# THEMCPCOMPANY: CREATING GENERAL-PURPOSE AGENTS WITH TASK-SPECIFIC TOOLS

**Anonymous authors** 

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031 032 033

034

037

038

040

041

042

043

044 045

046

047

048

051

052

Paper under double-blind review

# **ABSTRACT**

Since the introduction of the Model Context Protocol (MCP), the number of available tools for Large Language Models (LLMs) has increased significantly. These task-specific tool sets offer an alternative to general-purpose tools such as web browsers, while being easier to develop and maintain than GUIs. However, current general-purpose agents predominantly rely on web browsers for interacting with the environment. Here, we introduce TheMCPCompany, a benchmark for evaluating tool-calling agents on tasks that involve interacting with various real-world services. We use the REST APIs of these services to create MCP servers, which include over 18,000 tools. We also provide manually annotated ground-truth tools for each task. In our experiments, we use the ground truth tools to show the potential of tool-calling agents for both improving performance and reducing costs assuming perfect tool retrieval. Next, we explore agent performance using tool retrieval to study the real-world practicality of tool-based agents. While all models with tool retrieval perform similarly or better than browser-based agents, smaller models cannot take full advantage of the available tools through retrieval. On the other hand, GPT-5's performance with tool retrieval is very close to its performance with ground-truth tools. Overall, our work shows that the most advanced reasoning models are effective at discovering tools in simpler environments, but seriously struggle with navigating complex enterprise environments. The MCP-Company reveals that navigating tens of thousands of tools and combining them in non-trivial ways to solve complex problems is still a challenging task for current models and requires both better reasoning and better retrieval models.<sup>1</sup>

#### 1 Introduction

Since the introduction of the MCP protocol by Anthropic in November 2024 (Anthropic, 2024; FastMCP, 2025), there has been continuous explosive growth in the number of MCP servers. A June 2025 survey by Virustotal (Quintero, 2025) counted 17845 MCP server projects on GitHub. The awesome-mcp-servers list (Gizdov, 2025) contains over 7000 publicly available MCP servers. And this is just public servers; more and more, organizations are creating MCP servers to expose the functionality of internal tools to LLMs as well. This makes sense for a number of reasons. Using MCP servers, LLMs can directly call the specific tools needed for completing each task (e.g., create\_pr and merge\_pr) (Patil et al., 2023; Schick et al., 2023). MCP servers are relatively simple to create and maintain, and providing direct access to tool documentation provides a straightforward way for LLMs to interact with the environment.

Despite this proliferation of direct access to tools and API surfaces, however, general-purpose agents still predominantly rely on general-purpose tools such as web browsers and code interpreters to solve problems (Fourney et al., 2024). Here, we aim to understand the capabilities and performance of an alternative approach: general-purpose agents based on large, heterogeneous tool collections.

Although there are several prior papers studying specific aspects of tool-based agents, none of them provides a comprehensive view of the challenges that come with the combination of a large number of tools and complex tasks in a complex environment. First, there is a growing body of work on creating general-purpose AI agents. While these works represent the complexity of tasks and

<sup>&</sup>lt;sup>1</sup>We will release all our code and data after the double-blind review process.

environments that agents face in reality, they often incorporate a very small number of task-specific tools (e.g., a dedicated search tool) (Mozannar et al., 2025; Soni et al., 2025). Thus, it is unclear how AI agents behave when the number of available tools increases significantly. On the other hand, there is a rich literature that studies different challenges of tool calling with LLMs (Feng et al., 2025), such as complex function calls (Zhong et al., 2025) and large tool sets (Qin et al., 2023). However, tool calling works often rely on simple environments that are not representative of practical applications, like automating enterprise workflows. Our goal is to provide a realistic environment that includes challenging tasks, complex services, and a large and complex tool set for studying the potential and challenges of tool-based agents in practical scenarios.

We introduce TheMCPCompany, an extension of TheAgentCompany (Xu et al., 2024a) that simulates a software company where MCP tools are available for all operations in the company. In fact, this simulation represents our vision for enterprise environments in the future. To better represent complex enterprise workflows, we expand TheAgentCompany's environment by introducing the Microsoft Azure cloud platform. We then create a fully functional MCP server for each of the services (Azure, Plane, GitLab, ownCloud, and RocketChat) that exposes its full functionality through tools (more than 18,000 tools in total, of which almost 17,000 come from Azure). We adapt the existing tasks from TheAgentCompany to the MCP setting and create a new set of tasks specifically for Azure. These tasks range from relatively simple ones whose solutions can be found in a web search to complex, enterprise-level debugging (Fig. 1). Finally, we annotate a small set of required tools for each task, allowing us to evaluate tool use separately from tool selection.

We also create MCPAgent, a baseline agent that treats tool retrieval itself as a tool. MCPAgent has access to all 18k tools, but it must discover them by constructing queries and then reasoning about the results. This allows the agent to explore different solution trajectories and dynamically search for the required tools and their dependencies. We implement MCPAgent based on OpenHands' CodeAct agent (Wang et al., 2024c).

We evaluate six different LLMs on the tasks adapted from TheAgentCompany and show that task-specific tools are a practical and even preferred interface for interacting with the environment. Compared to OpenHands' CodeAct agent, which uses a text-based browser, an agent with access to the ground truth tools improves performance by 13.79 points and reduces costs by \$2.29 per task on average (54% reduction in costs). Even without the ground truth tools, our MCPAgent with the tool-finder function outperforms the alternative browser-based agent by 5.39 points and reduces costs by \$2.06 per task on average. On these tasks, GPT-5 performs almost as well with the tool finder as with ground-truth tools.

In contrast, on our hardest tasks in the Azure environment, even the most capable reasoning models fail almost completely. We find that agents mainly struggle with the diversity and complexity of Azure services. For example, they fail to correctly identify the issue with a broken application, do not consider all possible solutions when one fails, and often implement only part of the solution.

Our results show that agents can solve problems in enterprise environments that are more complex and contain far more tools than previously considered in the literature. They also show that MCP is a key facilitator: exposing tools to LLMs via a standardized protocol leads to better results than relying on browser-based agents. However, our results also reveal a key challenge going forward in this space: navigating thousands or more tools that must be combined in non-obvious ways to solve complex problems is both a retrieval and a reasoning problem. The most advanced reasoning models are capable of searching for tools, but more work is needed on both fronts to fully realize our vision for future enterprise environments. TheMCPCompany supports this work by inviting future contributions to explore more realistic and complex scenarios that agents face in practice.

## 2 RELATED WORK

AI Agents There is a growing body of work on AI agents (Handa et al., 2025; Shao et al., 2024; 2025; Xie et al., 2024). Although most of the first generation of agents are domain-specific, such as coding (Wang et al., 2024c; Xia et al., 2024; Yang et al., 2024) or browsing agents (Chezelles et al., 2024), more recently there has been a push toward general-purpose agents that can complete diverse tasks across multiple domains (Hu et al., 2025; Wu et al., 2023). Since a general-purpose agent needs to interact with different services depending on the given task, current agent frameworks

 predominantly interact with the environment via general-purpose tools such as a browser, shell, or Python interpreter (Soni et al., 2025). Recently, Song et al. (2024) proposed using REST API calls instead of browser interactions. However, compared to REST APIs, MCP tools are easier to create and are being actively developed by the machine learning community, and thus better suited for use with LLMs. Moreover, Song et al. (2024) use a small number of tools (less than a thousand) for each task and provide a short description of all tools in the prompt, which does not scale to large tool sets capable of performing in practical scenarios. For these cases, retrieval is necessary.

Agent benchmarks have also evolved in different directions. For example, there are many benchmarks that aim to create complex tasks (Mialon et al., 2023), simulate realistic environments (Xu et al., 2024a), or study the impact of agents on the workforce (Styles et al., 2024). However, similar to agent frameworks, these benchmarks are either limited to a small set of tools (Wang et al., 2024a; Yao et al., 2024) or mainly rely on the browser (Zhou et al., 2023) for agent interactions.

As a result, the challenges and opportunities for agents that primarily rely on large tool sets to interact with the environment are largely unknown. Here, we build on prior work (Xu et al., 2024a) and maintain the complexity and realism of the tasks and environment. However, we replace the few general-purpose tools with a large number of task-specific tools and investigate the challenges and opportunities that agents face in this new setup.

**Tool Use** The ability to call tools to interact with the environment is what makes the current generation of AI agents feasible. There is an extensive body of research studying various aspects of tool calling with LLMs (Chen et al., 2025; Liu et al., 2024b; Yuan et al., 2023), ranging from the complexity of tool calls (Zhong et al., 2025) to dependency between tools (Lumer et al., 2025). However, most works rely on a small set of tools and do not represent the growing scale of MCP tools available to LLMs (Dong et al., 2025; Feng et al., 2025; Li et al., 2025; Wang et al., 2025). While there are several works that investigate large tool sets (Qin et al., 2023), their environments are simple compared to what agent benchmarks provide (Fei et al., 2025; Gan & Sun, 2025; Liu et al., 2024a; Mo et al., 2025; Shi et al., 2025; Xu et al., 2024b). The tasks are also simple and often there is significant semantic overlap between the task description and tool specifications, which simplifies tool selection (Li et al., 2023). However, in practice, task descriptions (e.g., fix a broken app) often do not mirror the name and description of the required tools (e.g., list\_managed\_identities).

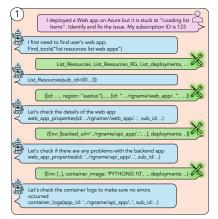
With the increasing popularity of MCP, there is a renewed interest in tool calling benchmarks but through MCP servers (Lei et al., 2025; Luo et al., 2025b). What sets MCP tools apart from traditional tool calling is the opportunity for massively scaling the number of tools by standardizing the communication protocol. However, current MCP benchmarks are generally limited to between a few hundred and a thousand tools (Gao et al., 2025; Liu et al., 2025; Luo et al., 2025a; Yin et al., 2025). Moreover, the related MCP servers for each task are manually selected for the agent prior to execution which ignores the impact of tool selection as one of the main challenges that agents face when dealing with large tool sets (Luo et al., 2025b).

Unlike prior work on tool calling, we take full advantage of MCP's main strength, scalability, and create more than 18,000 functional tools for interacting with different real-world services. Also, in our setup, we do not directly provide the related tools for each task to the agent. Instead, it needs to use a tool finder function to search for and discover the required tools on its own.

# 3 THEMCPCOMPANY

Considering the simplicity of developing and maintaining MCP servers and the growing interest of the community, we argue that in the near future, MCP tools will be LLMs' primary interface for interacting with the world. In other words, there will be an MCP tool for every operation and every application (e.g., GitLab); teams in an organization will also offer MCP servers for interacting with their services. In fact, this is already happening. Many services already offer MCP servers, and there are numerous efforts to further simplify widespread adoption of MCP. For example, Docker Desktop offers a dedicated toolkit that simplifies the deployment of MCP servers (Docker, 2025), and there is even a registry for keeping track of the growing number of MCP servers<sup>2</sup>.

<sup>&</sup>lt;sup>2</sup>gh/modelcontextprotocol/registry



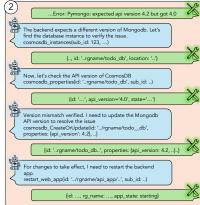


Figure 1: The correct solution trajectory for one of our Composite Azure tasks. Note that agents must use the tool finder function to discover each one of the tools used in the trajectory. But, due to space constraints, here, we only show the first call to find-tools.

Here, we describe TheMCPCompany, an extension of TheAgentCompany benchmark that simulates a realistic and tool-rich environment. We first include the Microsoft Azure cloud platform, which significantly increases the complexity of the environment and the number of available actions. Then, we create MCP servers for all services in TheMCPCompany that expose the full functionality of each service through a large collection of tools.

### 3.1 THEAGENTCOMPANY (BACKGROUND)

TheAgentCompany is a benchmark that simulates a small-scale software company for evaluating agents' ability to complete everyday enterprise tasks (Xu et al., 2024a). TheAgentCompany offers self-hosted services for project management (Plane), DevOps (GitLab), communication (RocketChat), and productivity (ownCloud) as docker images with pre-populated data. It also configures LLM-powered non-player characters that act as employees in the company. This provides a realistic environment, where the agent needs to interact with multiple services and simulated employees to successfully complete a task. For each task, TheAgentCompany provides an evaluation script, with multiple checkpoints, that assigns partial credit when the agent only completes part of the task.

We choose TheAgentCompany as the basis for our work since its environment approximates real-world applications more closely than other benchmarks. More importantly, the action space provided by the four hosted services is considerably larger than that of other agent benchmarks, which facilitates our goal of creating an environment with a large tool set for agent evaluation.

#### 3.2 AZURE TASKS

While TheAgentCompany uses real-world applications for simulating the environment, these self-hosted applications are simpler than many services used in production by most organizations. For example, most software companies use cloud computing platforms, such as Azure and AWS, as part of their workflow. These services are so complex that employees take dedicated courses just to be able to manage the infrastructure. Therefore, to have a realistic view of LLMs' potential for practical applications, we require an environment where agents directly work with services used in production, instead of simpler proxies commonly used for evaluation. To achieve this, we have created two small sets of tasks that require managing resources in the Microsoft Azure cloud platform. Our tasks exercise a range of activities that require interacting with different parts of Azure, including resource management, security, storage, compute and Cognitive Services like image recognition.

In the first category, we have created 10 *Primitive* tasks, where the agent only needs to take a very specific action on a very specific resource. Examples of primitive tasks are adding tags to a given resource or deleting a specific resource. These tasks mainly measure agents' ability to identify the correct tool for a given action from the large pool of Azure MCP tools and generate the correct tool call. For the second category, we have created seven *Composite* tasks that are intended to reproduce more challenging real-world scenarios that an Azure user would normally have to carry

out. The composite tasks involve an infrastructure with multiple services (e.g., CosmosDB, Key vault, Function app) that are configured for a specific application, like serving a TODO list web app. In this category, the agent is given higher-level goals, such as fixing a broken app, implementing a security policy, or adding a new feature. To successfully complete the composite tasks, the final state of the environment must meet the requirements of the task in addition to having a working application. The composite tasks are more difficult and measure the agent's ability to understand and navigate the complex logic of the Azure environment, such as coordinating code edits and environment configuration and understanding the space of possible solutions for a given problem.

**Task Details** To make evaluations more accessible, our tasks focus on the cheapest Azure resources and can be run practically for free using a free-tier Azure subscription (free Azure subscriptions come with a \$200 credit. During the development and troubleshooting of the tasks, which involved executing each task many times, we spent less than \$1 of this limit). For each task, we provide a task description, an evaluation script to judge whether the task was completed successfully, and a proof-of-concept script that solves the task using the available MCP tools. Moreover, to have a reproducible environment, we provide a Terraform<sup>3</sup> script for each task that initializes and tears down the execution environment on Azure.

#### 3.3 A LARGE AND COMPREHENSIVE TOOL SET

To provide an environment where the agent primarily relies on task-specific tools for interacting with other services, we create a large collection of tools that collectively expose the full functionality of each of the services in the environment. For example, we create dedicated tools for merging a PR on GitLab or listing the resources in an Azure subscription.

Most modern services come with comprehensive REST APIs that offer a dedicated endpoint for each individual operation (e.g. list available users). While prior work has proposed agents that directly call the REST APIs (Song et al., 2024), we argue that MCP tools are a more appropriate solution for large-scale adoption in long term. Thanks to libraries like FastMCP (FastMCP, 2025), MCP tools are easier to develop and maintain compared to REST APIs. More importantly, MCP tools are LLM friendly: each tool is accompanied by the description of its functionality and arguments, and MCP provides an easy and standard method for accessing these documentations. This allows LLMs to discover the required tools for each task and also learn how to use new tools on the fly. On the other hand, there is no standard method for providing the REST API documentations to LLMs. Therefore, we convert the REST APIs of Azure, GitLab, and RocketChat into dedicated MCP servers that provide a corresponding tool for each API endpoint. We also extract the description for each tool and its arguments from the API specifications provided by each service. For RocketChat, we use an LLM with access to RocketChat's online documentation to write more informative tool descriptions. See Appendix B for details.

Plane and ownCloud do not provide comprehensive REST API support. To overcome this, we treat own-Cloud as a file server and manually create an MCP server that provides basic file operations (e.g., download and upload). We observe that these file operations are sufficient for completing TheAgentCompany tasks, and the agent often uses Python libraries to manipulate the spreadsheet or presentation files on ownCloud. Finally, we adopt the official MCP server for Plane and manually add any missing tools that are required for completing the tasks. After creating the MCP servers, we manually go through all the tasks and make sure they are feasible with the available tools.

Service	#MCP Tools	Avg #Args	Complex Tools (%)
Plane	52	2.06	28.85
RocketChat	520	2.82	12.31
ownCloud	11	1.64	0.00
GitLab	1,085	5.47	10.69
Azure	16,837	5.63	22.50
Total	18,505	5.53	21.52

Table 1: The number and properties of tools provided by TheMCPCompany.

Furthermore, for each task, we manually annotate a small set of tools that are sufficient for the successful completion of the task. In our later analysis, we use these annotated tools to isolate the impact of tool selection and measure the upper bound on the performance of tool-based agents with

<sup>&</sup>lt;sup>3</sup>hashicorp/terraform

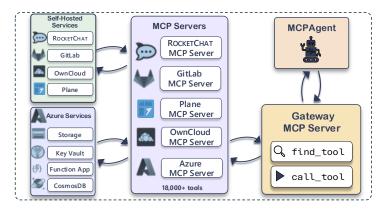


Figure 2: Our MCP servers expose the full functionality of each service through tools. Instead of directly providing the 18,000 tools to the agent, we provide it with a gateway MCP server with two tools, which the agent can use to search for and invoke the required tools at each step.

current models. We also update the task descriptions and evaluation scripts in TheAgentCompany, which are written for browser-use agents, to be compatible with tool-based agents.

**Tool Characteristics** In addition to providing a large number of tools, TheMCPCompany's tool set also represents the complexity of tool calls in practice (Table 1). On average, our tools accept more than five arguments and, in some cases, the agent has to provide up to 39 arguments for some tool invocations for Azure. For example, to create a virtual machine, the agent should provide detailed information about all the dependent resources (e.g., disk, network interface, virtual networks, OS image, role assignments, etc.). There is also a significant dependency between our tools. For instance, the agent has to first create all the dependent resources in order to be able to successfully call the tool for creating a virtual machine. Moreover, many of TheMCPCompany's tools require passing arguments with complex data types. Specially, for Azure and Plane, 22.5% and 28.85% of tools have at least one argument of type array or object. For Azure, this is more than 3K complex functions, and in our experience, most of the tools that change the environment state (e.g., create or modify resources) require deeply nested arguments.

Moreover, our tool set represents the chaotic nature of real-world applications. For example, there are similar tools with totally different purposes (e.g., send\_msg\_to\_room, send\_msg\_to\_individual). On the opposite side, often there are several tools for each action, with slight differences (e.g., gitlab\_search\_all, gitlab\_search\_issues). Similarly, there are different sequences of tool calls for accomplishing a goal, with some more efficient than others.

# 4 MCPAGENT

We also create a baseline agent to study the feasibility of tool-based agents with a large tool set (Fig. 2). Utilizing the extremely large number of tools is the main challenge for creating practical tool-based agents. Naive solutions are untenable; the context window of current LLMs does not fit the specification for all the tools in our benchmark (18,000+). To address this issue, prior work uses retrieval models to select the necessary tools based on the task description (Qin et al., 2023). However, for realistic and challenging tasks, such as those in TheAgentCompany, the task description often has little in common semantically with the description of the required tools. For example, while role assignment is necessary for managing storage accounts in Azure, there is little to no semantic similarity between tools related to role assignment and the description of a task for backing up a storage account.

Instead of selecting the tools prior to execution, we allow the agent to select the tools itself. Specifically, we create a gateway MCP server with a tool finder function that the LLM can use to search for required tools at each step using a text query. Under the hood, the tool finder uses a text embedding model to encode the JSON specification of the tools and also the agent's query. Then, based on the cosine similarity between query and tool embeddings, it returns the specification for the top-k tools. Since the LLM does not have direct access to the main tools, the gateway MCP server provides

Model	Browser				MCPAgent				Oracle Tool Set			
Wiodei	Score	Success (%)	Steps	Cost (\$)	Score	Success (%)	Steps	Cost (\$)	Score	Success (%)	Steps	Cost (\$)
Sonnet 4	45.06	34.86	31.16	5.02	48.79	39.43	30.82	2.75	56.36	47.43	26.97	2.13
Opus 4.1	41.16	31.43	24.07	14.58	48.68	39.43	22.53	7.29	57.26	48.00	23.65	7.17
GPT 4.1	31.71	22.99	22.71	1.72	37.10	27.43	20.48	0.75	46.76	36.00	16.05	0.56
03	30.53	22.86	21.92	1.17	45.39	37.14	23.41	0.83	50.63	40.57	22.53	0.65
GPT-5-mini GPT 5	33.36 50.24	24.57 40.00	31.74 28.75	0.41 2.20	32.11 52.32	22.86 42.29	29.27 19.39	0.26 0.85	49.33	38.86 44.57	22.33 17.54	0.17 0.66

Table 2: The performance of different LLMs on the 175 tasks adapted from TheAgentCompany. Browser: the LLM uses the browser for completing tasks. MCPAgent: the LLM uses the tool finder function to discover and invoke the required tools. Oracle Tool Set: the LLM is provided with the required tools for each task.

another function that takes the name and arguments of any of the retrieved tools, calls the tool for the LLM, and returns the results. This architecture keeps the number of tools manageable for the agent, and at the same time, it provides more flexibility by allowing the LLM to explore different solutions and choose the required tools dynamically. Moreover, it also provides a unified interface to a heterogeneous set of tools. Finally, except for the browser tool, our agent has access to all the standard tools in OpenHands' CodeAct agent (Wang et al., 2024b;c) (Think, Python, Shell, Web fetch, and File edit), which are necessary for completing TheAgentCompany tasks.

# 5 EXPERIMENTS

#### 5.1 SETUP

We build our agent based on OpenHands' CodeAct agent (Wang et al., 2024c), with a slightly modified system prompt that instructs the LLM to use tools instead of the browser. See Appendix A for details. We then evaluate GPT-4.1, o3, GPT-5-mini, GPT-5, Sonnet-4, and Opus-4.1 on TheAgent-Company and Azure tasks (Anthropic, 2025; OpenAI, 2025). We use OpenAI's text-embedding-3-large model to calculate the embeddings for the tool finder function (OpenAI, 2025). Unfortunately, because of incompatibility with OpenHands, we disable the thinking blocks for Opus-4.1. In our experiments with TheAgentCompany tasks, we use an earlier version of our MCP servers, with about 13,000 tools. However, it does not substantially impact our experiments since these tools are not needed for TheAgentCompany tasks.

**Evaluation** For The Agent Company tasks, we use the same evaluation metrics as Xu et al. (2024a). The score for each task consists of two parts. The obtained credit from evaluation checkpoints accounts for 50% of the final score. The other 50% is only assigned if the agent completes the task successfully. We also report the percentage of tasks completed successfully and the average steps and inference costs for each task. The inference costs are calculated based on the token usage for each task and prices published by LLM providers. Since there are many valid solution trajectories for Azure tasks, we only consider the successful completion for evaluation without partial credits.

#### 5.2 THEAGENTCOMPANY TASKS

**Potential of Task-specific Tools** First, we consider the question of whether task-specific tools are an appropriate interface for interacting with the environment. We directly provide the small oracle tool set to the agent for each task, excluding the impact of tool retrieval on performance. Compared to OpenHands' default CodeAct agent, which uses a text-based browser, using task-specific tools increases performance by 13.79 points on average across different models, with more than 20 points for o3 (Columns Browser and Oracle Tool Set in Table 2). Except for GPT-5 which has good performance in both cases, we observe that the reasoning models, Opus-4.1 and o3, benefit more from task-specific tools than do their non-reasoning counterparts (Sonnet4 and GPT-4.1).

While with a browser, the agent needs to navigate the web interface and process the entire content of each web page, task-specific tools allow the agent to take the necessary action directly and only process the required information, which reduces inference costs. Across different models, the agent with the oracle tool set reduces inference costs by \$2.29 on average per task compared to the

browser-based agent, with up to \$7.41 reduction in average costs per task for Opus-4.1. Moreover, for all models except for Opus-4.1 and o3, the number of required steps for each task also decreases, which directly translates to latency and usability of the resulting agents. The combination of better performance and reduced costs positions large sets of task-specific tools as a promising approach for developing general-purpose agents.

**Task-specific Tools in Practice** In real-world applications, we do not have access to the oracle tool set. To investigate the feasibility of creating general-purpose agents with task-specific tools in practice, we evaluate MCPAgent, which uses tool retrieval to discover the necessary tools for each task (Table 2). We find that even without the oracle tool set, using task-specific tools is preferred over the browser. Compared to the browser-based agent, MCPAgent improves performance by 5.39 points on average across all models, with a maximum improvement of 14.86 points for o3. Interestingly, the increases in performance are consistently larger for reasoning models compared to their non-reasoning counterparts.

Without the oracle tool set, the LLMs cannot take full advantage of the task-specific tools, and their performance is, on average, 8.4 points behind the agent with access to ground truth tools. We believe this gap would decrease in the future as the capabilities of LLMs improve. In fact, GPT-5 already closes the gap, and its performance without the oracle tool set only decreases by 2.13 points. However, this is the exact opposite for smaller and more affordable models like GPT-5-mini. In fact, the performance of GPT-5-mini without the oracle tool set is worse than its performance with the browser tool.

Interestingly, despite the additional calls to the tool finder function, MCPAgent provides similar cost savings to the agent with access to oracle tool set. Compared to OpenHands' CodeAct agent, MC-PAgent reduces inference costs by \$2.06 on average per task across all models. Our results show that even with current models, creating general-purpose agents with task-specific tools instead of a few general-purpose tools is practical and also provides significant benefits. These findings encourage future work to explore more effective agentic solutions for taking advantage of the growing number of task-specific tools available to LLMs.

#### 5.3 AZURE TASKS

Given the large action space of the Azure environment, we first use our Primitive Azure tasks to evaluate if LLMs can correctly find and invoke the correct tool to achieve a very specific and clear goal, such as deleting a virtual machine (Table 3). We find that GPT-5, Sonnet-4, and Opus-4.1 use the tool finder function effectively and achieve nearly perfect scores on our Azure tasks. However, GPT-4.1, o3, and GPT-5-mini struggle even with these simple tasks. Also, surprisingly, despite clear instructions to use tools, GPT-4.1 and o3 often insist on using command line tools, and after they fail, they just provide a high-level outline of the solution and give up.

Model	Primitive	Composite
Sonnet 4	9/10	1/7
Opus 4.1	9/10	1/7
GPT 4.1	5/10	0/7
o3	6/10	1/7
GPT-5-mini	2/10	0/7
GPT 5	9/10	1/7

Table 3: The number of successfully completed Azure tasks in each category using MCPAgent with different LLMs.

Evaluation on our Composite tasks shows that LLMs' problem-solving capabilities diminish when faced with

complex tasks in a complex environment, and all models consistently fail on almost all these tasks. We find that after failure, models do not explore alternative solutions. For instance, if the model does not have enough quota to deploy an Azure function, it does not try a different region or deploy the app on other resources like a container. Moreover, models do not follow a systematic approach for diagnosing and resolving problems. Instead, they focus on the most common cause for a given problem, often Identity and Access Management (IAM), and do not even check if their solution resulted in a functioning infrastructure.

## 5.4 TOOL CALLING PATTERNS

**TheAgentCompany Tasks** Table 4 reports the tool-use statistics of each model for TheAgent-Company tasks. LLMs effectively use the tool finder function and find the required tools after

retrieving only about 20 tools, which is well below the maximum number of tools allowed by inference APIs (often 128). Also, solving each task requires only a handful of calls to task-specific tools, which explains the reduced inference costs of tool-based agents.

We find that reasoning models are better suited for use with a large number of task-specific tools. First, reasoning models call the MCP tools more accurately and fail less often than non-reasoning models. Similarly, reasoning models use tool retrieval more effectively and consistently achieve better retrieval recall. Finally, among the models that we tested, GPT-5 generates the most comprehensive and longest queries, which could explain its superior performance with MCPAgent.

**Azure Tasks** Table 5 in the Appendix reports these statistics for Azure tasks, with similar patterns. One interesting observation is that the complexity of the tasks is also reflected in models' tool calling patterns. Except for GPT-4.1, o3, and GPT-5-mini that often fall back to command line tools and fail, other models consistently retrieve and call more tools for Composite tasks than Primary tasks. Also, calling the correct tools with correct requirements and arguments is more challenging for Composite tasks and consequently, the agent's tool calls fail more often. For composite tasks, identifying a solution and retrieving the required tools is

Model	#Retrieved	#MCP	Failed	Retrieval	Query
	Tools	Calls	Calls (%)	Recall	Length
Sonnet 4	15.7	9.9	10.7	60.0	34.5
Opus 4	25.8	7.3	8.5	69.7	32.6
GPT 4.1	13.5	9.1	29.7	44.9	31.6
o3	22.2	7.8	13.0	53.1	19.2
GPT-5-mini	20.2	8.1	22.2	32.8	44.6
GPT 5	15.3	11.5	8.3	58.7	52.9

Table 4: MCPAgent's tool calling statistics on the 175 adapted tasks from TheAgentCompany. Query length is measured in the number of characters.

also difficult, and agents use the tool finder function more often and with longer queries.

#### 5.5 ERROR ANALYSIS

To better understand the failure modes of tool-based agents, for GPT-5 experiments with retrieval and oracle tool set, we inspect the trajectories of 10 tasks where the agent receives zero points. We find that in retrieval mode, if the model does not find the required tools after a few attempts, it formulates an alternative solution even if it does not meet the task requirements. Therefore, investigating better tool retrieval methods is an important research direction for developing more capable tool-based agents. We also notice that for lengthy tasks with many steps or complicated requirements, the model often only completes part of the task before prematurely declaring victory. We see this pattern clearly both with the more difficult TheAgentCompany tasks and with the composite Azure tasks. Interestingly, during the course of this project, we noticed GPT-5's excellent performance is in part due to its perseverance. However, this can also cause it to eventually exceed its context window for long-horizon tasks.

# 6 CONCLUSION

In this work, we introduce TheMCPCompany, a benchmark for general-purpose agents that primarily use task-specific tools for interacting with the environment. We provide MCP servers with a large number of tools (more than 18,000) that expose the full functionality of several real-world services. Our tool set is created from existing REST APIs and thus closely simulates tool calling in the real world. In addition, we include Microsoft's Azure cloud platform in our environment and provide the necessary tools for all possible interactions with Azure, which significantly increases the environment's complexity. Through extensive experiments, we show the significant potential of task-specific tools for improving performance and reducing costs compared to browser-based agents. We also use tool retrieval to create a practical agent that automatically discovers the necessary tools for each task. We find that, even with imperfect retrieval, using task-specific tools still improves performance and reduces inference costs. Our results encourage future work to explore task-specific tools as an alternative approach for creating general-purpose agents. Also, the integration of Azure in our environment provides a valuable opportunity for future work to create more challenging tasks and further explore the agents' behavior in a real enterprise environment.

# LIMITATIONS

Unintended Consequences of Deploying LLM Agents in Practice While providing the full functionality of production services, like Azure, to LLM agents opens a whole new category of tasks that LLMs can accomplish, it also increases the risks. Without any restrictions, deploying LLM agents in practice comes with many risks, such as destroying critical resources, incurring unnecessary costs (e.g., deploying expensive services), or exposing sensitive information to unauthorized users. For example, in our Azure tasks, GPT-5 mistakenly deletes a virtual machine, which is an irreversible action. While our work mainly focuses on the ability of agents to complete a given task, this is not sufficient for using LLM agents in practice. In addition to improving LLMs' performance, we encourage future work to also investigate potential approaches for mitigating the side effects of LLM actions without limiting the available actions to the LLM, for example, through human-in-the-loop agentic systems (Mozannar et al., 2025). By incorporating Azure, TheMCPCompany provides a realistic environment for future work to investigate different aspects of LLM agents in practical applications.

**Number of Azure Tasks** Our Azure tasks reveal the weaknesses of LLM agents in navigating complex real-world environments. However, considering the numerous Azure services, there are many other types of problems and scenarios that are not included in our tasks. TheMCPCompany exposes the full functionality of Azure through tools. To better understand LLMs' behavior in enterprise workflows, we encourage future work to use TheMCPCompany's large tool set and investigate LLMs' behavior on other tasks and types of problems, such as multi-subscription governance, threat detection, and disaster recovery.

# ETHICS STATEMENT

Although the artifacts and methods presented in our work do not raise any immediate ethical concerns, incorporating LLM agents in actual production workflows requires extensive supervision and careful analysis, especially when interacting with user data. For example, in some of TheAgentComany tasks, the LLM is tasked to review several resumes and select the most qualified candidate. Delegating such tasks to LLM agents requires careful consideration since LLMs' biases could adversely impact parts of society (Bender et al., 2021).

# REPRODUCIBILITY STATEMENT

In our work, we use the same environment as TheAgentCompany (Xu et al., 2024a), which is based on publicly available docker images and creates the same container for all experiments. To create a reproducible environment for Azure tasks, we rely on the infrastructure-as-code paradigm. Specifically, we provide Terraform scripts for every task that create the same resources for each task every time and also destroy the resources at the end, to avoid extra costs. Moreover, we exclusively rely on the cheapest Azure services and the free credit assigned to all users, which ensures everyone can reproduce our results on Azure tasks. We use the default OpenHands (Wang et al., 2024c) parameters in our experiments and explain the exact version of OpenHands in our experiments as well as any modifications in Appendix A. Finally, to facilitate further progress in this direction, we will also release our data and code (including our MCP servers) to the public after the double-blind review process.

#### REFERENCES

- Anthropic. Introducing the model context protocol. https://www.anthropic.com/news/model-context-protocol, November 2024. Accessed: 2025-06-30.
- Anthropic. Model's overview. https://docs.claude.com/en/docs/about-claude/models/overview, September 2025. Accessed: 2025-09-23.
- Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pp. 610–623, 2021.

- Chen Chen, Xinlong Hao, Weiwen Liu, Xu Huang, Xingshan Zeng, Shuai Yu, Dexun Li, Shuai Wang, Weinan Gan, Yuefeng Huang, et al. Acebench: Who wins the match point in tool usage? arXiv preprint arXiv:2501.12851, 2025.
  - De Chezelles, Thibault Le Sellier, Sahar Omidi Shayegan, Lawrence Keunho Jang, Xing Han Lù, Ori Yoran, Dehan Kong, Frank F Xu, Siva Reddy, Quentin Cappart, et al. The browsergym ecosystem for web agent research. *arXiv* preprint arXiv:2412.05467, 2024.
  - Docker. Docker mcp catalog. https://docs.docker.com/ai/mcp-catalog-and-toolkit/catalog/, 2025. Accessed: 2025-09-23.
  - Guanting Dong, Yifei Chen, Xiaoxi Li, Jiajie Jin, Hongjin Qian, Yutao Zhu, Hangyu Mao, Guorui Zhou, Zhicheng Dou, and Ji-Rong Wen. Tool-star: Empowering llm-brained multi-tool reasoner via reinforcement learning. *arXiv preprint arXiv:2505.16410*, 2025.
  - FastMCP. Fastmcp: The fast, pythonic way to build mcp servers and clients. https://gofastmcp.com, September 2025. Accessed: 2025-9-18.
  - Xiang Fei, Xiawu Zheng, and Hao Feng. Mcp-zero: Proactive toolchain construction for llm agents from scratch. *arXiv preprint arXiv:2506.01056*, 2025.
  - Jiazhan Feng, Shijue Huang, Xingwei Qu, Ge Zhang, Yujia Qin, Baoquan Zhong, Chengquan Jiang, Jinxin Chi, and Wanjun Zhong. Retool: Reinforcement learning for strategic tool use in llms. *arXiv preprint arXiv:2504.11536*, 2025.
  - Adam Fourney, Gagan Bansal, Hussein Mozannar, Cheng Tan, Eduardo Salinas, Friederike Niedtner, Grace Proebsting, Griffin Bassman, Jack Gerrits, Jacob Alber, et al. Magentic-one: A generalist multi-agent system for solving complex tasks. *arXiv preprint arXiv:2411.04468*, 2024.
  - Tiantian Gan and Qiyao Sun. Rag-mcp: Mitigating prompt bloat in llm tool selection via retrieval-augmented generation. *arXiv preprint arXiv:2505.03275*, 2025.
  - Xuanqi Gao, Siyi Xie, Juan Zhai, Shqing Ma, and Chao Shen. Mcp-radar: A multi-dimensional benchmark for evaluating tool use capabilities in large language models. *arXiv preprint arXiv:2505.16700*, 2025.
  - Orislav Gizdov. awesome-mcp-servers. https://github.com/bgizdov/awesome-mcp-servers, September 2025. Accessed: 2025-09-23.
  - Kunal Handa, Alex Tamkin, Miles McCain, Saffron Huang, Esin Durmus, Sarah Heck, Jared Mueller, Jerry Hong, Stuart Ritchie, Tim Belonax, et al. Which economic tasks are performed with ai? evidence from millions of claude conversations. *arXiv preprint arXiv:2503.04761*, 2025.
  - Mengkang Hu, Yuhang Zhou, Wendong Fan, Yuzhou Nie, Bowei Xia, Tao Sun, Ziyu Ye, Zhaoxuan Jin, Yingru Li, Qiguang Chen, et al. Owl: Optimized workforce learning for general multi-agent assistance in real-world task automation. *arXiv preprint arXiv:2505.23885*, 2025.
  - Fei Lei, Yibo Yang, Wenxiu Sun, and Dahua Lin. Mcpverse: An expansive, real-world benchmark for agentic tool use. *arXiv preprint arXiv:2508.16260*, 2025.
  - Minghao Li, Yingxiu Zhao, Bowen Yu, Feifan Song, Hangyu Li, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. Api-bank: A comprehensive benchmark for tool-augmented llms. *arXiv* preprint arXiv:2304.08244, 2023.
  - Xuefeng Li, Haoyang Zou, and Pengfei Liu. Torl: Scaling tool-integrated rl. *arXiv preprint arXiv:2503.23383*, 2025.
    - Weiwen Liu, Xu Huang, Xingshan Zeng, Xinlong Hao, Shuai Yu, Dexun Li, Shuai Wang, Weinan Gan, Zhengying Liu, Yuanqing Yu, et al. Toolace: Winning the points of llm function calling. *arXiv* preprint arXiv:2409.00920, 2024a.
    - Zhiwei Liu, Jielin Qiu, Shiyu Wang, Jianguo Zhang, Zuxin Liu, Roshan Ram, Haolin Chen, Weiran Yao, Shelby Heinecke, Silvio Savarese, et al. Mcpeval: Automatic mcp-based deep evaluation for ai agent models. *arXiv preprint arXiv:2507.12806*, 2025.

Zuxin Liu, Thai Hoang, Jianguo Zhang, Ming Zhu, Tian Lan, Juntao Tan, Weiran Yao, Zhiwei Liu,
 Yihao Feng, Rithesh RN, et al. Apigen: Automated pipeline for generating verifiable and diverse
 function-calling datasets. Advances in Neural Information Processing Systems, 37:54463–54482,
 2024b.

Elias Lumer, Pradeep Honaganahalli Basavaraju, Myles Mason, James A Burke, and Vamse Kumar Subbiah. Graph rag-tool fusion. *arXiv preprint arXiv:2502.07223*, 2025.

- Zhiling Luo, Xiaorong Shi, Xuanrui Lin, and Jinyang Gao. Evaluation report on mcp servers. *arXiv* preprint arXiv:2504.11094, 2025a.
- Ziyang Luo, Zhiqi Shen, Wenzhuo Yang, Zirui Zhao, Prathyusha Jwalapuram, Amrita Saha, Doyen Sahoo, Silvio Savarese, Caiming Xiong, and Junnan Li. Mcp-universe: Benchmarking large language models with real-world model context protocol servers. *arXiv preprint arXiv:2508.14704*, 2025b.
- Grégoire Mialon, Clémentine Fourrier, Thomas Wolf, Yann LeCun, and Thomas Scialom. Gaia: a benchmark for general ai assistants. In *The Twelfth International Conference on Learning Representations*, 2023.
- Guozhao Mo, Wenliang Zhong, Jiawei Chen, Xuanang Chen, Yaojie Lu, Hongyu Lin, Ben He, Xianpei Han, and Le Sun. Livemcpbench: Can agents navigate an ocean of mcp tools? *arXiv* preprint arXiv:2508.01780, 2025.
- Hussein Mozannar, Gagan Bansal, Cheng Tan, Adam Fourney, Victor Dibia, Jingya Chen, Jack Gerrits, Tyler Payne, Matheus Kunzler Maldaner, Madeleine Grunde-McLaughlin, et al. Magentic-ui: Towards human-in-the-loop agentic systems. *arXiv preprint arXiv:2507.22358*, 2025.
- OpenAI. Models. https://platform.openai.com/docs/models, September 2025. Accessed: 2025-09-23.
- Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. Gorilla: Large language model connected with massive apis, 2023. *URL https://arxiv. org/abs/2305.15334*, 2023.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*, 2023.
- Bernardo Quintero. What 17,845 github repos taught us about malicious mcp servers. https://blog.virustotal.com/2025/06/what-17845-github-repos-taught-us-about.html, June 2025. Accessed: 2025-09-23.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:68539–68551, 2023.
- Yijia Shao, Vinay Samuel, Yucheng Jiang, John Yang, and Diyi Yang. Collaborative gym: A framework for enabling and evaluating human-agent collaboration. *arXiv preprint arXiv:2412.15701*, 2024.
- Yijia Shao, Humishka Zope, Yucheng Jiang, Jiaxin Pei, David Nguyen, Erik Brynjolfsson, and Diyi Yang. Future of work with ai agents: Auditing automation and augmentation potential across the us workforce. *arXiv preprint arXiv:2506.06576*, 2025.
- Zhengliang Shi, Yuhan Wang, Lingyong Yan, Pengjie Ren, Shuaiqiang Wang, Dawei Yin, and
   Zhaochun Ren. Retrieval models aren't tool-savvy: Benchmarking tool retrieval for large language models. arXiv preprint arXiv:2503.01763, 2025.
  - Yueqi Song, Frank Xu, Shuyan Zhou, and Graham Neubig. Beyond browsing: Api-based web agents. *arXiv preprint arXiv:2410.16464*, 2024.

- Aditya Bharat Soni, Boxuan Li, Xingyao Wang, Valerie Chen, and Graham Neubig. Coding agents with multimodal browsing are generalist problem solvers. *arXiv preprint arXiv:2506.03011*, 2025.
  - Olly Styles, Sam Miller, Patricio Cerda-Mardini, Tanaya Guha, Victor Sanchez, and Bertie Vidgen. Workbench: a benchmark dataset for agents in a realistic workplace setting. *arXiv* preprint *arXiv*:2405.00823, 2024.
  - Hongru Wang, Cheng Qian, Wanjun Zhong, Xiusi Chen, Jiahao Qiu, Shijue Huang, Bowen Jin, Mengdi Wang, Kam-Fai Wong, and Heng Ji. Acting less is reasoning more! teaching model to act efficiently. *arXiv* preprint arXiv:2504.14870, 2025.
  - Jize Wang, Ma Zerun, Yining Li, Songyang Zhang, Cailian Chen, Kai Chen, and Xinyi Le. Gta: a benchmark for general tool agents. *Advances in Neural Information Processing Systems*, 37: 75749–75790, 2024a.
  - Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji. Executable code actions elicit better llm agents. In *Forty-first International Conference on Machine Learning*, 2024b.
  - Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, et al. Openhands: An open platform for ai software developers as generalist agents. *arXiv preprint arXiv:2407.16741*, 2024c.
  - Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. Autogen: Enabling next-gen llm applications via multiagent conversation framework. *arXiv* preprint arXiv:2308.08155, 3(4), 2023.
  - Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. Agentless: Demystifying llm-based software engineering agents. *arXiv preprint arXiv:2407.01489*, 2024.
  - Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh J Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, et al. Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments. *Advances in Neural Information Processing Systems*, 37:52040–52094, 2024.
  - Frank F Xu, Yufan Song, Boxuan Li, Yuxuan Tang, Kritanjali Jain, Mengxue Bao, Zora Z Wang, Xuhui Zhou, Zhitong Guo, Murong Cao, et al. Theagentcompany: benchmarking llm agents on consequential real world tasks. *arXiv* preprint arXiv:2412.14161, 2024a.
  - Qiancheng Xu, Yongqi Li, Heming Xia, and Wenjie Li. Enhancing tool retrieval with iterative feedback from large language models. *arXiv preprint arXiv:2406.17465*, 2024b.
  - John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. Swe-agent: Agent-computer interfaces enable automated software engineering. *Advances in Neural Information Processing Systems*, 37:50528–50652, 2024.
  - Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan.  $\tau$ -bench: A benchmark for tool-agent-user interaction in real-world domains. *arXiv preprint arXiv:2406.12045*, 2024.
  - Ming Yin, Dinghan Shen, Silei Xu, Jianbing Han, Sixun Dong, Mian Zhang, Yebowen Hu, Shujian Liu, Simin Ma, Song Wang, et al. Livemcp-101: Stress testing and diagnosing mcp-enabled agents on challenging queries. *arXiv preprint arXiv:2508.15760*, 2025.
  - Lifan Yuan, Yangyi Chen, Xingyao Wang, Yi R Fung, Hao Peng, and Heng Ji. Craft: Customizing llms by creating and retrieving from specialized toolsets. *arXiv preprint arXiv:2309.17428*, 2023.
  - Lucen Zhong, Zhengxiao Du, Xiaohan Zhang, Haiyi Hu, and Jie Tang. Complexfuncbench: exploring multi-step and constrained function calling under long-context scenario. *arXiv preprint arXiv:2501.10132*, 2025.
  - Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, et al. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023.

Model	#Retrieved	#MCP	Failed	#Retrieval	Query
	Tools	Calls	Calls (%)	Attempts	Length
		Primi	tive		
Sonnet 4	19.1	9.5	22.1	4.2 3.7	42.8
Opus 4	33.0	8.2	8.5		39.1
GPT 4.1	10.8	5.6	39.3	2.7 2.9	33.6
o3	24.2	2.6	11.5		23.8
GPT-5-mini	15.4	2.8	25.0	0.9	44.9
GPT 5	22.5	9.7	17.5		67.0
		Compo	site		
Sonnet 4	37.7	12.0	23.8	8.7	45.4
Opus 4	59.0	10.7	16.0	6.9	44.3
GPT 4.1	14.0	4.4	54.8	3.4	49.7
o3	6.4	0.9	33.3		30.0
GPT-5-mini	15.8	1.3	11.1	1.1 6.0	92.0
GPT 5	29.6	13.6	25.3		89.4

Table 5: MCPAgent's tool calling statistics on our Primitive and Composite Azure tasks. Query length is measured in the number of characters.

# A IMPLEMENTATION DETAILS

We implement our agent based on the OpenHands 0.48.0 CodeAct agent, with slight changes (Wang et al., 2024c). We remove the browser tool from the environment and instead provide the agent with the gateway MCP server, described in Section 4. In our experiments, we notice that LLMs often call the MCP tools directly and do not use the call\_tool function from the gateway MCP server. To avoid runtime errors, we allow the agent to call the MCP tools directly. Then, we post-process the LLM response and replace direct MCP tool calls with calls to call\_tool function.

We also extend the system prompt and provide the agent with additional guidance for using the MCP tools and interacting with the environment. Specifically, we append the information in Table 8 to the end of the original OpenHands CodeAct agent's system prompt. For fair comparisons, we also update the system prompt for browser-based agent and the agent with access to ground tools and include any information from Table 8 that is applicable to other agents. See Table 6 and Table 7 for the exact information that is added to the system prompt of the browser-based agent and agent with access to ground truth tools, respectively.

In our experiments, we disable the vision capabilities of models and evaluate the tasks solely based on the models' text understanding and generation capabilities. For all other configurations and hyperparameters, we use the default values from OpenHands.

# B MCP TOOL DOCUMENTATIONS

The API specifications for Azure and GitLab APIs provide high-quality documentation for each endpoint. However, RocketChat's OpenAPI specifications do not provide good descriptions for many of the endpoints. To improve the documentation quality, we use the original OpenHands' CodeAct agent to rewrite the description for each endpoint based on the documentation available on the web. Specifically, we prompt GPT-4.1 with the user prompt in Table 9 to generate new descriptions for each RocketChat endpoint.

```
759
760
761
762
763
764
765
766
767
768
769
       <COMPANY_ENVIRONMENT>
770
       - Everyone in this company is very responsive. People often respond
771
       to your messages immediately. The good thing is that you do not need
772
       to wait a long time for other's response. You just check your messages
773
       immediately and often times they have already responded to you.
774
       - **Very important** If you need a response from an employee, check
       if they have replied before finishing the task. You should never (I
775
       emphasize NEVER) finish the task without checking if they have responded
776
777
       - Our company hosts an internal version of Owncloud, GitLab, Plane, and
778
       RockChat. Do **NOT** access the public version of these services.
779
        </COMPANY_ENVIRONMENT>
       <GITLAB_INSTRUCTIONS>
781
        - You should always try to use the browser to interact with our
782
       internal GitLab instance. But, if it is absolutely necessary
783
       to call the GitLab REST APIs directly, you might do so using
784
       curl like the following: 'curl -H "PRIVATE-TOKEN: root-token"
        "http://the-agent-company.com:8929/api/v4/REST/API/PATH"
785
       - If you need to clone a repo from gitlab, use the following
786
       credentials:
787
        - username:
                    root
788
       - password:
                    theagentcompany
789
       - For some tasks, it is easier to clone the repo and work locally than
       working with the the repo in the browser. For example, if you need to
790
       explore the structure of a repo, read many files, etc., it is easier to
791
       clone the repo and work with its local version.
792
       </GITLAB_INSTRUCTIONS>
793
794
```

Table 6: The additional information appended to OpenHands (Wang et al., 2024c) CodeAct system prompt for the agent that uses the browser tool.

852 853

854

```
813
815
816
817
818
819
        <COMPANY_ENVIRONMENT>
820
        - Everyone in this company is very responsive. People often respond
821
        to your messages immediately. The good thing is that you do not need
822
        to wait a long time for other's response. You just check your messages
823
        immediately and often times they have already responded to you.
        - **Very important** If you need a response from an employee, check
824
        if they have replied before finishing the task. You should never (I
825
        emphasize NEVER) finish the task without checking if they have responded
826
        or not.
827
        - Our company hosts an internal version of Owncloud, GitLab, Plane,
828
        RockChat, and Azure. You can interact with these internal services
        using tools. Do **NOT** access the public version of these services.
829
        </COMPANY_ENVIRONMENT>
830
831
        <GITLAB_INSTRUCTIONS>
832
        - You must always use tools to interact with GitLab.
833
        - Remember, you should not access 'gitlab.com' which is the public
        version. Instead you should use tools to access our internal GitLab
834
835
        - If you need to clone a repo, first use tools to find the http url of
836
        the repo for cloning. Then use this internal url with the git command
837
        as usual.
838
        - Do not try to guess the web address of the internal GitLab.
839
        use tools to get the precise url for each GitLab project if needed.
        - You should always try to use tools to interact with our
840
        gitlab instance. But, if it is absolutely necessary to call
841
        the GitLab REST APIs directly, you might do so using curl like the following: 'curl -H "PRIVATE-TOKEN: root-token"
842
        "http://the-agent-company.com:8929/api/v4/REST/API/PATH"'
843
        - If you need to clone a repo from gitlab, use the following
844
        credentials:
845
        - username: root
846
        - password: theagentcompany
847
        - For some tasks, it is easier to clone the repo and work locally than
848
        calling many tools. For example, if you need to explore the structure
        of a repo, read many files, etc., it is easier to clone the repo and
849
        work with its local version.
850
        </GITLAB INSTRUCTIONS>
851
```

Table 7: The additional information appended to OpenHands (Wang et al., 2024c) CodeAct system prompt for the agent that has access to the oracle tool set.

```
864
          <TOOL_USE_INSTRUCTIONS>
865
          - In addition to the tools that are given to you in the current context window, there
866
          are tens of thousands of other external tools that you can use. However, they are not
          immediately available to you.
867
          - You can use the external tools to interact with RocketChat, Owncloud, Plane project
          management platform, gitlab, azure, etc.
868
          - To use external tools, you first have to find the tools that you need. You should use the "find-tools" tool to search for useful tools. Think of "find-tools" as a search engine for
869
          tools. Given a query, it returns the useful or related tools for that query.
870
          - Once you find the tools that you need, you can call them as you call any other tool.
871
          </TOOL_USE_INSTRUCTIONS>
          <TOOL_USE_BEST_PRACTICES>
872
          - You should come up with a plan for solving the task step by step. Then follow the plan
873
          step by step and potentially use external tools if needed to complete each step.
          - External tools empower you with new capabilities. Make full use of them. For example when
874
          the user asks you "find the cheapest iphone", although you currently have no way of knowing
875
          the price of an iphone, you can search for tools that help you with this step. For instance, you can call "find_tools("electronic price list")" and it could return tools that can provide
876
          you with the information that you need.
877
            If you fail to find the correct tools the first time, change the query and search again.
          - If you find a useful tool but you do not have the exact input arguments that it requires,
878
          do not give up. You can search for other tools that help you obtain the input arguments for
879
          that tool.
          - For example, if you want to check the price of an item based on its name but you find
880
          a tool that returns the price but needs the inventory ID, you should search and find an
          additional tool that helps you find the inventory ID from product name.
          - If you find an external tool but you are not able to successfully invoke the tool (e.g.,
882
          you get errors desipte multiple attempts), you should not give up.
883
          find another tool that provides a similar functionality.
          - Often there are multiple trajectories that could solve a task. If you were not able to
884
          solve the task with your current approach (e.g., did not find the correct tools or were not
885
          able to successfully call the tools), you should try again. Find new tools that could do the
          same thing and try again.
886
           - For example, if you want to check the price of a product but the tool that returns the
887
          prices raises a permission error, you could try to find a tool that returns recent purchase
          receipts for that item and extract its price from the receipts.
888
           - You should attempt 3-4 different potential trajectories with different tools and try to
889
          find a feasible solution for the task based on the available tools before giving up.
          - If you fail at any step, regardless of whether you have used external tools in that step,
890
          you should search for potential external tools that could help you accomplish that step
891
          successfully.
           For example if you tried to access a service directly by URL and failed, you should try to
892
          find an external tool for completing that step.
893
          </TOOL_USE_BEST_PRACTICES>
          <COMPANY_ENVIRONMENT>
894
           - Everyone in this company is very responsive. People often respond to your messages
895
          immediately. The good thing is that you do not need to wait a long time for other's
                     You just check your messages immediately and often times they have already
          response.
          responded to you.
897
           **Very important** If you need a response from an employee, check if they have replied
          before finishing the task. You should never (I emphasize NEVER) finish the task without
898
          checking if they have responded or not.
899
          - Our company hosts an internal version of Owncloud, GitLab, Plane, RockChat, and Azure.
          You can interact with these internal services using external tools as explained above. Do
900
          **NOT** access the public version of these services.
901
          </COMPANY_ENVIRONMENT>
          <GITLAB INSTRUCTIONS>
902
          - You must always use the external tools (explained above) to interact with GitLab.
903
          - Remember, you should not access 'gitlab.com' which is the public version. Instead you
          should use tools to access our internal gitlab instance.
904
          - If you need to clone a repo, first use external tools to find the http url of the repo for
905
          cloning. Then use this internal url with the git command as usual.
          - Do not try to guess the web address of the internal GitLab.
                                                                             Instead use the external tools
906
          to get the precise url for each GitLab project if needed.
907
          - You should always try to use the external tools to interact with our gitlab
          instance. But, if it is absolutely necessary to call the GitLab REST APIs directly, you might do so using curl like the following: `curl -H "PRIVATE-TOKEN: root-token"
908
909
          "http://the-agent-company.com:8929/api/v4/REST/API/PATH"
          - If you need to clone a repo from gitlab, use the following credentials:
910
          - username: root
911
                        theagentcompany
          - password:
          - For some tasks, it is easier to clone the repo and work locally than calling many external
912
          tools. For example, if you need to explore the structure of a repo, read many files, etc.,
913
          it is easier to clone the repo and work with its local version.
          </GITLAB_INSTRUCTIONS>
914
```

Table 8: The additional information appended to OpenHands (Wang et al., 2024c) CodeAct system prompt for MCPAgent, which uses tool retrieval to discover the required tools for each task.

967

```
921
922
       Your task is to create a summary and description for a RocketChat REST
923
       API endpoint.
924
       <RELATED RESOURCES>
925
       - RocketChat OpenAPI specifications:
926
       https://github.com/RocketChat/Rocket.Chat-Open-API
927
       - RocketChat API documentation website:
928
       https://developer.rocket.chat/apidocs
929
       </RELATED RESOURCES>
930
       <INPUT FORMAT>
931
       You will get an endpoint formatted as "HTTP_METHOD API_PATH"
932
       You also get a category that helps you find the documentation or
933
       specification for the endpoint.
       </INPUT FORMAT>
934
935
       <OUTPUT FORMAT>
936
       The output must be a json file (api_info.json) with three keys, endpoint,
937
       summary and description. Like the following:
938
       "endpoint": "endpoint given in the input task",
939
       "summary": "short summary",
940
       "description": "longer description of what the API does plus any
941
       additional information."
942
943
       </OUTPUT FORMAT>
944
       <NOTES>
945
       "summary" is only **ONE** sentence that very briefly describes what the
946
       endpoint does.
947
       "description" is often longer but not too long. It can contain any
948
       extra details that helps to use the endpoint correctly once the user
949
       decided to use it.
       </NOTES>
950
951
       952
       Task:
953
       create a summary and description for "POST /api/v1/channels.create" in
       the "rooms" category
954
       OUTPUT (content of api_info.json):
955
956
       "endpoint": "POST /api/v1/channels.create",
957
       "summary": "Create a public channel",
958
       "description": "Create a public channel. You can also include
       specified users, set permissions, and more."
959
960
       961
962
       ## Task
963
       Create a summary and description for "${method} ${path}" in the
964
       "${category}" category.
965
966
```

Table 9: The task description used to prompt GPT-4.1 to rewrite RocketChat tool descriptions.