

An Evaluation Mechanism of LLM-based Agents on Manipulating APIs

Anonymous ACL submission

Abstract

LLM-based agents can greatly extend the abilities of LLMs and thus attract sharply increased studies. An ambitious vision – serving users by manipulating massive API-based tools – has been proposed and explored. However, we find a widely accepted evaluation mechanism for generic agents is still missing. This work aims to fill this gap. We decompose tool use capability into seven aspects and form a thorough evaluation schema. In addition, we design and release an instruction dataset and a toolset – the two sides that the agents bridge between – following the principle of reflecting real-world challenges. Furthermore, we evaluate multiple generic agents. Our findings can inspire future research in improving LLM-based agents and rethink the philosophy of API design.

1 Introduction

Large Language Models (LLMs) exhibit remarkable capabilities across a variety of tasks, such as language, mathematics, coding, and etc (Bubeck et al., 2023). However, they still face some limitations, such as having frozen knowledge, being bad at some specialized tasks like calculation, and not being able to ground their generated solution outlines to the real world. Meanwhile, there are existing systems or models that can perform very well on domain-specific tasks. Therefore, a mechanism that links LLMs with the existing ecosystem of tools can bring the ability of LLM-based AI to another level.

To stretch the ability of LLMs, a sharply increasing number of works have studied LLM-based agents¹ which can manipulate API-based tools. A very ambitious vision is to build a new AI ecosystem that connects LLMs with millions of APIs, assessed via an API platform, for task completion.

¹The term *agent* denotes *LLM-based agent* by default in this work.

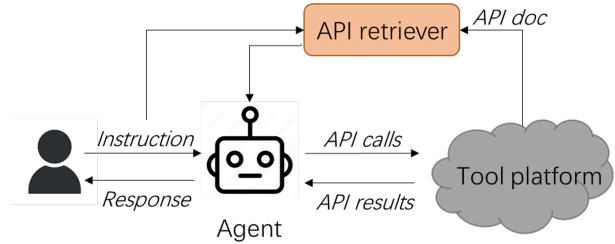


Figure 1: Agent connecting a user with massive APIs.

As shown in Fig. 1, the agent acts like a super-APP. It eases user interaction with language rather than GUI, manipulates massive API-based tools and thus supports a mass of functionalities. This requires the agent to have, on one hand, rich knowledge to deal with different user needs, and on the other hand developer’s skills for manipulating APIs given documentation and understanding the results of API execution.

Though the research community has made an effort to build up generic LLM-based agents, we found that a widely accepted evaluation mechanism for LLM-based agents on tool use is still lacking. This prevents researchers from making fair comparisons between different agent systems, as well as gaining insights into the challenges of designing agents. In this work, we aim to narrow this gap by providing an evaluation mechanism, putting special emphasis on discovering the limitations of existing agents as well as the problems of current API design philosophy.

It is very challenging to evaluate an agent because of the complex process of an agent performing user tasks. Once a user issues an instruction, the agent first decomposes it into solvable subtasks with available tools, depending on the complexity of the task. Then, it may need to collect user needs by interacting with the user and call multiple APIs. Eventually, it responds to the user as per the results of APIs. The problems this complex

process brings to agent evaluation are: (1) the possible involvement of users makes the evaluation very hard. (2) the agent may fail in any step of the sequential process, making the samples in the later stage rare. (3) executing a complex process may involve multiple aspects of capability. However, an end-to-end performance cannot help locate the weaknesses of agents and gain more insights. To solve these problems, we dissect the whole process into intermediate decision behaviors to get a thorough view of the involved capabilities. According to this anatomy, we include 7 aspects of tool use capability in our evaluation schema. Each aspect corresponds to one separate evaluation subtask without involving human users.

Under the guide of our evaluation schema, we build a toolset and a dataset of instructions, following the principle of reflecting real-world challenges. We construct our toolset by addressing a series of concerns as demonstrated below rather than assembling some random tools. (1) To ensure reflection of real-world challenges, we collect tools from the real API platform. It matters to reflect the properties of API design and documentation. (2) We intentionally control the diversity of tools regarding functionalities and API structures. According to our observation, a tool may have a single API, a list of APIs, or several collections of APIs. These organizational structures can indicate the variety of functionality and expose different difficulties in calling. (3) We find some APIs depend on other APIs². We include this kind of dependency relationship in our toolset. (4) To enable high-level tasks applied to the toolset, we take measures to increase the coherence of tools regarding application scenarios. (5) We noticed that the affordability of toolset can be one potential problem for individual researchers. To avoid this problem, we devoted lots of engineering work to make the toolset usable at low or even no cost.

On top of our toolset, we construct a set of user instructions that may use the contained tools to solve. We analyze how humans ask questions and summarize five types of instructions, varying in user intentions and complexities. These different types of instructions can be used to produce evaluation data required by our evaluation schema. Additionally, we emphasize that the instructions should be in the real way of user expression. Only by

²This is caused by the principle of API design – being simple and general.

this can the evaluation data imply the mismatch between user expression and the form of information required by APIs.

We make our dataset and toolset publically available at [[placeholder of GitHub repo]]³.

2 Related Works

2.1 LLM-based Agents

LLM is the core component of an LLM-based agent. In the LLM domain, ChatGPT (Ouyang et al., 2022) is the most typical proprietary LLM and represents the SOTA LLM while many open-sourced LLMs like LLaMA (Touvron et al., 2023) are also very competitive. These LLMs have shown impressive language ability, rich knowledge, great potential in reasoning, and unbelievable generality in Question Answering tasks (Bubeck et al., 2023).

These characteristics of LLMs naturally inspire researchers to use LLMs as the brain of agents, which are designed to interact with complex environments and are closer to general AI.

In the surging literature of LLM-based agents, the agents have been explored to (1) manipulate external tools to solve more complex tasks (Nakano et al., 2021; Song et al., 2023; Shen et al., 2023); (2) play games (Zhu et al., 2023; Xu et al., 2023b); (3) form a multi-agent system which can do big projects collaboratively (Qian et al., 2023; Taleb-rad and Nadiri, 2023); (4) and even be embedded in robots to interact with the physical world (Wang et al., 2023; Ichter et al., 2022). This work focuses on the tool use ability of agents.

Some works have explored connecting LLMs with a pre-specified set of tools. By enabling LLMs to manipulate tools, the LLMs can access more information than that frozen in the weights (Nakano et al., 2021), overcome the shortcomings of LLMs like calculation (Schick et al., 2023), and complete more complex tasks than QA (Zhou et al., 2023; Shen et al., 2023). As proof of concept, these works demonstrate that equipping LLMs with tool use ability is a promising direction.

We distinguish between close-world settings and open-world settings. The close-world settings usually have a few special properties: (1) the number of tools is usually limited; (2) the design of APIs tends to be simplified to ease the calling by LLMs. (3) considering the toolset will not be updated frequently, optimizing LLMs for the toolset is feasible,

³Our dataset and toolset will be released upon acceptance.

for example, by constructing toolset-specific training data to fine-tune the LLMs.

On the opposite, in the open-world settings (Liang et al., 2023; Patil et al., 2023), (1) the number of APIs API platform can be massive and may keep increasing; (2) the APIs are designed in an LLM-agnostic way and documented by following a unified schema required by the API platform. (3) the tools available in the API platform always keep changing. Our work is for evaluating agents designed for the second setting.

Different forms of tools have been considered in the literature, such as APIs (Patil et al., 2023; Qin et al., 2023), websites (Deng et al., 2023; Yao et al., 2022) and mobile APPs (Zhang et al., 2023; Hong et al., 2023; Rawles et al., 2023). We divide these tools into two categories: API-based tools and UI-based APPs (e.g. websites, desktop software and mobile APPs), as per the different challenges they pose to the LLMs. This work aims to serve the investigation of agents on manipulating API-based tools.

An agent system basically consists of an LLM and an inference pipeline. The LLM is injected with the ability of manipulating tools by fine-tuning (Tang et al., 2023; Qin et al., 2023; Patil et al., 2023) or in-context learning (Shen et al., 2023; Xu et al., 2023a). In addition, considering the complexity of tool use tasks, the inference process is usually enhanced with more sophisticated reasoning (Yao et al., 2023), searching method of solution path (Qin et al., 2023), and etc. To facilitate the creation of agent systems, a few open-sourced frameworks have been released (Li et al., 2023a; Qin et al., 2023).

2.2 Evaluation of LLM-based Agents

Many early works evaluated their agent systems with their own evaluation suites, making it hard to compare different agents. In addition, these works usually only report the end-to-end performance, which is inadequate for gaining insight into the tool use ability. To address these problems, more and more effort has been put into benchmarking the existing agents.

T-Eval (Chen et al., 2023) is close to our work. they designed an inference process composed of several steps and compared the overall and step-wise performance of several LLMs as the backbones. Our goal is to design a universal evaluation mechanism which is not coupled with the design of agents.

AgentBoard (Ma et al., 2024) created one evaluation toolkit to test the generality of agents across different types of environments, such as Embodied environments, game environments, and tool environments. It is not dedicated to in-depth evaluation of tool use ability.

3 Methodology

We make our evaluation schema first to make clear our targets. Then, we take measures to construct a toolset inheriting real-world challenges and supporting tasks with varying complexities. Extra engineering work makes its usage affordable. Following this, we design five types of instruction data and align them to the evaluation schema. Finally, we determine the metrics of each evaluation sub-task.

3.1 Evaluation Schema

We go through the process of an agent performing a task and point out the capabilities involved in different stages. As shown in Fig. 2, the process starts with a user sending an instruction. Depending on its complexity⁴, the agent may need to make a solution outline via planning.

- **Planning**, i.e. decomposing a complex task into several simple subtasks, each of which is solvable with a single API. Considering the varying complexity of user instructions as well as the mismatch between user needs and the design of APIs, planning would be very commonly used by agents.

When dealing with a simple task, the agent first needs to decide whether to use external tools. If yes, it then retrieves a few candidates of potentially useful tools and selects one of them according to the documentation of tools. Otherwise, it replies to the user directly.

- **Deciding whether to use tools**. Failing to trigger tool use when needed makes tasks not solved, while misusing tools can hurt LLM’s performance in ordinary QA tasks.
- **Selecting useful tools**. If tools are required, the agent should be able to figure out the useful ones.
- **Responding with intrinsic knowledge**. There is a risk that fine-tuning an LLM with tool-use data makes it lose its original ability. Therefore,

⁴We measure the complexity of a certain task with the number of tools required to complete this task.

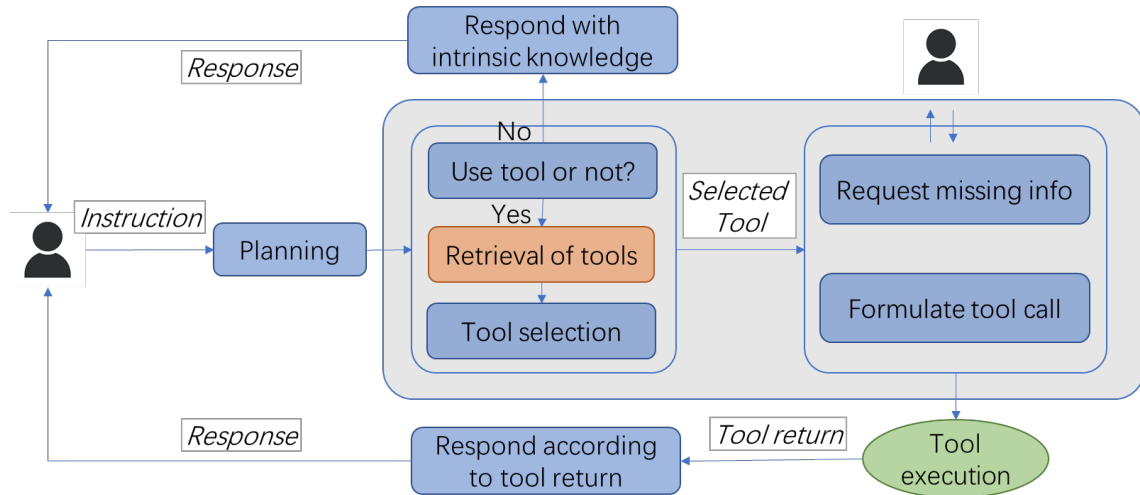


Figure 2: The process of an agent performing a task. Seven decision behaviors can potentially be involved. We examine the performance of an agent in each of them to gain a thorough understanding of its tool-use capability.

we also check whether LLMs still retain their intrinsic knowledge.

To call a certain tool, the agent needs to parse required parameter information from the context (i.e. conversation). If required information is not provided, the agent should ask the user to make clarifications.

- **Requesting missing parameter information.** It's common for users to initiate a dialogue with partial information. In this case, the agent should have the consciousness of requesting clarification rather than hallucinating.

- **Formulating tool calls.** When sufficient information is provided, the agent should be able to parse it and convert it to a valid format as per the specification of APIs. Here, one challenge to overcome is the mismatch between the user expression and the required format of parameters. Sometimes, commonsense reasoning is required.

Eventually, after receiving the execution results of tools, the agent synthesizes a final response to the user.

- **Responding according to the tool returns.** The variety of tools demands the agent to have great generality so that it can interpret the results of APIs, which are usually in JSON instead of natural language, and eventually generate a concise answer.

In summary, our evaluation schema includes seven types of capabilities potentially involved in the process of tool use.

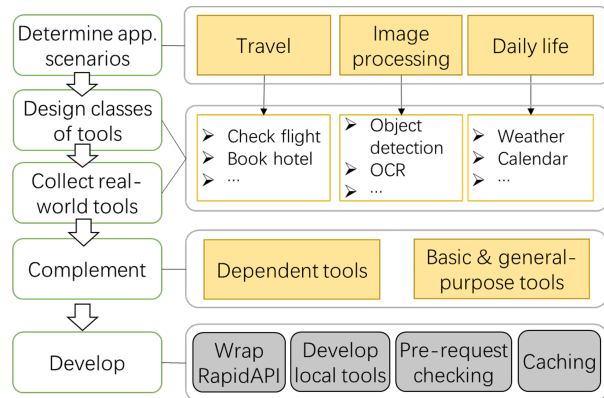


Figure 3: The process of building our toolset. The left side shows the steps while the right side illustrates their corresponding details.

3.2 Toolset

Fig. 3 illustrates the process of developing our toolset, as detailed below.

Determining application scenarios. To achieve high coherence of tools, we start with selecting a few application scenarios (e.g. travel, image processing) of agents. Within each scenario, the tools have a relatively high chance of being combined to solve user's needs. In addition, having different scenarios helps ensure the diversity of tools.

Designing classes of tools. For each scenario, we think about the potentially useful tools. Then, instead of collecting the tools directly, we design the classes of tools to provide an umbrella, under which the tools from different sources and implemented by different people can be integrated. Here, each tool class is defined with a set of main func-

308 tionalities.

309 **Collecting real-world tools with documentation.**

310 For each tool class, we look for its real-world imple-

311 mentation from a well-known API platform Rap-

312 idAPI ⁵, where massive APIs are deployed and

313 documented with a unified schema. This is to make

314 sure the design and documentation of APIs can

315 reflect **real-world problems**.

316 In RapidAPI, most tools do not only have a sin-

317 gle API (very typical in existing works) but mul-

318 tiple ones organized with a list or multiple collec-

319 tions. These multi-API tools may not only make the

320 documentation of different APIs entangled but also

321 comprise **dependencies** between APIs. For exam-

322 ple, an API for checking flights takes airport codes

323 as input, while checking these codes given airport

324 names need to use another API (see Appendix A

325 for a detailed example). This raises the difficulty

326 level of manipulating APIs for agents. Therefore,

327 we intentionally include tools with different API

328 structures.

329 **Complementing tools.** The **dependencies** be-

330 tween APIs occur not only within a single tool

331 but also across the boundary of tools. This under-

332 lying reason is that the philosophy of designing

333 APIs is to make them simple and general utilities

334 of many different APPs. However, for a toolset, a

335 certain API’s functionality will not be usable unless

336 its dependent APIs are also included. To avoid this

337 problem, we analyze the dependent APIs of already

338 collected APIs and collect them in our toolset.

339 Additionally, we add a few **basic and general-**

340 **purpose tools** (e.g. calculator, search engine, code

341 interpreter).

342 **Development.** We first wrap remote tool services

343 deployed on RapidAPI and develop a few local

344 tools, forming an initial toolset. The tool services

345 deployed on API platforms are usually not free.

346 Even though we have tried to select the tool with

347 the most free quota when collecting tools, a portion

348 of them provide very limited free requests. Fre-

349 quent access to the tool services can cause high

350 subscription fees – an obstacle for research. To

351 address these problems, we take the following mea-

352 sures: (1) developing **free alternatives** to some

353 tools while reusing their documentation and API

354 designs. (2) adding a **caching** mechanism to avoid

355 repeated requests. (3) **check the validity** of API

356 calls before sending them to the remote services.

⁵<https://rapidapi.com/hub>

3.3 Instruction Data 357

358 We design five types of instructions that can be

359 used to evaluate all aspects of tool use capability in

360 our evaluation schema.

361 **Types of Instructions.** Our first three types are

362 low-level instructions, which can be solved mainly

363 with the functionality of a certain API. We con-

364 struct these instructions for each API in turn.

365 Type-I: Instructions that do not need tools to solve

366 but may mislead agents to call tools. For ex-

367 ample, for the question "What’s the weather

368 **usually** like in London", one agent may call

369 real-time weather API if they cannot understand

370 the nuance caused by "usually". This type of

371 instruction can be used to test two abilities: *de-*

372 *deciding whether to use tool*, and *Responding with*

373 *intrinsic knowledge*.

374 Type-II: Instructions that need to use tools and

375 provide sufficient information for formulating

376 function calls. With this type of instruction,

377 we can check whether an agent can parse or

378 infer parameter information correctly from user

379 questions. Also, because this type of instruction

380 is relatively easy, an agent has a higher chance

381 of getting a final assistant response (rather than

382 being interrupted by invalid function calls). We

383 thus can check whether an agent can make a

384 proper response according to the return of a

385 tool.

386 Type-III: Instructions that need to use tools but

387 provide insufficient information for filling pa-

388 rameters. For example, "Can you check the

389 weather for me?". The agents would need to ask

390 the user for his location. This type of instruc-

391 tion is very common and thus very important

392 for evaluating agents. It can be used to check

393 whether an agent can make multi-round inter-

394 actions with the user consciously to solve the

395 user’s need.

396 From the Type-II and Type-III instructions, we fil-

397 ter out the data produced for APIs having depen-

398 dencies. Then, we collect the Type-II ones from

399 them to form Type-IV instructions.

400 Type-IV: instructions that are not complex but still

401 need to use multiple APIs to solve because of

402 the dependencies of APIs.

403 Lastly, we create high-level instructions issuing

404 complex tasks:

Table 1: Eight evaluation tasks and their used instruction data.

	Type-I	Type-II	Type-III	Type-IV	Type-V
Task-1: Deciding whether to use tools	✓	✓			
Task-2: Tool selection		✓			
Task-3: Requesting user to clarify missing info			✓		
Task-4: Filling parameter values		✓			
Task-5: Responding with intrinsic knowledge	✓				
Task-6: Responding according to tool returns		✓			
Task-7: Planning for resolving dependency				✓	
Task-8: Planning for high-level task					✓

Type-V: instructions that require to be decomposed into solvable sub-tasks by APIs. For example, "Plan a seven-day trip in Dubai for me". To complete this task, an agent would, for example, check the weather and search for interesting spots to visit.

Generating Instructions. We generate initial instruction data by prompting GPT4. Apart from the special requirements for each instruction type, we include the following general rules: (1) the instructions should be asked in the real way of human speaking. (2) human users do not mention API in their questions. For each type of instruction, we use one generator to generate instructions first and then use one discriminator to filter out invalid ones. Additionally, human annotators double-check the instructions to ensure high quality and avoid ethical issues ⁶.

Aligning instructions to evaluation schema. In Table 1, we enumerate the evaluation tasks and the corresponding types of instructions to use. In tasks-1,3,4,5, the agent is given an instruction and its corresponding API specification. In task-2, the agent is given an instruction, a correct API along with a few perturbing APIs. In task-6, the agent is given a conversation history including a user message, an assistant message containing a function call and a tool result. In task-7 and 8, the agent is given an instruction and multiple APIs, expecting tool use response and chat response respectively.

3.4 Assessment & Metrics

We assess the performance of an agent for each data instance as below:

- Task-1: whether correct decision has been made for the two types of instruction. Overall precision, recall and macro-F1 score can be computed.

⁶Two of our authors participate in data annotation.

- Task-2: whether the right tool is chosen from the given candidates.
- Task-3: whether the response is to request clarification of missing information.
- Task-4: percentage of correct function call.
- Task-5: whether the response is related to the question. Answer quality is our concern.
- Task-6: whether the response is based on the tool results and whether desired information in the results is interpreted precisely.
- Task-7: the chain of function calls is compared with a ground-truth order of actions. The rate of progress is used as metrics.
- Task-8: whether the solution outline is sound – the coverage of provided APIs and whether the functionality of each API is correctly understood.

The assessing scripts, implemented by mixing rules and GPT-4 usage, are included in our evaluation mechanism too.

4 Experiment

In this section, we first demonstrate more details of our dataset and toolset. Then, we apply our evaluation mechanism to examine several well-known LLMs equipped with generic tool use ability, including ChatGPT series – GPT-3.5-turbo and GPT-4-8k (abbreviated as GPT-3.5 and GPT-4 below), and Qwen1.5 series with sizes 7b, 14b and 72b (abbreviated as Qwen-7b, Qwen-14b and Qwen-72b) (Bai et al., 2023). The new findings can show the value of our evaluation mechanism.

4.1 Dataset & Toolset

Dataset. The numbers of different types of instructions are shown in Tab. 2, while the size of data for each evaluation task is shown in Tab. 3.

Table 2: Number of each type of instruction.

Type	I	II	III	IV	V
Num	372	326	195	141	50

Table 3: Data number of each evaluation task.

Task	1	2	3	4	5	6	7	8
Num	698	326	195	326	372	242	141	50

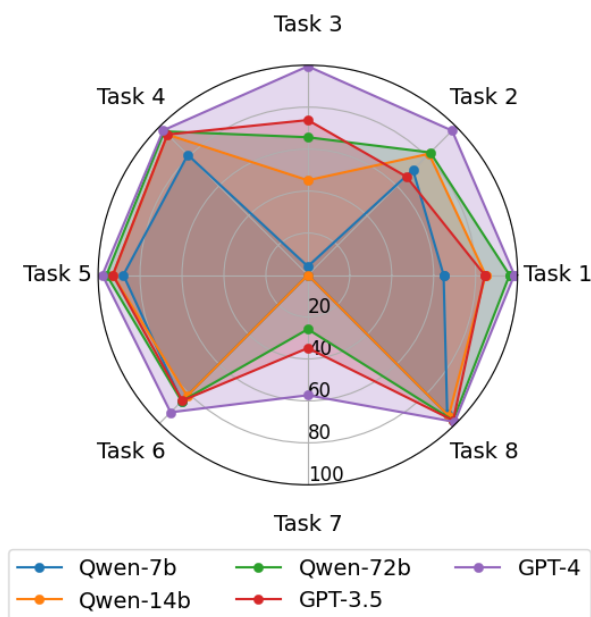


Figure 4: Comparing the performance of five generic agents in eight evaluation tasks. Metrics values can be found in Tab. 5.

Toolset. Our toolset consists of 66 APIs organized into 27 tools. 83% of these APIs are implemented based on API services from RapidAPI, while the other 17% are developed from scratch. We recognize 28 pairs of (API, dependent APIs). In addition, we combine 9 groups of coherent APIs for supporting high-level tasks. See Appendix B for concrete tool classes, and functionalities.

4.2 Evaluation of Generic LLM-based Agents

In Fig. 4, we compare the performance of agents in 8 evaluation tasks.

Task-1: On the decision of tool utilization. (1) We found Qwen-7b and Qwen-14b have the problem of misusing tools – tending to use tools once given. This leads to relatively low precision in their tool-use decision. (2) On the contrary, GPT-3.5 is conservative in tool use – tending not to use tools even needed – resulting in a low recall. (3) GPT-4 and Qwen-72b can make proper decisions, above

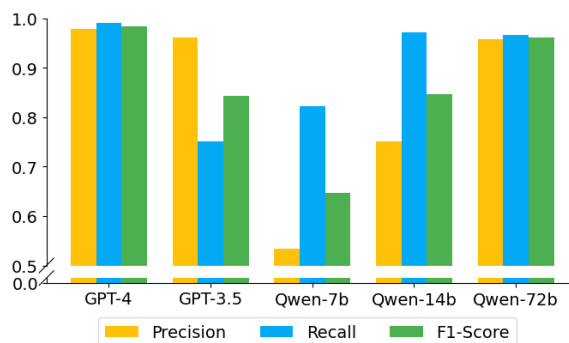


Figure 5: Performance of agents in tool-use decision: precision, recall and macro-F1.

0.96 in macro-F1 scores.

Task-2: On tool selection. To check whether the agents can figure out the right API to use, we provide the agents with one correct API along with four perturbing ones ⁷.

GPT-3.5 and Qwen-7b perform the worst in tool selection, however, for different reasons. Among the mistakes made by GPT-3.5, 77% is because of its conservativeness again – did not call any API, while only 23% are incorrect selections. Qwen-7b always selects the wrong tools, showing its weakness in understanding API specifications.

Compared with GPT-4, the Qwen-14b and Qwen-72b achieved accuracies less than 83%, having a big gap from GPT-4.

Task-3: On the awareness of requesting clarification. It is very often that one user initiates a dialogue with partial information. This requires the agent to figure out the missing parameter information for calling a certain tool and ask the user to clarify. However, we found that the three open-sourced LLMs – Qwen series – are bad at this, worse than both GPT-3.5 and GPT-4. GPT-4 performs almost perfectly while GPT-3.5 still has big space for improvement.

In this case, some typical mistakes include: (1) hallucinating parameter values (no evidence can be found from the user’s questions). (2) using a placeholder-like value (e.g. /path/to/image) instead of a real value. (3) imprecise parameter values (e.g. a location parameter requires a city name but is given a country name) are used, leading to exceptions in executing APIs. (4) required parameters are missing in the function calls.

⁷We select the perturbing APIs which have top-4 highest similarities with the target API regarding sentence embeddings encoded with gte-base-en-v1.5 (Li et al., 2023b).

528	Task-4: On the correctness of function calls.	ferent challenges. For LLM-based agents, talking	578
529	When sufficient information is contained by the	is much easier than doing.	579
530	user’s questions, most LLMs including Qwen-14b	Next, we discuss a bit more from other angles.	580
531	have over 94% correct function calls, except Qwen-		
532	7b achieving around 80%. These numbers are	On scaling law. Though not a new finding any-	581
533	pretty decent. We reckon the reason is most ex-	more, the scaling law still applies in API manipu-	582
534	isting works focus on this setting while neglecting	lation scenarios. The performance of Qwen series	583
535	the others.	reveals larger LLMs have better performance re-	584
		garding almost every aspect of tool use capability.	585
536	Task-5: On the utilization of intrinsic knowl-		
537	edge. We empirically notice that, in some exist-	On API design and quality of documentation.	586
538	ing works (e.g. Qin et al., 2023), fine-tuning LLMs	The effect of API design and documentation quality	587
539	with tool use data makes the LLMs forget their orig-	deserves more attention. A few concrete examples	588
540	inal capabilities in ordinary Question Answering	are: (1) A tool with multiple APIs may introduce	589
541	(QA) tasks. Fortunately, this did not happen in the	its functionalities in its tool-level description while	590
542	generic agents we evaluated. Intuitively, it is not	giving very unclear API-level descriptions. (2) The	591
543	hard to achieve since QA is a more basic task for	execution results of APIs have poor readability or	592
544	generic LLMs. Even though, we still consider keep-	are too verbose. (3) The APIs with dependencies	593
545	ing this aspect in our evaluation schema to remind	seem too complex for the LLMs to use.	594
546	the phenomenon of over-fitting.	Despite getting some insights, we believe more	595
		research on the API side needs to be done. One	596
547	Task-6: On interpreting tool results. Overall,	question already inspired by our observations is: in	597
548	these agents are good at interpreting tool results.	the era of LLMs, should we design new standards	598
549	However, we still noticed a few typical errors by	for API design and documentation? It is a complex	599
550	them. In some cases, the LLMs fail to locate the de-	problem and deserves dedicated research. We treat	600
551	sired information, to some extent because of poor	it as future work.	601
552	readability of results. In addition, we find LLMs		
553	have shortcomings in a kind of copy&paste capa-	5 Conclusion	602
554	bility of target information. This makes some infor-	The LLM community is driving towards an ambi-	603
555	mation that is sensitive to character-level precision	tious vision: building a new AI ecosystem in which	604
556	(e.g. URLs, longitude and latitude, long decimal	LLM-based agents serve users by manipulating	605
557	values, etc.) not useful anymore. Furthermore, we	millions of APIs. However, an evaluation mech-	606
558	find some APIs, e.g. searching APIs, return very	anism for such agents is still missing, preventing	607
559	long results exceeding the max context length of	studies from proceeding in this area. In this work,	608
560	LLMs.	we narrow this gap by proposing a new evaluation	609
		mechanism for generic LLM-based agents. We de-	610
561	Task-7&8: On planning capability. We exam-	signed a thorough evaluation mechanism schema	611
562	ine the planning ability of agents with two folds of	aiming to examine seven different aspects of tool	612
563	experiments. In the first fold of experiments, we	use capability. Also, we release one dataset and	613
564	examine whether an agent can complete a low-level	one toolset, both designed to reflect real-world chal-	614
565	task by manipulating APIs with dependencies. All	lenges. These resources can support the studies	615
566	the evaluated agents have poor performance – even	on improving LLM-based agents as well as a new	616
567	GPT-4 has a success rate lower than 60%. We find	philosophy of API design. We evaluated five influ-	617
568	the Qwen models, even Qwen-72b, rarely have the	ential LLMs and shared findings and insights into	618
569	sense of starting with more basic APIs. It is very	their tool use capability. The found weaknesses of	619
570	challenging for the agents to manipulate APIs with	LLMs can indicate the future directions to go.	620
571	dependencies.		
572	In the second fold of experiments, we check	6 Limitations	621
573	whether an agent can outline a solid plan involving	When designing evaluation tasks, we did not in-	622
574	tool use for a high-level task. We find these LLMs’	clude data involving multi-turn interactions with	623
575	performance in decomposing a high-level task is	users. Can the agents still formulate correct func-	624
576	always decent. Though both settings require the	tion calls by parsing information from multi-turn	625
577	planning capability of LLMs, they impose very dif-		

626 dialogue? What the performance will be like if
627 the agents need to continuously request clarifica-
628 tion from the users? These problems cannot be
629 answered by our evaluation mechanism. Our eval-
630 uation mechanism does not provide an end-to-end
631 performance or an overall performance score.

632 In addition, in the landscape of generic agent
633 research, API retriever is critical. We assume the
634 existence of a good third-party API retriever. More
635 studies dedicated to API retrievers are suggested to
636 be done so that we can be closer to estimating the
637 overall performance of LLM-based agent systems.

638 7 Disclaimer

639 The toolset released by us is only for research pur-
640 poses. It is not for the usage of solving real-life
641 problems (e.g. checking flight prices).

642 References

643 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,
644 Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei
645 Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin,
646 Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu,
647 Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren,
648 Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong
649 Tu, Peng Wang, Shijie Wang, Wei Wang, Sheng-
650 guang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang,
651 Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu,
652 Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingx-
653 uan Zhang, Yichang Zhang, Zhenru Zhang, Chang
654 Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang
655 Zhu. 2023. Qwen technical report. *arXiv preprint*
656 *arXiv:2309.16609*.

657 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan,
658 Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter
659 Lee, Yin Tat Lee, Yuanzhi Li, Scott M. Lundberg,
660 Harsha Nori, Hamid Palangi, Marco Túlio Ribeiro,
661 and Yi Zhang. 2023. *Sparks of artificial general*
662 *intelligence: Early experiments with GPT-4*. *CoRR*,
663 abs/2303.12712.

664 Zehui Chen, Weihua Du, Wenwei Zhang, Kuikun
665 Liu, Jiangning Liu, Miao Zheng, Jingming Zhuo,
666 Songyang Zhang, Dahua Lin, Kai Chen, and Feng
667 Zhao. 2023. *T-eval: Evaluating the tool utilization*
668 *capability step by step*. *CoRR*, abs/2312.14033.

669 Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen,
670 Samuel Stevens, Boshi Wang, Huan Sun, and Yu Su.
671 2023. *Mind2web: Towards a generalist agent for the*
672 *web*. *CoRR*, abs/2306.06070.

673 Wenyi Hong, Weihang Wang, Qingsong Lv, Jiazheng
674 Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan
675 Wang, Yuxuan Zhang, Juanzi Li, Bin Xu, Yuxiao
676 Dong, Ming Ding, and Jie Tang. 2023. *Cogagent:*
677 *A visual language model for GUI agents*. *CoRR*,
678 abs/2312.08914.

Brian Ichter, Anthony Brohan, Yevgen Chebotar, 679
Chelsea Finn, Karol Hausman, Alexander Herzog, 680
Daniel Ho, Julian Ibarz, Alex Irpan, Eric Jang, Ryan 681
Julian, Dmitry Kalashnikov, Sergey Levine, Yao Lu, 682
Carolina Parada, Kanishka Rao, Pierre Sermanet, 683
Alexander Toshev, Vincent Vanhoucke, Fei Xia, 684
Ted Xiao, Peng Xu, Mengyuan Yan, Noah Brown, 685
Michael Ahn, Omar Cortes, Nicolas Sievers, Clayton 686
Tan, Sichun Xu, Diego Reyes, Jarek Rettinghouse, 687
Jornell Quiambao, Peter Pastor, Linda Luu, Kuang- 688
Huei Lee, Yuheng Kuang, Sally Jesmonth, Nikhil J. 689
Joshi, Kyle Jeffrey, Rosario Jauregui Ruano, Jasmine 690
Hsu, Keerthana Gopalakrishnan, Byron David, Andy 691
Zeng, and Chuyuan Kelly Fu. 2022. *Do as I can, not*
692 *as I say: Grounding language in robotic affordances*.
693 In *Conference on Robot Learning, CoRL 2022, 14-18*
694 *December 2022, Auckland, New Zealand*, volume
695 205 of *Proceedings of Machine Learning Research*,
696 pages 287–318. PMLR. 697

Chenliang Li, He Chen, Ming Yan, Weizhou Shen, 698
Haiyang Xu, Zhikai Wu, Zhicheng Zhang, Wen- 699
meng Zhou, Yingda Chen, Chen Cheng, Hongzhu 700
Shi, Ji Zhang, Fei Huang, and Jingren Zhou. 2023a. 701
Modelscope-agent: Building your customizable
702 *agent system with open-source large language mod-*
703 *els*. In *Proceedings of the 2023 Conference on Em-*
704 *pirical Methods in Natural Language Processing,*
705 *EMNLP 2023 - System Demonstrations, Singapore,*
706 *December 6-10, 2023*, pages 566–578. Association
707 for Computational Linguistics. 708

Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, 709
Pengjun Xie, and Meishan Zhang. 2023b. Towards 710
general text embeddings with multi-stage contrastive 711
learning. *arXiv preprint arXiv:2308.03281*. 712

Yaobo Liang, Chenfei Wu, Ting Song, Wenshan Wu, 713
Yan Xia, Yu Liu, Yang Ou, Shuai Lu, Lei Ji, 714
Shaoguang Mao, Yun Wang, Linjun Shou, Ming 715
Gong, and Nan Duan. 2023. *Taskmatrix.ai: Com-*
716 *pleting tasks by connecting foundation models with*
717 *millions of apis*. *CoRR*, abs/2303.16434. 718

Chang Ma, Junlei Zhang, Zhihao Zhu, Cheng Yang, 719
Yuju Yang, Yaohui Jin, Zhenzhong Lan, Lingpeng 720
Kong, and Junxian He. 2024. *Agentboard: An an-*
721 *alytical evaluation board of multi-turn LLM agents*.
722 *CoRR*, abs/2401.13178. 723

Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, 724
Long Ouyang, Christina Kim, Christopher Hesse, 725
Shantanu Jain, Vineet Kosaraju, William Saunders, 726
Xu Jiang, Karl Cobbe, Tyna Eloundou, Gretchen 727
Krueger, Kevin Button, Matthew Knight, Benjamin 728
Chess, and John Schulman. 2021. *Webgpt: Browser-*
729 *assisted question-answering with human feedback*.
730 *CoRR*, abs/2112.09332. 731

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, 732
Carroll L. Wainwright, Pamela Mishkin, Chong 733
Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, 734
John Schulman, Jacob Hilton, Fraser Kelton, Luke 735
Miller, Maddie Simens, Amanda Askell, Peter Welin- 736
der, Paul F. Christiano, Jan Leike, and Ryan Lowe. 737

738	2022. Training language models to follow instructions with human feedback . In <i>Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022</i> .	793
739		794
740		795
741		796
742		797
743		798
744	Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. 2023. Gorilla: Large language model connected with massive apis . <i>CoRR</i> , abs/2305.15334.	799
745		800
746		801
747		802
748	Chen Qian, Xin Cong, Cheng Yang, Weize Chen, Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong Sun. 2023. Communicative agents for software development . <i>CoRR</i> , abs/2307.07924.	803
749		804
750		805
751		
752	Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2023. Toollm: Facilitating large language models to master 16000+ real-world apis . <i>CoRR</i> , abs/2307.16789.	806
753		807
754		808
755		809
756		810
757		811
758		812
759	Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy P. Lillicrap. 2023. Android in the wild: A large-scale dataset for android device control . <i>CoRR</i> , abs/2307.10088.	813
760		814
761		
762		
763	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 . <i>CoRR</i> , abs/2302.04761.	815
764		816
765		817
766		818
767		819
768	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 huggingface . <i>CoRR</i> , abs/2303.17580.	820
769		821
770		822
771		823
772	Yifan Song, Weimin Xiong, Dawei Zhu, Cheng Li, Ke Wang, Ye Tian, and Sujian Li. 2023. Restgpt: Connecting large language models with real-world applications via restful apis . <i>CoRR</i> , abs/2306.06624.	824
773		825
774		
775		
776	Yashar Talebirad and Amirhossein Nadiri. 2023. Multi-agent collaboration: Harnessing the power of intelligent LLM agents . <i>CoRR</i> , abs/2306.03314.	826
777		
778		
779	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 . <i>CoRR</i> , abs/2306.05301.	827
780		828
781		829
782		830
783	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura,	831
784		832
785		833
786		834
787		835
788		836
789		
790		
791		
792		
	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models . <i>CoRR</i> , abs/2307.09288.	837
		838
	Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2023. Voyager: An open-ended embodied agent with large language models . <i>CoRR</i> , abs/2305.16291.	839
		840
		841
		842
		843
		844
		845
	Qiantong Xu, Fenglu Hong, Bo Li, Changran Hu, Zhengyu Chen, and Jian Zhang. 2023a. On the tool manipulation capability of open-source large language models . <i>CoRR</i> , abs/2305.16504.	
	Yuzhuang Xu, Shuo Wang, Peng Li, Fuwen Luo, Xiaolong Wang, Weidong Liu, and Yang Liu. 2023b. Exploring large language models for communication games: An empirical study on werewolf . <i>CoRR</i> , abs/2309.04658.	
	Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022. Webshop: Towards scalable real-world web interaction with grounded language agents . In <i>Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022</i> .	
	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan Cao. 2023. React: Synergizing reasoning and acting in language models . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	
	Chi Zhang, Zhao Yang, Jiakuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. 2023. Appagent: Multimodal agents as smartphone users . <i>CoRR</i> , abs/2312.13771.	
	Xuanhe Zhou, Guoliang Li, and Zhiyuan Liu. 2023. LLM as DBA . <i>CoRR</i> , abs/2308.05481.	
	Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyu Yang, Gao Huang, Bin Li, Lewei Lu, Xiaogang Wang, Yu Qiao, Zhaoxiang Zhang, and Jifeng Dai. 2023. Ghost in the minecraft: Generally capable agents for open-world environments via large language models with text-based knowledge and memory . <i>CoRR</i> , abs/2305.17144.	

846 **A An Example Illustrating Dependency** 847 **between APIs**

848 Here, we show the specifications of three
849 APIs from RapidAPI - skyscanner80⁸ in List-
850 ing 1,2 and 3. Only necessary information
851 for showing the dependencies between APIs
852 is kept in these doc examples. To use API
853 flights_search_one_way, we need to first check
854 API flights_auto_complete for the IDs of the
855 origin and destination. Afterwards, because
856 flights_search_one_way may not be able to re-
857 turn all the results at one time. More requests
858 to API flights_search_incomplete need to be
859 done to finalize fetching all flight data.

860 **B Details of Toolset**

861 The design of application scenarios, tool classes
862 and their functionalities can be found in Tab. 4.

863 **C Experimental settings**

864 We access GPT-4-8K and GPT-3.5-turbo via API
865 and access Qwen1.5 series LLMs via local running.

866 We only ran the experiments of evaluating
867 generic agents once. The metrics values are av-
868 eraged within the evaluation data for each task.

869 **D Performance of Agents**

870 Tab. 5 contains the performance of two GPT ver-
871 sions and three Qwen1.5 versions in 8 evaluation
872 tasks.

⁸<https://rapidapi.com/datastore/api/skyscanner80>

Listing 1: Documentation of flights_auto_complete API

```

1 {
2   "name": "flights_auto_complete",
3   "description": "This endpoint is responsible for providing a list of
4     airports for the location",
5   "parameters": {
6     "query": {
7       "type": "STRING",
8       "description": "Name of the location where the Airport is
9         situated. Ex: New York"
10    }
11  }
12 }

```

Listing 2: Documentation of flights_search_one_way API

```

1 {
2   "name": "flights_search_one_way",
3   "description": "This API helps to get the list of one-way flights.
4     Note:- In the event that the status is incomplete (data->context
5     ->status=incomplete), you must utilize the api/v1/flights/search
6     -incomplete endpoint to retrieve the complete data until it's
7     complete (data->context->status=complete).",
8   "parameters": {
9     "fromId": {
10      "type": "STRING",
11      "description": "`fromId` can be retrieved from `
12        flights_auto_complete` (data->id) Ex:
13        eyJzIjoiTl1lDQSIIsImUiOiIyNzUzNzU0MiIsImgiOiIyNzUzNzU0MiJ9 (
14        New York)"
15    },
16    "toId": {
17      "type": "STRING",
18      "description": "`toId` can be retrieved from `
19        flights_auto_complete` (data->id) Ex: eyJzIjoiTE9
20        ORCIIsImUiOiIyNzU0NDAwOCIsImgiOiIyNzU0NDAwOCJ9 (London)"
21    },
22    "departDate": {
23      "type": "Date",
24      "description": "Format: YYYY-MM-DD. Ex: 2024-06-01"
25    }
26  }
27 }
28 }

```


Listing 3: Documentation of flights_search_incomplete API

```

1 {
2   "name": "flights_search_incomplete",
3   "description": "Obtain complete data for the endpoint of
4     flights_search_one_way, flights_search_roundtrip. Until the item
5     's status is complete (data->context->status=complete), you must
6     call the API multiple times",
7   "parameters": {
8     "sessionId": {
9       "type": "STRING",
10      "description": "sessionId can be retrieved from
11        flights_search_one_way or flights_search_roundtrip (data->
12        context->sessionId)"
13    }
14  }
15 }

```

Table 4: Design of scenarios, tool classes, and API functionalities.

Scenario	Tool Class	Functionalities
daily life	weather	realtime weather, weather forecast, astronomy info
	news	news search, headlines
	calendar	public holidays, check month calendar
	recipe	search recipe
image processing	object detection	recognize objects in image
	ocr	extract text in image
	image translation	translate text in image
	image file processing	compression, format conversion, resize
	removing background	remove background
	web capture	take image screenshot
travel	flight	search one-way flights, search round-way flights, check flight details and prices
	accommodation	search hotels, check hotel details, prices, and reviews
	tourist attraction	search attractions, check details, photos and reviews of attractions
	currency	exchange rate
	airport	check airport info
	check codes	language codes, country codes,
	geocoding	convert between address and coordinates
basic & general-purpose	search	web search, image search, video search, news search
	python interpreter	python interpreter
	calculator	math calculation
	translation	translation
	ip lookup	check ip address
	access user info	user profile, location
	agent equipments	get current time

Table 5: Performance of generic agents on eight evaluation tasks.

Tasks \ Agents		GPT-4	GPT-3.5	Qwen-7b	Qwen-14b	Qwen-72b
Task 1	Precision	0.98	0.96	0.53	0.75	0.96
	Recall	0.99	0.75	0.82	0.97	0.97
	F1-score	0.98	0.84	0.65	0.85	0.96
Task 2	Accuracy	0.97	0.66	0.71	0.82	0.83
Task 3	Percentage	0.99	0.74	0.04	0.45	0.66
Task 4	Precision	0.98	0.95	0.81	0.95	0.97
Task 5	Relatedness	0.98	0.93	0.88	0.93	0.96
Task 6	Passing rate	0.93	0.85	0.85	0.81	0.85
Task 7	Progress	0.57	0.35	0.00	0.00	0.26
Task 8	Passing rate	0.99	0.98	0.94	0.95	0.98