

# Enhancing Tool Retrieval with Iterative Feedback from Large Language Models

Qiancheng Xu, Yongqi Li<sup>†</sup>, Heming Xia, Wenjie Li

Department of Computing, The Hong Kong Polytechnic University, China

{qiancheng.xu, he-ming.xia}@connect.polyu.hk

liyongqi@gmail.com cswjli@comp.polyu.edu.hk

## Abstract

Tool learning aims to enhance and expand large language models’ (LLMs) capabilities with external tools, which has gained significant attention recently. Current methods have shown that LLMs can effectively handle a certain amount of tools through in-context learning or fine-tuning. However, in real-world scenarios, the number of tools is typically extensive and irregularly updated, emphasizing the necessity for a dedicated tool retrieval component. Tool retrieval is nontrivial due to the following challenges: 1) complex user instructions and tool descriptions; 2) misalignment between tool retrieval and tool usage models. To address the above issues, we propose to enhance tool retrieval with iterative feedback from the large language model. Specifically, we prompt the tool usage model, i.e., the LLM, to provide feedback for the tool retriever model in multi-round, which could progressively improve the tool retriever’s understanding of instructions and tools and reduce the gap between the two standalone components. We build a unified and comprehensive benchmark to evaluate tool retrieval models. The extensive experiments indicate that our proposed approach achieves advanced performance in both in-domain evaluation and out-of-domain evaluation<sup>1</sup>.

## 1 Introduction

Large language models (LLMs) have demonstrated remarkable success in language-related tasks and are considered a potential pathway to achieving artificial general intelligence (Zhao et al., 2023). However, despite their powerful capabilities, LLMs are still limited in many aspects, such as knowledge update and mathematical reasoning. A promising way to overcome these limitations is to empower LLMs with external tools, known as tool

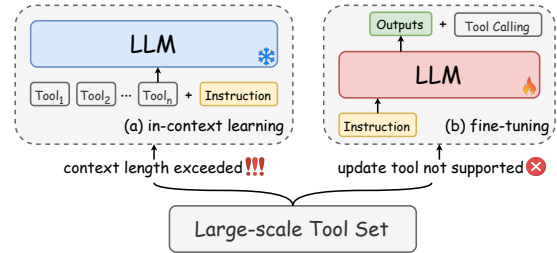


Figure 1: Illustration of two tool-learning approaches in LLMs: (a) in-context learning and (b) fine-tuning. The challenges posed by the extensive and frequently updated tools require the external tool retrieval component.

learning (Qin et al., 2023; Qu et al., 2024a). Tool learning not only enhances LLMs’ performance on existing tasks but also allows them to tackle tasks that were previously beyond their reach. Besides, the ability to use tools is a crucial hallmark on the path to advanced intelligence.

Existing tool learning methods have preliminarily demonstrated that LLMs could effectively utilize specific tools to complete corresponding tasks. They either leverage LLMs’ in-context learning ability to facilitate tool usage with tool descriptions (Shen et al., 2023) or fine-tune LLMs to integrate tool learning capabilities into parameters, e.g., Toolformer (Schick et al., 2023). However, as illustrated in Figure 1, existing methods still face significant challenges in real-world scenarios due to the following reasons. 1) The number of tools is usually vast, making it impossible for LLMs to handle them all with the limited input length of in-context learning. 2) Tools would frequently and irregularly update, rendering finetuning-based approaches costly and impractical. Therefore, a tool retrieval component, which aims to select appropriate tools from a large-scale tool set, is essential for LLMs.

Despite the practicality and necessity, tool retrieval has been inadequately studied. Some approaches have adopted traditional document re-

<sup>†</sup>Corresponding author.

<sup>1</sup>Code available at <https://github.com/travis-xu/TR-Feedback>.

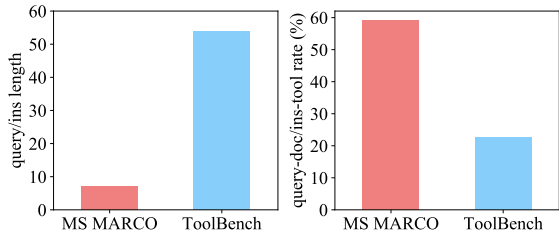


Figure 2: Comparison between the document retrieval and tool retrieval datasets. Tool retrieval presents more challenges due to the complex instructions (in the left figure) and the lower reputation rate (in the right figure).

retrieval methods to retrieve tools for LLMs (Li et al., 2023a; Qin et al., 2024). However, we argue that they overlook the unique challenges of tool retrieval for LLMs: 1) Complex user instructions and tool descriptions. As illustrated in Figure 2, compared with document retrieval, user instructions are usually ambiguous and complex, and the reputation rate between instructions and corresponding tool descriptions is much lower. Unfortunately, the retriever model is typically limited in its capacities because of the efficiency requirements, which makes tool retrieval more difficult and challenging. 2) Misalignment between tool retrieval and tool usage models. Previous approaches deploy the tool retriever separately from the downstream tool-usage model, which hinders the LLM from knowing which tools are really useful from the tool-usage perspective. Thus, it will result in a tool recognition gap between the tool retriever and tool usage model, degrading the tool-use performance further.

To address the above issues, we propose to enhance tool retrieval with iterative feedback. Our motivation is to utilize the LLM to enhance the comprehension ability of the tool retriever and bridge the gap between the two independent models. At each iteration, we conduct a feedback generation process by asking the LLM to provide feedback step-by-step, conditioned on the user instruction and retrieved tools from the retriever. The LLM will first comprehend the instruction and tool functionalities thoroughly, and then assess the effectiveness of those retrieved tools. According to the assessment, the LLM will refine the user instruction to improve the tool retrieval process. The refined instruction will substitute previous user instruction and be used to retrieve a new list of tools from the tool set. In the next iteration, the new candidate tool list will be fed into the LLM for a new round of LLMs’ feedback. During this it-

erative process, the tool retriever is expected to provide more appropriate tools for the tool-usage model. In this manner, the comprehension capability and tool preference of LLMs could be progressively incorporated into the retriever, and thus the tool retriever’s performance could be continuously enhanced. We build a comprehensive tool retrieval benchmark, named TR-bench. The benchmark takes into account real-world practices with updated tools, and therefore encompasses both in-domain and out-of-domain settings. The experimental results show our approach achieves the best performance among the current methods with both in-domain and out-of-domain settings.

The key contributions are summarized:

- We identify the importance of tool retrieval in tool learning and present the distinct challenges of tool retrieval.
- We propose to enhance tool retrieval with iterative feedback from the LLM. By leveraging iterative feedback, the tool retriever model gets continual improvements, ultimately reducing the misalignment between them.
- We build a comprehensive tool retrieval benchmark with in-domain and out-of-domain settings, which will also aid future tool retrieval research. The extensive experiments demonstrate superior performance of our approach.

## 2 Related Work

### 2.1 Tool Learning in LLMs

Tool learning aims to equip LLMs with external tools to enhance and expand their capabilities (Ruan et al., 2023; Wang et al., 2024b; Huang et al., 2024c). Generally, existing tool learning methods could be categorized into in-context learning and fine-tuning approaches. The former approach encourages LLMs to use tools with descriptions, documentation, or demonstrations (Yuan et al., 2024; Du et al., 2024), while the latter one trains the parameters of LLMs using specially created tool-use datasets (Hao et al., 2023; Tang et al., 2023; Gao et al., 2024). However, no matter whether the in-context learning or fine-tuning approach encounters severe challenges in real-world scenarios, where the candidate tools are extensive and frequently updated. Therefore, it is crucial to equip LLMs with a tool retrieval component to select appropriate tools from a large-scale tool

set. Recent works have proposed a stopgap measure through traditional document retrieval (Patil et al., 2023; Qin et al., 2024; Zheng et al., 2024), task decomposition (Anantha et al., 2023; Huang et al., 2024b) and graph-based methods (Qu et al., 2024b). In this work, we aim to develop a method specialized for enhancing the tool retriever.

## 2.2 Document Retrieval

Early popular document retrieval methods rely on sparse retrieval that calculates the relevance of documents to a query based on the frequency of query terms in each document, e.g., BM25 (Robertson and Zaragoza, 2009). With the development of language models (Devlin et al., 2019), the dense retrieval paradigm has gained considerable attention in the research community (Mitra and Craswell, 2017; Li et al., 2023b; Zhao et al., 2024; Li et al., 2024). By encoding queries and documents into high-dimensional vector representations and computing their relevance scores through inner product calculations, the paradigm can capture semantic relationships between queries and documents, thereby enhancing retrieval performance (Karpukhin et al., 2020). However, tool retrieval presents unique challenges, rendering traditional document retrieval methods suboptimal. We address these challenges by harnessing LLMs’ feedback to iteratively refine the tool retrieval process.

## 3 Preliminaries

### 3.1 Task Definition

Given a user’s instruction, tool retrieval aims to select a small number of tools, which could aid the LLM in answering the instruction, from a large-scale tool set. Formally, we define the user instruction as  $q$  and the tool set as  $D = \{d_1, d_2, \dots, d_N\}$ , where  $d_i$  represents the description of each tool and  $N$  is the total number of tools. The retriever model  $R$  needs to measure the relevance  $R(q, d_i)$  between the instruction  $q$  and each tool description  $d_i$ , and return  $K$  tools, denoted as  $D = \{d_1, d_2, \dots, d_K\}$ .

### 3.2 Dense Retriever

Dense retriever usually leverages the encoder-based LLM to encode the user instruction  $q$  and a tool description  $d$  into dense embeddings  $E(q)$  and  $E(d)$ , respectively. Then, it could measure the relevance between  $q$  and  $d$  by calculating the similarity

score between these two embeddings, denoted as  $R(q, d) = \text{sim}(E(q), E(d))$ .

Dense retriever is trained via the contrast learning objective, which is designed to minimize the distance between the instruction embedding and embeddings of positive tools (the instruction’s ground-truth tools) while maximizing the distance between the instruction embedding and embeddings of negative tools. The objective can be formulated as follows,

$$\mathcal{L} = -\frac{1}{B} \sum_{i=1}^B \log \frac{e^{R(q_i, d_i^+)}}{e^{R(q_i, d_i^+)} + \sum_j e^{R(q_i, d_{ij}^-)}}, \quad (1)$$

where  $B$  denotes the batch size,  $d_i^+$  denotes the positive tool, and  $d_{ij}^-$  represents the  $j$ -th negative tool to the instruction  $q_i$ .

However, due to the efficiency requirements, dense retrieval utilizes a dual-encoder architecture, which has limited ability to understand instructions. In this study, our goal is to improve the tool retrieval process with the feedback from the tool-usage model, i.e., the LLM.

## 4 Methodology

### 4.1 Overview

Recent studies have found that LLMs show a great capability in acting as a critic (Zheng et al., 2023) and could provide comprehensive feedback to improve performance across a range of tasks (Madaan et al., 2023; Asai et al., 2024). Inspired by those observations, we propose an innovative framework that leverages the LLM’s feedback to improve the tool retrieval process iteratively. Different from approaches which focus on feedback from execution results after tool execution step (Yao et al., 2023; Wang et al., 2024a), we obtain LLMs’ feedback before the actual tool execution step, i.e., right after the tool retrieval step.

As illustrated in Figure 3, at each iteration, the LLM will provide feedback on the current-turn retrieval results. Specifically, the LLM will first comprehend the user instruction and tool functionalities thoroughly. Then, it will assess the effectiveness of those retrieved tools for handling the instruction. Based on the assessment, the LLM could provide a refinement to the retrieval model, refining the user instruction if necessary. To ensure that the retriever model is aware of the iteration round, we conduct an iteration-aware feedback training process to adapt the retriever model with continuously refined user instructions.

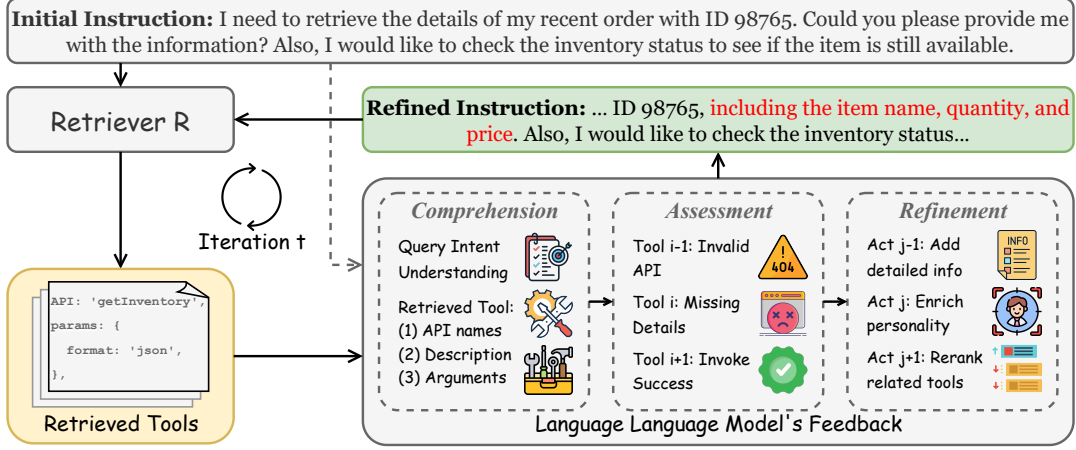


Figure 3: Illustration of our proposed iterative tool retrieval method. At each iteration, the LLM follows a three-step feedback generation process, which includes comprehension, assessment, and refinement, to improve the instruction.

## 4.2 Feedback Generation

Assuming at the iteration step  $t$ , given the refined instruction  $q^t$ , we could utilize retriever model  $R$  to retrieve a list of top- $K$  tools  $\{d_1^t, \dots, d_K^t\}$ . We then conduct a three-step feedback generation process by feeding those retrieved tools and associated tool descriptions into the LLM as follows.

**Comprehension.** Firstly, the LLM is prompted to give comprehension on both the given instruction and retrieved tools. The prompt provided to LLM includes two parts: (1) summarize the abstract user goals by ignoring detailed entity information in the given instruction; (2) understand the functionalities of retrieved tools, focusing on the category, name, description, input and output parameters of given tools. This step can be formulated as,

$$F_C = LLM(P_C, q^t, \{d_1^t, \dots, d_K^t\}), \quad (2)$$

where  $F_C$  denotes LLM’s comprehension output and  $P_C$  denotes the prompt provided to LLM.

**Assessment.** The LLM will assess the effectiveness of retrieved tools for handling the instruction based on its comprehension of the user’s intent and tool functionalities. The assessment is conducted from two perspectives: 1) identify which of the user’s goals could and could not be solved by the retrieved tools with corresponding reasons; and 2) analyze whether the ranked order of retrieved tools corresponds with their significance in addressing the user’s intent with specific reasons. The step can be formulated as,

$$F_A = LLM(P_A, q^t, \{d_1^t, \dots, d_K^t\}, F_C), \quad (3)$$

where  $F_A$  denotes the LLM’s assessment output.

**Refinement.** Lastly, the LLM will refine user instruction based on its assessment. Specifically, we ask the LLM to determine whether the refinement is necessary based on the two following questions: 1) Whether all the user’s goals have been solved by currently retrieved tools, 2) and whether all existing appropriate tools are given the highest ranking priorities by the retriever. If one of the answers is not “yes”, we prompt the LLM to provide a potential refinement for retrieval improvement. Otherwise, the LLM will directly return a special token “N/A” without conducting any refinement.

The feedback from the LLM is finalized made on the current user instruction  $q^t$ . Specifically, we prompt the LLM to generate refined instruction with enriched information in two dimensions: 1) more detailed and personalized content about those user’s intent which have not been solved by current tools, helping the retriever explore other relevant tools; (2) more scenario-specific tool-usage information about existing appropriate tools, helping the retriever give higher ranking priority to those tools. This step can be formulated as,

$$F_R = LLM(P_R, q^{t-1}, \{d_1^{t-1}, \dots, d_K^{t-1}\}, F_A), \quad (4)$$

where  $P_R$  is the corresponding prompt and  $F_R$  denotes LLM’s refinement output, i.e., the new refined instruction  $q^{t+1}$ .

## 4.3 Iteration-Aware Feedback Training

We concatenate a special iteration-aware token “Iteration  $t$ ” in front of the instruction, where  $t$  is the instruction’s iteration step (e.g., “Iteration  $t - 1$ ” for  $q^{t-1}$  and “Iteration  $t$ ” for  $q^t$ ).

We also employ the hard negative sampling in training. Concretely, for each given instruction, we randomly sample an incorrect tool from the retrieved top- $K$  tool list. The high similarity scores of those tools indicate that they are prone to be mistaken as correct tools by the retriever. In feedback training, we utilize those tool-instruction pairs as hard negative samples. Then the loss function for each iteration could be calculated as,

$$\mathcal{L}(q) = -\frac{1}{B} \sum_{i=1}^B \log \frac{e^{R(q_i, d_i^+)}}{e^{R(q_i, d_i^+)} + \sum_{j \neq i} e^{R(q_i, d_{ij}^-)} + \sum e^{M(q_i, d_{ij}^H)}}, \quad (5)$$

where  $B$  denotes the batch size and  $d_{ij}^H$  denotes the hard negative sample. By distinguishing the subtle differences in the tool descriptions, the retriever could achieve a deeper understanding of the tool functionalities and their relation with user instructions.

Then we conduct joint training on the retriever model to solve the refined user instructions across all  $T$  iterations. The final training objective could be formulated as the sum of losses in each iteration as follows,

$$\mathcal{L}_{feedback} = \sum_{t=1}^T \alpha^t \mathcal{L}(q^t), \quad (6)$$

where  $\alpha^t$  is a balancing factor and  $L(q^t)$  is the loss function calculated by Equation 5 based on the refined user instructions  $q^t$  in the  $t$ th iteration. In this way, the LLM’s comprehensive knowledge of the user requirements could be injected into the retriever through those refined instructions. Besides, with the aid of iteration-aware tokens and joint-training manner, the tool retriever could address refined instructions in new iterations while also remembering addressing instructions in previous iterations, ensuring a continuous capability to solve user instructions across all iterations.

#### 4.4 Inference

At the time of inference, the feedback generation process keeps working while the feedback training process ceased. The retriever will update the candidate tool list based on the refined user instruction from LLM’s feedback iteratively, until output the final retrieved tools.

Concretely, assume that we have obtained a retriever  $R$  after the feedback training. For each

	scenarios	# instructions	# tool set
Training Set	ToolBench-I1	86,643	-
	ToolBench-I2	84,270	-
	ToolBench-I3	25,044	-
	ToolBench-All	195,937	-
In-domain Evaluation	ToolBench-I1	796	10,439
	ToolBench-I2	573	13,142
	ToolBench-I3	218	1,605
	ToolBench-All	1,587	13,954
Out-of-domain Evaluation	T-Eval	553	50
	UltraTools	1,000	498

Table 1: Statistics of the TR-bench, which is conducted from ToolBench (Qin et al., 2024), T-Eval (Chen et al., 2024), and UltraTools (Huang et al., 2024a).

initial test instruction  $q_{test}^0$ , we add a special token “Iteration 0” in front of the instruction. Then we use the trained retriever  $R$  to retrieve an initial tool list  $D_{test}^0$ , containing  $K$  candidate tools  $\{d_1, d_2, \dots, d_K\}$ . The retrieved  $D_{test}^0$  and  $q_{test}^0$  will be fed to the LLM for feedback generation, including instruction refinement, as discussed in Section 4.2. After obtaining the refined instruction  $q_{test}^1$ , we add a token “Iteration 1” to it and then input it to  $R$  for the next-round tool retrieval. Then, we can get an updated tool list  $D_{test}^1$  for a new round of feedback generation. As such, we could obtain a final tool list  $D_{test}^T$  after  $T$  iterations.

## 5 Experiments

### 5.1 Setup

**Datasets and evaluation.** To assess the tool retrieval performance of models, we conduct an experiment on tool retrieval benchmark, referred to as **TR-bench**, based on three datasets, including ToolBench (Qin et al., 2024), T-Eval (Chen et al., 2024), and UltraTools (Huang et al., 2024a). To address real-world requirements, we conduct evaluations in both *in-domain* and *out-of-domain* settings. Specifically, the training set is from ToolBench, while the test set of ToolBench is employed for in-domain evaluation, and the test sets from T-Eval and UltraTools are used for out-of-domain evaluation. The statistics of TR-bench are summarized in Table 1.

Following ToolBench, we adopt the Normalized Discounted Cumulative Gain (NDCG) (Järvelin and Kekäläinen, 2002), an ideal metric for tool retrieval to evaluate the quality of retrieved tools. In our evaluation, we report  $NDCG@m$  ( $m = 1, 3, 5, 10$ ), calculated according to the position of each golden tool among top- $m$  candidates tools retrieved by the tool retriever. Thus, the more ac-

Methods	SINGLE-TOOL (I1)			CATEGORY (I2)			COLLECTION (I3)			ALL		
	N@1	N@3	N@5	N@1	N@3	N@5	N@1	N@3	N@5	N@1	N@3	N@5
BM25	18.37	17.97	19.65	11.97	9.85	10.95	25.23	18.95	20.37	15.84	13.98	15.63
Ada Embedding	57.52	54.90	58.83	36.82	28.83	30.68	54.59	42.55	46.83	46.59	41.06	43.95
ToolRetriever	84.20	89.59	89.65	68.24	77.43	77.90	81.65	87.24	87.13	75.73	83.19	83.06
<b>Ours</b>	<b>90.70</b>	<b>90.95</b>	<b>92.47</b>	<b>89.01</b>	<b>85.46</b>	<b>87.10</b>	<b>91.74</b>	<b>87.94</b>	<b>90.20</b>	<b>88.53</b>	<b>87.00</b>	<b>88.83</b>
<b>% improve</b>	7.72%	1.52%	3.15%	30.44%	10.37%	11.81%	12.36%	0.80%	3.52%	16.90%	4.58%	6.95%

Table 2: In-domain evaluation on TR-bench in terms of NDCG@ $m$  under scenarios including single-tool (I1), intra-category multi-tool (I2), intra-collection multi-tool (I3), and the whole data (All). % improve represents the relative improvement achieved by our method over the previously best tool retrieval method.

Methods	T-EVAL				ULTRATOOLS			
	N@1	N@3	N@5	N@10	N@1	N@3	N@5	N@10
BM25	52.12	43.19	45.23	52.91	15.10	14.13	16.03	18.34
Ada Embedding	80.11	69.11	71.95	79.62	31.46	33.75	39.91	46.40
ToolRetriever	82.10	72.03	74.15	<b>80.76</b>	48.20	<b>47.73</b>	53.01	58.93
<b>Ours</b>	<b>84.45</b>	<b>73.31</b>	<b>74.45</b>	80.25	<b>49.30</b>	47.50	<b>54.30</b>	<b>59.92</b>
<b>% improve</b>	2.86%	1.78%	0.40%	-0.06%	2.28%	-0.48%	2.43%	1.68%

Table 3: Out-of-domain evaluation on TR-bench in terms of NDCG@ $m$  under two scenarios, T-Eval (Chen et al., 2024) and UltraTools (Huang et al., 2024a). % improve represents the relative improvement achieved by our method over the previously best tool retrieval method.

curately the tool retriever can retrieve correct tools, the higher the NDCG@ $m$  score will be.

**Baselines.** We compare our method against representative retrieval methods. 1) BM25 (Robertson and Zaragoza, 2009): the classical sparse retrieval method; 2) Ada Embedding: the closed-sourced OpenAI’s text-embedding-ada-002 model<sup>2</sup>; 3) ToolRetriever (Qin et al., 2024): a dense retrieval approach specifically finetuned on tool retrieval datasets.

**Implementation details.** We employ SentenceBERT (Reimers and Gurevych, 2019) to train our retriever model based on BERT-base (Devlin et al., 2019). We set the learning rate to  $2e-5$  with 500 warm-up steps. The batch size in training is set to 64. We utilize ChatGPT (gpt-3.5-turbo-0125)<sup>3</sup> as the LLM for giving feedback. The number of tool candidates  $K$ , the balancing factor  $\alpha$ , and the iteration round  $T$  are set to 10, 1, and 3, respectively. We have trained the model several times to confirm that the improvement is not a result of random chance and present the mid one. Our experiments were conducted on four NVIDIA A6000 GPUs with 48 GB of memory.

<sup>2</sup><https://platform.openai.com/docs/guides/embeddings/embedding-models>.

<sup>3</sup><https://openai.com/index/introducing-chatgpt-and-whisper-apis/>.

## 5.2 Main Results

**In-domain evaluation.** The results of the in-domain evaluation are reported in Table 2. It is observed that non-finetuned retrieval methods, i.e., BM25 and Ada Embedding, perform much worse than other finetuned methods. This is reasonable since non-finetuned methods have not been specifically adopted for tool retrieval. While Tool Retriever outperforms non-finetuned methods, the performance is still not satisfying. In comparison, our proposed method consistently outperforms all finetuned and non-finetuned baselines. Significantly, our method maintains strong performance in the intra-category multi-tool (I2) scenario, even as other methods’ performance declines, demonstrating the robustness of our proposed method across different scenarios. The above results prove the effectiveness of our method in enhancing tool retrieval accuracy, particularly in challenging scenarios with multi-tools.

**Out-of-domain evaluation.** Since the tools are usually frequently updated in real-world, we further test all methods in the out-of-domain setting, where the training data from ToolBench and the test data from T-Eval and UltraTools are used. The experimental results are shown in Table 3. We could observe that our method significantly outperforms other baselines across both scenarios. This demonstrates that our method not only excels in in-domain benchmarks but also maintains robust

Methods	N@1	N@3	N@5	N@10
Ours	89.01	85.46	87.10	88.41
<i>w/o warm-up</i>	85.51	81.36	84.47	86.92
<i>w/o hard-negative</i>	86.04	80.41	84.00	85.98
<i>w/o joint</i>	85.38	81.55	83.79	86.20
<i>w/o joint &amp; hard-neg</i>	83.77	77.67	81.21	83.69

Table 4: Ablation study of our method under the intra-category multi-tool (I2) scenario.

Iteration	N@1	N@3	N@5	N@10	Efficiency
1	85.69	80.48	83.94	86.27	6.12s
2	87.78	83.48	86.31	88.26	8.59s
3	<b>89.01</b>	<b>85.46</b>	<b>87.10</b>	<b>88.41</b>	10.30s

Table 5: Analysis on iteration round under the intra-category multi-tool (I2) scenario. The efficiency is measured by the time consumption to complete one user instruction.

performance across varied scenarios, revealing its generalization ability of tool retrieval.

We further compare the tool usage performance of our method with ToolRetriever in the I2 scenario. We adopt ToolLLaMA (Qin et al., 2024) which is trained on LLM-annotated solution path as the tool usage model, and use “pass rate” and “win rate” as evaluation metrics. Our method achieves 75.6% for pass rate compared to ToolRetriever’s 68.5%, and 65.9% for win rate compared to ToolRetriever’s 60.8%. The results demonstrate the performance improvement in tool usage, benefiting the entire tool learning process.

### 5.3 Ablation Study

We conduct ablation studies to investigate the efficacy of different components in our methods. First, we remove the warm-up training by directly conducting our method on a retriever based on Sentence-BERT. Then, we analyze the contribution of hard negative sampling in our method by removing the hard-to-distinguish samples from the training. In addition, we assess the efficacy of joint training in our method, by substituting it with a loss  $\mathcal{L}_{feedback} = \mathcal{L}(q^t)$ , with respect to only the refined instructions  $q^t$  at current iteration  $t$ . Table 4 reports the ablation test performance (i.e., NDCG@ $m$  ( $m = 1, 3, 5, 10$ )) under the intra-category multi-tool instructions (I2) scenario on ToolBench.

From the results, we can observe that our method achieves comparably high NDCG scores even without warm-up training, indicating that it does not

Methods	N@1	N@3	N@5
ToolRetriever (BERT-based)	68.24	77.43	77.90
<b>Ours (BERT-based)</b>	<b>89.01</b>	<b>85.46</b>	<b>87.10</b>
ToolRetriever (RoBERTa-based)	76.61	69.81	74.99
<b>Ours (RoBERTa-based)</b>	<b>88.13</b>	<b>85.41</b>	<b>86.75</b>

Table 6: Analysis on different base models under the intra-category multi-tool (I2) scenario.

Embedding Size	N@1	N@3	N@5	N@10
300	87.61	83.49	85.20	86.50
512	87.61	82.85	84.67	85.81
<b>768</b>	<b>89.01</b>	<b>85.46</b>	<b>87.10</b>	<b>88.41</b>
1024	88.66	83.91	85.94	87.04
2048	88.74	83.95	85.98	87.43

Table 7: Analysis on embedding sizes under the intra-category multi-tool (I2) scenario.

heavily rely on prior tool-use knowledge. When hard negative sampling is removed, the performance degradation illustrates that hard negative sampling could enable the model to discriminate between similar tool functionalities. Besides, the model’s performance further declines when joint training is removed, demonstrating that the model could balance new and previous knowledge in this joint-training manner.

### 5.4 In-depth Analysis

**Analysis on iteration round.** The iteration round is an important factor in our method. We conduct experiments to investigate changes in effectiveness and efficiency with different iteration round  $T$ . The results are presented in Table 5, and the efficiency is measured by the cost of time to complete one user instruction on average.

By analyzing the results in Table 5, we gain two findings. 1) We could observe a continuous improvement as the iteration round increases. This shows that the tool retriever progressively enhances its performance with the aid of LLMs’ feedback. 2) In terms of time efficiency, we find that adding one additional round of refinement takes an average of 6.12s/instruction, primarily resulting from the time waiting for LLM’s feedback when calling the OpenAI API. As the number of iterations increases, we can see that the extra inference time required for each instruction decreases. This is due to the fact that there will be fewer instructions requiring refinement as retrieval performance improves.

**Analysis on base models.** We further analyze the impact of different base models on the perfor-



Figure 4: Case study on the effect of user instruction refinement through 3 iterations. The original instruction is revised step-by-step, leading to improved retrieval results.

mance. Specifically, we replace the base model BERT in our method with another classic language model, RoBERTa (Liu et al., 2019). The results are shown in Table 6. As we can see, our method still achieves significant improvement over the baseline with the same RoBERTa model. Another observation is that RoBERTa is more effective in serving as a base model for the retrieval application, which benefits from its effective training strategies. The improvements demonstrate the robustness of our method with different base models.

**Analysis on embedding sizes.** Since the retriever model  $R$  encodes the textual instruction and tool description into dense vectors, we explore the impact of the embedding size on retrieval performance. as shown in Table 7. From the table, we can find that larger embedding sizes result in greater performance improvements compared to smaller embedding sizes. This is probably due to the fact that embeddings with larger sizes could accommodate more knowledge. However, when the embedding size increases from 768 to 2048, there is a slight decrease in performance. This suggests that a specific embedding size is sufficient, and larger embedding sizes may pose challenges to training. It is worth noting that larger embedding sizes neces-

sitate higher training costs and increased inference memory. Therefore, we recommend an optimal embedding size of 768.

## 5.5 Case Study

As shown in Figure 4, we conduct case study by using an example of instruction refinement to take a closer look at the effect of our method.

In the *1st* iteration, we can observe that the refined instruction has included more detailed information (i.e., “total number”) about the user’s requirements than the original instruction, enabling the retriever to identify more appropriate tools (e.g., Check residential proxies service status). This reveals that the comprehension capabilities of LLMs could be instilled into the retrieval process through feedback. In the *2nd* iteration, our method further refines the instruction by omitting irrelevant content (i.e., “information”) which may mislead the retriever into retrieving incorrect tools (e.g., Retrieve Proxy Information). Another benefit of the refinement is that some correct tools (e.g., Bash Code Compiler) will move up in positions of the top- $K$  rankings, improving the overall retrieval performance. In the *3rd* iteration, our method showcases great decision-aware capabilities, where the iterative process could be terminated if no further



refinement is deemed necessary.

## 6 Conclusion and Future Work

In this study, we concentrate on the crucial tool retrieval in the tool learning of LLMs. We have identified the bottleneck in the tool retrieval-usage pipeline as the limited tool retrieval model. We propose the unique challenges of the tool retrieval compared with document retrieval. To improve the current tool retrieval process, we propose leveraging the LLM’s feedback to assess the retrieval results and provide detailed suggestions for refining user instructions. In order to integrate the retriever model into this iterative process, we implement iteration-aware feedback training. This will improve the tool retriever’s capabilities and close the gap between tool retrieval and usage models. We conduct the TR-benchmark to comprehensively evaluate the models’ ability in real-world tool retrieval scenarios. Our method demonstrates the best performance in both in-domain and out-of-domain settings.

In the future, we aim to improve this work from the following aspects. 1) Limited by the training speed, we have applied the offline feedback generation, where feedback is generated before training the tool retriever. We will also assess whether online feedback generation yields further improvements in the future. 2) Furthermore, as the tool retriever serves the subsequent tool usage model in tool learning, we intend to conduct further evaluations of the tool retriever models based on the subsequent tool usage results.

### Limitations

1) Undoubtedly, our iterative refinement will reduce the inference speed of the tool retrieval. The efficiency issue is inherent in approaches involving LLMs’ interaction. We have evaluated the efficiency as the number of iterative rounds increases. Fortunately, we observed that the retrieval model can achieve a significant performance improvement after just a single round of LLMs’ feedback compared to without feedback. Furthermore, the performance enhancement of the tool retrieval is crucial for the subsequent tool usage model, ensuring that the correct tools are retrieved and lays the foundation for all subsequent steps of tool usage. Therefore, we believe that performance improvement is worthwhile despite some efficiency loss. We will also pay more attention to this issue in the future.

2) Similar to document retrieval, the used datasets in our work also contain “false negative” samples. For instance, some tools may be capable of handling the user’s instruction but are not labeled as positive. This can disrupt the training and evaluation of tool retrieval and is a common limitation in many retrieval scenarios.

### Ethics Statement

The datasets used in our experiment are publicly released and labeled through interaction with humans in English. In this process, user privacy is protected, and no personal information is contained in the dataset. The scientific artifacts that we used are available for research with permissive licenses. And the use of these artifacts in this paper is consistent with their intended use. Therefore, we believe that our research work meets the ethics of the conference.

### Acknowledgments

The work described in this paper was supported by National Natural Science Foundation of China (62076212), Research Grants Council of Hong Kong (PolyU/15207821, PolyU/15207122, PolyU/15213323, and PolyU/15209724), and PolyU internal grants (ZVQ0).

### References

- Raviteja Anantha, Bortik Bandyopadhyay, Anirudh Kashi, Sayantan Mahinder, Andrew W Hill, and Srinivas Chappidi. 2023. Protip: Progressive tool retrieval improves planning. [arXiv preprint arXiv:2312.10332](https://arxiv.org/abs/2312.10332).
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024. [Self-RAG: Learning to retrieve, generate, and critique through self-reflection](#). In [The Twelfth International Conference on Learning Representations](#).
- Zehui Chen, Weihua Du, Wenwei Zhang, Kuikun Liu, Jiangning Liu, Miao Zheng, Jingming Zhuo, Songyang Zhang, Dahua Lin, Kai Chen, and Feng Zhao. 2024. T-eval: Evaluating the tool utilization capability of large language models step by step. In [Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics \(Volume 1: Long Papers\)](#), pages 9510–9529. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. 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 4171–4186. Association for Computational Linguistics.
- Yu Du, Fangyun Wei, and Hongyang Zhang. 2024. Anytool: Self-reflective, hierarchical agents for large-scale API calls. In Forty-first International Conference on Machine Learning.
- Shen Gao, Zhengliang Shi, Minghang Zhu, Bowen Fang, Xin Xin, Pengjie Ren, Zhumin Chen, Jun Ma, and Zhaochun Ren. 2024. Confucius: Iterative tool learning from introspection feedback by easy-to-difficult curriculum. In Proceedings of the AAAI Conference on Artificial Intelligence, pages 18030–18038.
- Shibo Hao, Tianyang Liu, Zhen Wang, and Zhitong Hu. 2023. Toolkengpt: Augmenting frozen language models with massive tools via tool embeddings. In Advances in Neural Information Processing Systems, volume 36, pages 45870–45894. Curran Associates, Inc.
- Shijue Huang, Wanjun Zhong, Jianqiao Lu, Qi Zhu, Jiahui Gao, Weiwen Liu, Yutai Hou, Xingshan Zeng, Yasheng Wang, Lifeng Shang, Xin Jiang, Ruifeng Xu, and Qun Liu. 2024a. Planning, creation, usage: Benchmarking LLMs for comprehensive tool utilization in real-world complex scenarios. In Findings of the Association for Computational Linguistics ACL 2024, pages 4363–4400. Association for Computational Linguistics.
- Tenghao Huang, Dongwon Jung, Vaibhav Kumar, Mohammad Kachuee, Xiang Li, Puyang Xu, and Muhao Chen. 2024b. Planning and editing what you retrieve for enhanced tool learning. In Findings of the Association for Computational Linguistics: NAACL 2024, pages 975–988. Association for Computational Linguistics.
- Yue Huang, Jiawen Shi, Yuan Li, Chenrui Fan, Siyuan Wu, Qihui Zhang, Yixin Liu, Pan Zhou, Yao Wan, Neil Zhenqiang Gong, and Lichao Sun. 2024c. Meta-tool benchmark for large language models: Deciding whether to use tools and which to use. In The Twelfth International Conference on Learning Representations.
- Kalervo Järvelin and Jaana Kekäläinen. 2002. Cumulated gain-based evaluation of ir techniques. ACM Transactions on Information Systems, 20(4):422–446.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 6769–6781. Association for Computational Linguistics.
- Minghao Li, Yingxiu Zhao, Bowen Yu, Feifan Song, Hangyu Li, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. 2023a. 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. Association for Computational Linguistics.
- Yongqi Li, Wenjie Wang, Leigang Qu, Liqiang Nie, Wenjie Li, and Tat-Seng Chua. 2024. Generative cross-modal retrieval: Memorizing images in multimodal language models for retrieval and beyond. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 11851–11861. Association for Computational Linguistics.
- Yongqi Li, Nan Yang, Liang Wang, Furu Wei, and Wenjie Li. 2023b. Generative retrieval for conversational question answering. Information Processing and Management, 60(5):103475.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. Self-refine: Iterative refinement with self-feedback. In Advances in Neural Information Processing Systems, volume 36, pages 46534–46594. Curran Associates, Inc.
- Bhaskar Mitra and Nick Craswell. 2017. Neural models for information retrieval. arXiv preprint arXiv:1705.01509.
- Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. 2023. Gorilla: Large language model connected with massive apis. arXiv preprint arXiv:2305.15334.
- Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen, Ning Ding, Ganqu Cui, Zheni Zeng, Yufei Huang, Chaojun Xiao, Chi Han, et al. 2023. Tool learning with foundation models. arXiv preprint arXiv:2304.08354.
- 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. 2024. ToolLLM: Facilitating large language models to master 16000+ real-world APIs. In The Twelfth International Conference on Learning Representations.
- Changle Qu, Sunhao Dai, Xiaochi Wei, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, Jun Xu, and Ji-Rong Wen. 2024a. Tool learning with large language models: A survey. arXiv preprint arXiv:2405.17935.

- Changle Qu, Sunhao Dai, Xiaochi Wei, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, Jun Xu, and Ji-Rong Wen. 2024b. Towards completeness-oriented tool retrieval for large language models. In Proceedings of the 33rd ACM International Conference on Information and Knowledge Management.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3982–3992. Association for Computational Linguistics.
- Stephen Robertson and Hugo Zaragoza. 2009. The probabilistic relevance framework: Bm25 and beyond. Foundations and Trends in Information Retrieval, 3(4):333–389.
- Jingqing Ruan, Yihong Chen, Bin Zhang, Zhiwei Xu, Tianpeng Bao, Hangyu Mao, Ziyue Li, Xingyu Zeng, Rui Zhao, et al. 2023. Tptu: Task planning and tool usage of large language model-based ai agents. In NeurIPS 2023 Foundation Models for Decision Making Workshop.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. In Advances in Neural Information Processing Systems, volume 36, pages 68539–68551. Curran Associates, Inc.
- Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. 2023. Hugging-gpt: Solving ai tasks with chatgpt and its friends in hugging face. In Advances in Neural Information Processing Systems, volume 36, pages 38154–38180. Curran Associates, Inc.
- Qiaoyu Tang, Ziliang Deng, Hongyu Lin, Xianpei Han, Qiao Liang, and Le Sun. 2023. Toolalpaca: Generalized tool learning for language models with 3000 simulated cases. arXiv preprint arXiv:2306.05301.
- Boshi Wang, Hao Fang, Jason Eisner, Benjamin Van Durme, and Yu Su. 2024a. LLMs in the imagination: Tool learning through simulated trial and error. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 10583–10604. Association for Computational Linguistics.
- Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lifan Yuan, Hao Peng, and Heng Ji. 2024b. MINT: Evaluating LLMs in multi-turn interaction with tools and language feedback. In The Twelfth International Conference on Learning Representations.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. ReAct: Synergizing reasoning and acting in language models. In International Conference on Learning Representations (ICLR).
- Siyu Yuan, Kaitao Song, Jiangjie Chen, Xu Tan, Yongliang Shen, Kan Ren, Dongsheng Li, and Deqing Yang. 2024. EASYTOOL: Enhancing LLM-based agents with concise tool instruction. In ICLR 2024 Workshop on Large Language Model (LLM) Agents.
- Wayne Xin Zhao, Jing Liu, Ruiyang Ren, and Ji-Rong Wen. 2024. Dense text retrieval based on pretrained language models: A survey. ACM Transactions on Information Systems, 42(4):1–60.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. arXiv preprint arXiv:2303.18223.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. In Advances in Neural Information Processing Systems, volume 36, pages 46595–46623. Curran Associates, Inc.
- Yuanhang Zheng, Peng Li, Wei Liu, Yang Liu, Jian Luan, and Bin Wang. 2024. ToolRerank: Adaptive and hierarchy-aware reranking for tool retrieval. In Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 16263–16273. ELRA and ICCL.