [†]		63	69	60
ToolAlpaca-7B [†]		70	73	70
LLaMA-3.1 8B		74	80	74
	Graph Plan			
	✓ ✓	85	88	84
Ours	✓ ✗	80	85	80
	✗ ✓	81	83	79
	✗ ✗	78	83	77

Table 6: Results on ToolAlpaca Dataset. **Proc.** and **Resp.** stand for *Procedure* and *Response*, respectively. Table values are shown as a percentage. [†] indicates results derived from the original paper.

on this dataset.

The dataset evaluates the accuracy of the tool call *Procedure* and the model’s final *Response*. The procedure is considered accurate when the model’s call matches the ground truth. The response is considered accurate when the model’s response can satisfy the user’s instruction. If they are both accurate, the model’s *Overall* performance is considered accurate. This evaluation was conducted by GPT-4. We presented the results in Table 6.

TOOLFLOW’s *Procedure* accuracy reached 85%, surpassing GPT-3.5’s 77%. In *Procedure* evaluation, correcting errors is considered as redundant actions and therefore judged as incorrect. Hence, this accuracy implies that in most cases, the first tool call of TOOLFLOW is accurate. On the other hand, TOOLFLOW’s *Response* accuracy of 88% is higher than the *Procedure* accuracy of 85%, indicating that TOOLFLOW corrected errors in some test cases. This suggests that TOOLFLOW has the ability to self-correct based on error messages, even though error correction samples are not included in the training data.

5.2.4 TOOLFLOW’s General Ability Is NOT Compromised by Fine-tuning.

The fine-tuned model risks catastrophic forgetting, where the capability for tool call is enhanced, but other abilities decline. As an AI assistant, LLM’s reasoning and conversational abilities are equally important. Therefore, we tested the tuned model on the MMLU, BBH, and MTBench datasets to examine whether catastrophic forgetting issues have occurred. The results are shown in Table 7.

The test results on MMLU and BBH show that there is no significant difference in performance between the models before and after training. How-

Setting	MTBench			MMLU	BBH	
	Graph Plan	Turn 1	Turn 2 Avg.			
-	-	7.84	6.85	7.34	69.3%	63.2%
✓	✓	8.08	7.16	7.62	69.8%	63.3%
✓	✗	7.39	6.43	6.91	69.3%	63.5%
✗	✓	7.59	7.04	7.31	69.1%	63.2%
✗	✗	7.16	6.85	7.01	68.9%	63.5%

Table 7: Results of general abilities test on MTBench, MMLU, and BBH. The results in the first row are from the model before tuning. The best results are marked in **bold**. **Red** indicates a decrease in performance.

ever, on the MTBench dataset, models trained on data without Graph-based Sampling or Plan-Generation exhibited a decline in performance. Notably, in the evaluation of *Turn 2*, models trained on synthetic data using the Plan-Generation strategy exhibited a slight performance improvement. This improvement is due to Plan-Generation enhancing the naturalness and coherence of synthetic dialogues, thereby boosting the model’s conversational abilities.

6 Correlation Analysis

To further investigate the impact of the diversity and coherence of the dialogue data on model performance, we conducted additional correlation analysis. In the previous experiments, we synthesized a total of $8,000 \times 4 = 32,000$ dialogues by using TOOLFLOW and its ablation settings. We randomly sampled from these data 10 times, each time selecting 4,000 dialogues to form 10 new training sets. We calculated the diversity metrics $D-3$ and H , and the coherence metrics SS and EnR for each dataset. Then, we used this data to fine-tune the Llama3.1-8B-Instruct. We tested these ten fine-tuned models on the BFCL and MTBench. Finally, we calculated the Pearson correlation coefficient between the evaluation metrics of the training data and the model performance and reported it in Table 8.

The average results on the BFCL show that both diversity and coherence of the training data contribute a lot to enhancing the model’s tool-calling capabilities. MTBench results show a strong positive correlation between data coherence and the model’s conversational performance, consistent with our assumption. Notably, while we use entropy and Distinct-N scores to assess diversity, the inconsistent correlation between these metrics and model performance suggests they may reflect different dimensions of diversity. On the other hand,

Metrics	D-3	H	SS	EnR
BFCL				
Simple	0.174	-0.272	0.347	0.118
Parallel	0.250	0.315	-0.168	0.072
Multiple	0.336	-0.042	0.644	0.612
Parallel Multiple	-0.274	0.479	-0.310	-0.102
Avg.	0.250	0.273	0.255	0.378
MTBench				
Turn 1	-0.087	0.221	0.262	0.298
Turn 2	-0.185	0.045	0.708	0.415
Avg.	-0.179	0.146	0.650	0.454

Table 8: The Pearson correlation between evaluation metrics of data and model’s performance. P-values > 0.2 and < 0.2 are marked in green and red, respectively.

coherence does not appear to positively impact the parallel test sets in the BFCL, likely due to the nature of these tests involving multiple calls within a single turn. Nevertheless, while our TOOLFLOW has demonstrated the benefits of increasing diversity and coherence in Section 4.2, the correlation results in this section further validate their positive effects on overall performance.

7 Dataset Overlap Analysis

To ensure the reliability of the evaluation, we conducted an overlap analysis between training and test datasets. This examination helps verify the independence of these test data and prevents potential data leakage issues. We employed both N-gram-based and similarity-based methods to demonstrate that there is no significant data leakage in the TOOLFLOW dataset. We also included the well-known xLam agent training set (Zhang et al., 2024) as a control group for comparison.

N-gram-based method Following the approach used in LLaMA-2 (Touvron et al., 2023), we considered a token contaminated if it appeared in any token n-gram longer than 10 tokens in both the evaluation sample and the training set. A tool was classified as leaked if more than 10% of the tokens in its JSON string were contaminated.

Similarity-based method We defined a tool as leaked if the cosine similarity between the given tool and any tool in the evaluation dataset exceeded 0.9. We used the all-MiniLM-L12-v2 encoder from HuggingFace² to obtain representations for all tools.

We present the proportions of data leakage across different evaluation metrics in Table 9.

²<https://huggingface.co/>

These results suggest that there is no severe data leakage between TOOLFLOW as a training set and the test sets.

Training	ToolFlow		xLam	
	N-gram	Similarity	N-gram	Similarity
Test				
BFCL	0.239%	2.922%	0.656%	5.247%
APIBank	0.000%	1.887%	0.000%	3.774%
ToolAlpaca	0.000%	0.000%	0.000%	0.000%

Table 9: Overlap between Training and Test Sets

8 Conclusion

In this work, we propose Graph-based Sampling and Planned Generation strategies to enhance the diversity and coherence of synthetic data. Based on these two strategies, we introduce a pipeline called TOOLFLOW for synthesizing tool calling data and generate 8,000 training samples. Using this dataset, we conduct SFT on Llama3.1-8B-Instruct, resulting in improved tool calling capability of the model. Subsequently, we conduct correlation analysis to demonstrate the influence of data diversity and coherence on model performance. This provides a reference for the composition of training data for the tool-enhanced agent.

Limitations

We summarize the limitations in two points.

As described in Section 3.4, the seed data is a pre-collected tool set including 16,000 APIs. Although our TOOLFLOW can synthesize more diverse data, it is undeniable that the size and diversity of the tool set also affect the diversity of the data. However, how to enrich the seed data has not yet been studied in this work.

On the other hand, TOOLFLOW utilizes GPT-4 for data synthesis, and then uses this data to train a 8B-model. Therefore, it still falls under the paradigm of using strong models to train weak models. Whether the model can be improved by training on its own synthesized data is still unknown. We believe that this weak-to-strong setting is more challenging but also more meaningful.

Ethic Statement

In this research, GPT-4 was employed as an evaluator and generator in a manner consistent with ethical guidelines. Transparency about its usage, accountability for its outputs, and mitigation of

potential biases were prioritized. Data privacy and security were strictly maintained, and the AI's limitations were acknowledged, ensuring it supplemented rather than replaced human judgment. This approach aimed to enhance the research quality while upholding academic integrity and ethical standards.

Acknowledgements

This work was partially supported by Hong Kong RGC GRF No. 14206324, CUHK direct grant No. 4055209, and CUHK Knowledge Transfer Project Fund No. KPF23GWP20.

References

- Ibrahim Abdelaziz, Kinjal Basu, Mayank Agarwal, Sadhana Kumaravel, Matthew Stallone, Rameswar Panda, Yara Rizk, GP Bhargav, Maxwell Crouse, Chulaka Gunasekara, Shajith Iqbal, Sachin Joshi, Hima Karanam, Vineet Kumar, Asim Munawar, Sumit Neelam, Dinesh Raghu, Udit Sharma, Adriana Meza Soria, Dheeraj Sreedhar, Praveen Venkateswaran, Merve Unuvar, David Cox, Salim Roukos, Luis Lastras, and Pavan Kapanipathi. 2024. [Granite-function calling model: Introducing function calling abilities via multi-task learning of granular tasks](#). *Preprint*, arXiv:2407.00121.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, and Anna Chen. 2022. [Constitutional ai: Harmlessness from ai feedback](#). *Preprint*, arXiv:2212.08073.
- Kinjal Basu, Ibrahim Abdelaziz, Subhajt Chaudhury, Soham Dan, Maxwell Crouse, Asim Munawar, Vernon Austel, Sadhana Kumaravel, Vinod Muthusamy, Pavan Kapanipathi, and Luis Lastras. 2024. [API-BLEND: A comprehensive corpora for training and benchmarking API LLMs](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12859–12870, Bangkok, Thailand. Association for Computational Linguistics.
- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2018. [MultiWOZ - a large-scale multi-domain Wizard-of-Oz dataset for task-oriented dialogue modelling](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5016–5026, Brussels, Belgium. Association for Computational Linguistics.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, and Alan Schelten. 2024. [The llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.
- Nouha Dziri, Ehsan Kamalloo, Kory Mathewson, and Osmar Zaiane. 2019. [Evaluating coherence in dialogue systems using entailment](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3806–3812, Minneapolis, Minnesota. Association for Computational Linguistics.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). *Preprint*, arXiv:2009.03300.
- Cheng-Yu Hsieh, Si-An Chen, Chun-Liang Li, Yasuhisa Fujii, Alexander Ratner, Chen-Yu Lee, Ranjay Krishna, and Tomas Pfister. 2023. [Tool documentation enables zero-shot tool-usage with large language models](#). *Preprint*, arXiv:2308.00675.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. [A diversity-promoting objective function for neural conversation models](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 110–119, San Diego, California. Association for Computational Linguistics.
- Minghao Li, Yingxiu Zhao, Bowen Yu, Feifan Song, Hangyu Li, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. 2023. [API-bank: A comprehensive benchmark for tool-augmented LLMs](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 3102–3116, Singapore. Association for Computational Linguistics.
- Weiwen Liu, Xu Huang, Xingshan Zeng, Xinlong Hao, Shuai Yu, Dexun Li, Shuai Wang, Weinan Gan, Zhengying Liu, Yuanqing Yu, Zezhong Wang, et al. 2024a. [Toolace: Winning the points of llm function calling](#). *Preprint*, arXiv:2409.00920.
- Zuxin Liu, Thai Hoang, Jianguo Zhang, Ming Zhu, Tian Lan, Shirley Kokane, Juntao Tan, Weiran Yao, Zhiwei Liu, Yihao Feng, Rithesh Murthy, Liangwei Yang, Silvio Savarese, Juan Carlos Niebles, Huan Wang, Shelby Heinecke, and Caiming Xiong. 2024b. [Apigen: Automated pipeline for generating verifiable and diverse function-calling datasets](#). *Preprint*, arXiv:2406.18518.
- Grégoire Mialon, Roberto Dessì, Maria Lomeli, Christoforos Nalmpantis, Ram Pasunuru, Roberta Raileanu, Baptiste Rozière, Timo Schick, Jane Dwivedi-Yu, Asli Celikyilmaz, Edouard Grave, Yann LeCun, and Thomas Scialom. 2023. [Augmented language models: a survey](#). *Preprint*, arXiv:2302.07842.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, et al. 2024. [Gpt-4 technical report](#). *Preprint*, arXiv:2303.08774.

- Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. 2023. [Gorilla: Large language model connected with massive apis](#). *Preprint*, arXiv:2305.15334.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2023. [Toolllm: Facilitating large language models to master 16000+ real-world apis](#). *Preprint*, arXiv:2307.16789.
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-bert: Sentence embeddings using siamese bert-networks](#). *Preprint*, arXiv:1908.10084.
- Jingqing Ruan, Yihong Chen, Bin Zhang, Zhiwei Xu, Tianpeng Bao, Guoqing Du, Shiwei Shi, Hangyu Mao, Ziyue Li, Xingyu Zeng, and Rui Zhao. 2023. [Tptu: Large language model-based ai agents for task planning and tool usage](#). *Preprint*, arXiv:2308.03427.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. [Toolformer: Language models can teach themselves to use tools](#). *Preprint*, arXiv:2302.04761.
- C. E. Shannon. 1948. [A mathematical theory of communication](#). *The Bell System Technical Journal*, 27(3):379–423.
- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, , and Jason Wei. 2022. Challenging big-bench tasks and whether chain-of-thought can solve them. *arXiv preprint arXiv:2210.09261*.
- Qiaoyu Tang, Ziliang Deng, Hongyu Lin, Xianpei Han, Qiao Liang, Boxi Cao, and Le Sun. 2023. [Toolalpaca: Generalized tool learning for language models with 3000 simulated cases](#). *Preprint*, arXiv:2306.05301.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, and Soumya Batra. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). *Preprint*, arXiv:2307.09288.
- Boshi Wang, Hao Fang, Jason Eisner, Benjamin Van Durme, and Yu Su. 2024. [Llms in the imaginarium: Tool learning through simulated trial and error](#). *Preprint*, arXiv:2403.04746.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. [Self-instruct: Aligning language models with self-generated instructions](#). *Preprint*, arXiv:2212.10560.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. [Wizardlm: Empowering large language models to follow complex instructions](#). *Preprint*, arXiv:2304.12244.
- Rui Yang, Lin Song, Yanwei Li, Sijie Zhao, Yixiao Ge, Xiu Li, and Ying Shan. 2023. [Gpt4tools: Teaching large language model to use tools via self-instruction](#). *Preprint*, arXiv:2305.18752.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. [React: Synergizing reasoning and acting in language models](#). *Preprint*, arXiv:2210.03629.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok, Zhenguang Li, Adrian Weller, and Weiyang Liu. 2024. [Metamath: Bootstrap your own mathematical questions for large language models](#). *Preprint*, arXiv:2309.12284.
- Jianguo Zhang, Tian Lan, Ming Zhu, Zuxin Liu, Thai Hoang, Shirley Kokane, Weiran Yao, Juntao Tan, Akshara Prabhakar, Haolin Chen, Zhiwei Liu, Yihao Feng, Tulika Awalgaonkar, Rithesh Murthy, Eric Hu, Zeyuan Chen, Ran Xu, Juan Carlos Niebles, Shelby Heinecke, Huan Wang, Silvio Savarese, and Caiming Xiong. 2024. [xlam: A family of large action models to empower ai agent systems](#). *Preprint*, arXiv:2409.03215.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2024. Judging llm-as-a-judge with mt-bench and chatbot arena. In *Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS '23*, Red Hook, NY, USA. Curran Associates Inc.

A Appendix

A.1 Details of Data Quality Assessment

Following Dziri et al. (2019), we converted the dialogue data into a Natural Language Inference (NLI) format. In this format, the request and response from the previous dialogue round serve as the **premise**, and the current round’s request serves as the **hypothesis**. We then use a trained BERT (Reimers and Gurevych, 2019) to predict the relationship between the two, and we calculate the ratio of entailment predictions (EnR). A higher proportion indicates greater coherence between consecutive dialogue rounds. Additionally, we measure the semantic similarity (SS) between the premise and hypothesis. We extract sentence representations using BERT and compute their cosine similarity. A higher similarity score indicates a more coherent dialogue.

Regarding diversity, we calculate the text’s Shannon entropy (H) based on the word frequency. We also compute the Distinct-N Score (Li et al., 2016) for the dataset, with $N = 3$ ($D-3$). Higher entropy or Distinct-N Score indicates that the dataset contains more information and has greater diversity. In addition, we sampled 200 dialogues in each set. We used GPT-4 to carefully evaluate each dialogue based on four dimensions: naturalness (NAT), coherence (COH), helpfulness ($HELP$), and accuracy (ACC). The prompt for GPT-4 evaluation is shown in Table 15.

Automatic evaluation indicates that the Planned-Generation strategy enhances conversation coherence. Both coherence metrics, SS and EnR, reflect this improvement. Intuitively, the plan is carefully designed by the model in advance, leading to more coherent dialogue. On the other hand, the Graph Sampling strategy can increase data diversity. This is because the strategy samples tools with strong associations, and the combination of these tools enhances data diversity.

GPT-4’s evaluation indicates that the Planned-Generation strategy enhances the naturalness of dialogue. This metric assesses whether a dialogue could realistically occur in the real world. In data synthesized without a plan, most user requests are tool calls with little chitchat, which is uncommon in real-world scenarios, resulting in lower naturalness. GPT-4’s coherence evaluation closely aligns with automatic assessments. In terms of helpfulness, both Graph Sampling and Planned-Generation strategies show improvement. The low

scores in Helpfulness are mainly due to the assistant frequently asking follow-up questions about parameters, which consumes dialogue turns. These strategies help reduce such behavior. The accuracy of synthesized data is high across all four settings, likely due to pre-filtering with quality control tools.

A.2 Reliability Analysis of GPT-4 as an Evaluator

To validate the reliability of GPT-4’s evaluations, we conducted additional human evaluations. Specifically, we sampled 50 examples from the data previously evaluated by GPT-4 for human assessment. We recruited four Computer Science PhD students to rate these 50 samples using the same criteria as GPT-4: naturalness, coherence, helpfulness, and accuracy. We collected the human evaluation results and calculated both the Cohen’s Kappa interrater agreement scores among evaluators and the Pearson correlation coefficients between human and GPT-4 evaluations.

Metric	NAT	COH	HELP	ACC
Cohen’s Kappa	0.63	0.71	0.92	0.95
Pearson Correlation	0.61	0.76	0.84	0.86

Table 10: Correlation between Human and GPT-4 Ratings

The results in Table 10 demonstrate high interrater agreement and strong correlations with GPT-4’s evaluations, particularly for helpfulness and accuracy metrics. While the agreement and correlations for naturalness and coherence were relatively lower compared to the other two metrics, they still maintained a minimum correlation of 0.61, indicating a strong positive correlation between human and GPT-4 evaluations. We believe these results substantiate the reliability of GPT-4’s evaluations.

A.3 GPT-4 v.s. Open Source LLM for Data Synthesis

To evaluate the reliance of synthesis on closed-source, high-performance LLMs (e.g., GPT-4), we conducted comparative experiments using open-source alternatives. Following the same pipeline but replacing GPT-4 with LLaMA-3.1-8B-instruct, we conducted preliminary synthesis experiments. While still in the initial validation phase, we have synthesized 4,235 dialogue instances. We replicated experiments using this data to fine-tune LLaMA-3.1-8B-instruct, with BFCL test results shown in Table 11. Here, "N/A" represents the

Data Source	Simple	Multiple	Parallel	Parallel multiple	Avg.
N/A	90.36	89.50	73.50	73.50	83.40
LLaMA-3.1-8B-instruct	89.88	90.10	85.55	80.15	86.42
GPT-4	91.23	91.85	87.10	84.45	88.66

Table 11: Performance comparison of models trained on data synthesized by GPT-4 and LLaMA.

Models	Simple	Multiple	Parallel	Parallel Multiple	Avg.
Mistral-7B-Instruct-v0.1	61.27%	54.00%	50.50%	47.50%	53.32%
+ TOOLFLOW	86.91%	81.00%	89.00%	80.50%	84.35%
Qwen2-7B	76.73%	63.50%	82.00%	55.50%	69.43%
+ TOOLFLOW	90.18%	81.50%	89.00%	78.50%	84.80%

Table 12: The results of Mistral 7B and Qwen2 7B on BFCL. '+TOOLFLOW' represents the results of the model after training with our data.

Models	L1	L2	Avg.
Mistral-7B-Instruct-v0.1	59.14	32.48	45.81
+ TOOLFLOW	68.42	45.19	56.81
Qwen2-7B	58.63	25.19	41.91
+ TOOLFLOW	64.97	36.93	50.95

Table 13: The results of Mistral 7B and Qwen2 7B on the API-Bank dataset. '+TOOLFLOW' represents the results of the model after training with our data.

baseline performance of LLaMA-3.1-8B-instruct (consistent with results reported in the main text). Notably, we leveraged existing results from Section 6’s Correlation analysis, where we had downsampled to ten subsets of 4,000 dialogues each, comparable to our LLaMA-3.1-8B-instruct synthesized dataset size. The table reports the average results across these ten subsets. While data synthesized by LLaMA-3.1-8B-instruct indeed shows lower performance compared to GPT-4-synthesized data, the gap is not substantial. Moreover, compared to LLaMA-3.1-8B-instruct’s initial performance, training on its self-synthesized data demonstrates improvements in tool-calling capabilities.

A.4 Train with other base model with TOOLFLOW

Table 12 and 13 display the results of two other base models fine-tuned with data generated with TOOLFLOW. Results demonstrate that TOOLFLOW enhances tool-calling performance across multiple LLMs, validating the generalizability of our synthetic data approach.

A.5 Prompts and Demonstrations

See Table 14 to 18 for details.

Prompt for Plan Generation

You are a conversation planning assistant. Given an available tools list and the target turn, you need to devise a dialogue flow between the user and the assistant centered around the use of these tools.

Here are definitions of the tools:

[tool_defs]

The number of target turns: [tgt_turn_num]

The planned conversation flow should adhere to the following requirements:

1. Only output the user's general requests for each interaction, without specifying the names of tools to be called, allowing users the freedom to expand upon their requests.
2. Each request must be categorized as one of the following types: a tool call request or chitchat.
3. Not every user request can be a tool call request; Chitchat can be appropriately included as a transition between conversations. But be careful not to include too much chitchat.
4. There must be a connection between consecutive requests, ensuring a smooth progression of the conversation. For instance, initiating with a request to book a meeting room, followed by request to send meeting invitation emails.
5. Please try to design more complex requests, which either calls multiple tools or requires multiple calls to resolve.
6. The outputted conversation flow must align with the anticipated number of turns for the interaction.

Here is an example:

The number of target turns: 5

The planned conversation flow:

1. Tool Call Request: The user expresses the need to reserve a meeting room.
2. Tool Call Request: The user requests to modify the reserved meeting room.
3. Chitchat: The user engages in casual conversation with the assistant about the topic of having too many meetings.
4. Tool Call Request: The user wishes to book multiple meetings, add corresponding schedules, and then send email reminders to the meeting participants.

Please directly output your planned dialogue flow without any additional analysis or explanation:

1. ...

2. ...

...

Table 14: Prompt for plan generation

Prompt for Data Evaluation

You are asked to evaluate some synthetic dialogue data. These synthetic dialogue occur between the user, the AI assistant, and the tool. Please evaluate the data based on the following criteria, assigning a score from 1 to 5 for each category. Use the detailed descriptions below to guide your assessment:

- Naturalness (1-5 points): Only evaluate whether the user's request and response is natural and realistic. Focus more on the natural flow of the conversation and less on the choice of words. For example, pay attention to whether users will ask similar questions in real scenarios. And assess whether user behavior is natural. For example, real users rarely ask similar questions consecutively or ask longer questions.
- Coherence (1-5 points): Evaluate the overall flow and logical connection between the turns in the conversation. Focus on checking whether the user's previous and subsequent rounds of requests are relevant.
- Helpfulness (1-5 points): Determine the effectiveness and value of the AI assistant's responses in addressing the user's needs.
- Accuracy (1-5 points): Check for the accuracy and consistency of the information provided. Everything returned by the tool is assumed to be accurate.

Please use a more ****CRITICAL**** and ****STRICT**** evaluation method. After scoring, please provide brief comments or feedback for each category to explain your ratings.

Please provide your evaluation in the following format:

Evaluation of Synthetic Dialogue Data

1. Naturalness: [Score] / 5
- Comments: [Brief comments or feedback]
2. Coherence: [Score] / 5
- Comments: [Brief comments or feedback]
3. Helpfulness: [Score] / 5
- Comments: [Brief comments or feedback]
4. Accuracy: [Score] / 5
- Comments: [Brief comments or feedback]

The dialogue you need to evaluate are as follows:

Table 15: Prompt for GPT-4 Data Evaluation

Prompt for User Agent

Please continue the next turn of dialogue based on the given tool definitions, the history of the dialogue, and the dialogue flow plan.

The requirements are as follows:

1. Only continue one turn, and the turn to be continued must be a user turn.
2. If the last turn of the history is the assistant asking for clarification, the user turn to be continued must provide the necessary parameters for the tool call as much as possible.
3. If the last turn of the history is the assistant fulfilling the requirement and returning the result, the user should make a new requirement. The new requirement should strictly adhere to the planned dialogue flow.
4. Please use [style] language style for the user turn to be continued. You can use references, transitions, and other methods to make the conversation more natural.
5. The requirement for the user turn to be continued must be clear and a problem that the assistant can immediately solve using the tool. The necessary parameters for the assistant to ask for clarification should also be provided immediately and not later.
6. Before continuing the writing, please clarify which step of the planned dialogue flow the current conversation is in, and output it.

Here is an example for reference:

```
<Start of Tool Definition>
{"name": "book_flight", "description": "Flight booking tools", "arguments": {"type": "object", "properties":
{"from_city_name": {"description": "departure city", "type": "string"}, "to_city_name": {"description": "arrival city",
"type": "string"}, "depart_date": {"description": "Departure date, in the format of YYYY-MM-dd, for example: 2024-02-24",
"type": "string"}, "flightNo": {"description": "flight number", "type": "string"}, "carbinClass": {"description": "Class of
service, such as Economy Class Y", "type": "string"}}, "required": ["from_city_name", "to_city_name", "depart_date"]}}}
<End of Tool Definition>
<Start of Dialogue Flow Plan>
1. Tool call requirement: The user proposes a request to book a flight ticket.
2. Tool call requirement: The user adds a request to book a flight ticket.
3. Chitchat: The user chats with the assistant about business trip issues.
<End of Dialogue Flow Plan>
<Start of Dialogue History>
<system>The current time is 2023-5-24 08:00:00.</system>
<system>The user is currently in Shenzhen.</system>
<user>I need help booking a flight to Beijing on June 1st. Can you assist me?</user>
<assistant>Sure, I will help you book a flight. Please wait a moment. book_flight{"from_city_name": "Shenzhen",
"to_city_name": "Beijing", "depart_date": "2023-06-01"}</assistant>
<tool>{"flight_message": "Flight booking successful. The details are as follows: Flight number YX002, flight time 2023-06-
01 11:00:00 -- 15:00:00, from Shenzhen Bao'an Airport to Beijing Capital International Airport."}</tool>
<assistant>I have booked a flight for you on June 1st, 2023, from 11:00 AM to 3:00 PM, from Shenzhen Bao'an Airport to
Beijing Capital International Airport. The flight number is YX002.</assistant>
<End of Dialogue History>
```

The expected output in this example is as follows:

```
<Start of Completed Process>
1. Tool call requirement: User requests to book a flight ticket.
<End of Completed Process>
<Start of Current Process>
2. Tool call requirement: User needs to add a flight ticket booking requirement.
<End of Current Process>
<Start of Continuation>
<user>Please book another return ticket, coming back on June 5th.</user>
<End of Continuation>
```

Here is the dialogue you need to continue:

```
<Start of Tool Definition>
[tool_definition]
<End of Tool Definition>
<Start of Dialogue Flow Plan>
[conversation_plan]
<End of Dialogue Flow Plan>
<Start of Dialogue History>
[conversation_begin]
<End of Dialogue History>
```

Table 16: Prompt for User agent

Prompt for Assistant Agent

You are an experienced assistant who resolves user issues by calling tools when necessary, and responds appropriately when tool call is not required. When you receive the return value from the tool, reply to the user based on the return value.

There are three examples,

Example 1.

<user>I need help booking a flight to Beijing on June 1st. Can you assist me?</user>

<assistant>Sure, I will help you book a flight. Please wait a moment. book_flight|{"from_city_name": "Shenzhen", "to_city_name": "Beijing", "depart_date": "2023-06-01"}</assistant>

Example 2.

<tool>{"flight_message": "Flight booking successful. The details are as follows: Flight number YX002, flight time 2023-06-01 11:00:00 -- 15:00:00, from Shenzhen Bao'an Airport to Beijing Capital International Airport."}</tool>

<assistant>I have booked a flight for you on June 1st, 2023, from 11:00 AM to 3:00 PM, from Shenzhen Bao'an Airport to Beijing Capital International Airport. The flight number is YX002.</assistant>

Example 3.

<user> All this traveling recently is really taking a toll on my sleep. </user>

<assistant>This might be due to jet lag. It could be helpful if you rest a bit earlier tonight. Before going to bed, you might want to try some relaxing activities, like reading a physical book, listening to gentle music, or practicing meditation, to help ease the stress that has built up during the day. </assistant>

Here are the guidelines for tool calls:

1. Inquiry for missing required parameters from the user only when necessary. If the required parameters are either empty or fully provided, proceed with the tool call directly.',
 2. Do not fabricate information for historical data or parameters not specified by the system or the user. Avoid using any vague pronouns such as "my XX," "the user," etc.',
 3. Any information already provided by the user should be included as parameter values where applicable, even if they are optional parameters.
-

Table 17: Prompt for Assistant agent

Prompt for Tool Agent

Now, you are to act as a tool. Given the definition of the tool and the corresponding tool call statement, you need to simulate the tool's response based on your own knowledge, and then return the result to the user in a specific format.

Here is the definition of the tool:

[tool_defs]

The tool call statement you receive is as follows:

[tool_calls]

Please return the simulated tool-generated response in the following JSON format, where each list element represents the return of a single call. Ensure the entire result can be read by json.loads. Please don't output any additional analysis or explanation.

```
[
  {
    "name": ...,
    "results": ...
  },
  {
    "name": ...,
    "results": ...
  }
  ...
]
```

Table 18: Prompt for Tool agent

```

1 {
2   "type": "function",
3   "function": {
4     "name": "getcurrency",
5     "description": "Get the current exchange rate for a specific
6     currency pair",
7     "parameters": {
8       "type": "object",
9       "properties": {
10        "basecurrency": {
11          "type": "string",
12          "description": "The base currency code, e.g., USD"
13        },
14        "targetcurrency": {
15          "type": "string",
16          "description": "The target currency code, e.g., EUR"
17        }
18      },
19      "required": ["basecurrency", "targetcurrency"]
20    },
21    "results": {
22      "type": "object",
23      "properties": {
24        "exchangerate": {
25          "type": "number",
26          "description": "The current exchange rate from base
27          currency to target currency"
28        },
29        "last_updated": {
30          "type": "string",
31          "description": "The date and time when the exchange
32          rate was last updated"
33        }
34      }
35    }
36  }
37 }

```

Figure 2: Example tool in JSON format.