# **ToolComp: A Multi-Tool Reasoning & Process Supervision Benchmark**

#### Anonymous Author(s)

Affiliation Address email

# Abstract

Despite recent advances in AI, the development of systems capable of executing 2 complex, multi-step reasoning tasks involving multiple tools remains a significant challenge. Current benchmarks fall short in capturing the real-world complex-3 ity of tool-use reasoning, where verifying the correctness of not only the final answer but also the intermediate steps is important for evaluation, development, 5 and identifying failures during inference time. To bridge this gap, we introduce 6 ToolComp, a comprehensive benchmark designed to evaluate multi-step tool-use reasoning. ToolComp is developed through a collaboration between models and 8 human annotators, featuring human-edited/verified prompts, final answers, and 9 process supervision labels, allowing for the evaluation of both final outcomes and 10 intermediate reasoning. Evaluation across six different model families and 20 12 total models demonstrates the challenging nature of our dataset, with an average accuracy of 55% among the frontier models. <sup>1</sup>

# Introduction

11

13

Recent advancements in large language models (LLMs) have demonstrated remarkable progress in a 15 range of natural language processing tasks. These models have achieved state-of-the-art performance 16 across diverse benchmarks, including question answering, summarization, and reasoning tasks. In 17 order to further increase the usefulness of LLMs, a growing area of research is centered around the 18 development of agentic capabilities, particularly their ability to autonomously interact with external 19 tools to solve complex, multi-step tasks as well as to interact with human systems such as the web or 20 mobile devices. 21

However, evaluating the effectiveness of these tool-use capabilities remains a pressing challenge. 23 While there have been notable efforts in developing benchmarks for tool-use capability, these often assess isolated instances of tool use, focusing on whether the model can invoke the correct tool 24 at the right time (Huang et al., 2024; Zhuang et al., 2023; Peng et al., 2021). Additionally, while 25 benchmarks for multi-step tool usage exist, most focus only on scoring the correctness of the final 26 answer (Mialon et al., 2023), despite that the complex nature of multi-step reasoning often requires 27 the evaluation for partial correctness or step-wise correctness of the reasoning trajectories. This can 28 29 be valuable for both understanding model failure modes and developing systems that can improve upon these intermediate reasoning flaws. 30

To address these shortcomings, we introduce ToolComp, a benchmark comprising 493 complex, 31 human-verified prompts that require language models to chain together multiple tool calls, accompanied by human-edited step-wise and final answers. By demanding intricate tool interactions and

<sup>&</sup>lt;sup>1</sup>Code and data is publicly available. A few data examples are shown in the supplementary materials.

providing human verification, ToolComp offers a rigorous assessment of a model's ability to perform complex, multi-step reasoning and tool use. We evaluate the current landscape of state-of-the-art models on their ability to chain together tool calls to reach the final answer, as well as their step-wise reasoning ability.

#### 1.1 Contributions and Key Takeaways

39 Our key contributions and takeaways are summarized as follows:

- **Introduction of ToolComp** We introduce ToolComp, a multi-tool reasoning and process supervision benchmark with 493 human-edited/verified prompts and final answers, designed to evaluate a model's ability to perform multi-step tool-use tasks (**Section 3**).
- **Step-by-Step Process Annotations** ToolComp includes 1716 detailed per-step supervision labels, enabling a comprehensive assessment of a model's intermediate reasoning when performing complex, multi-step tool-use tasks (**Section 3**).
- Assessment of State-of-the-Art Models We evaluate 20 models across 6 different model families on their ability to perform complex multi-step tool-use tasks as well as their intermediate reasoning ability. We find that GPT-5 has the best final answer performance, achieving 79.81% against human-verified final answers, and Gemini 2.5 Pro has the best performance against process supervision labels, achieving 83.42%. (Section 4 and Section A).

#### 2 Related Works

38

40

41

42

43

44

45

46

47

49

50

51

52

53

54

55

56

57

60

61 62

63

64

65

66

69

70

71

72

73

75

76

77

79

80

81

Benchmarks for Complex Tool Use Planning With rising interest in tool-augmented LLMs (Schick et al., 2023; Patil et al., 2023; Qin et al., 2023), several benchmarks have been introduced to assess their abilities. Earlier benchmarks were designed to assess a model's ability to do proper retrieval, execution, and extraction of one tool call for specific tasks such as general question answering (Yang et al., 2018; Joshi et al., 2017), fact verification (Thorne et al., 2018), or answering temporal queries (Chen et al., 2021; Kasai et al., 2024; Zhang & Choi, 2021; Dhingra et al., 2022; Vu et al., 2023). However, these benchmarks fail to assess a model's ability to plan and chain together multiple tool calls to answer more complex queries. More recent benchmarks aimed at evaluating multiple tool calls are often placed within or dependent on state-full systems (such as a code-base and/or a dynamic database) (Yan et al., 2024; Jimenez et al., 2024; Liu et al., 2023). Although these types of benchmarks assess a language model's ability to chain together multiple tool calls, the evaluation may penalize general-purpose language models that are not familiar with the given environments. Other benchmarks primarily rely on state-based evaluations, where the final state of the system is assessed against the desired state (Li et al., 2023; Peng et al., 2021), or win-rates against another reference state-of-the-art model (Qin et al., 2023), both of which lack the rigour of human-verified ground truth final answers. Closest to our work, the GAIA benchmark is a collection of complex tool-use queries that require multi-step tool-use reasoning and associated ground-truth answers (Mialon et al., 2023). Crucially, it does not contain step-wise labels, which can be important for identifying where an error occurred and providing precise feedback. Additionally, a significant portion of GAIA requires specialized capabilities such as web browsing, multi-modality, and diverse file-type reading. In our work, we focus on text-only tasks in order to disentangle specialized capabilities and multi-step reasoning, allowing us to focus on the latter.

**Process Reward Models** Recent work has shown the power of utilizing process supervision signals, which are granular signals on the step-wise correctness of a solution, as opposed to outcome supervision signals, which are broad signals on the correctness of the entire solution. Utilizing these signals, Lightman et al. (2023) and Wang et al. (2024a) have shown dramatic improvements in performance in ranking solutions to mathematical reasoning tasks and using these signals to further improve performance in traditional RLHF algorithms such as Proximal Policy Optimization (PPO) (Schulman et al., 2017).

In this work, through a hybrid human-AI annotation workflow, we generate per-step process supervision labels, which uniquely enable us to rigorously evaluate a model's intermediate reasoning capability. Table 1 provides a comparative overview of popular tool-use benchmarks, including our work, ToolComp.

Table 1: The contributions and metadata of popular benchmarks in Tool Use. Our work, ToolComp, is shown in the first column. From left to right, we include work from Mialon et al. (2023), Yan et al. (2024), Qin et al. (2023), Li et al. (2023), and Xu et al. (2023). \* Although 2 of the 8 tools are not evaluated by simply matching a verified final answer, the remaining 6 have verified final answers.

Resource	ToolComp	GAIA	BFCL	ToolBench	API-Bank	ToolBench
Real-World API Calls	✓	✓	✓	✓	✓	<b>✓</b>
Multi-Tools Scenario	✓	✓	✓	✓	X	X
Multi-Step Reasoning	✓	✓	$\checkmark$	✓	✓	✓
Step-Wise Labels	✓	X	X	X	X	X
Verified Final Answer	✓	✓	X	X	X	<b>/</b> *
Number of Tools	11	23	3	3451	53	8

# 86 3 ToolComp

### 87 **3.1 Tools**

88

89

90

91

92

93

94

95

For the creation of this benchmark and evaluation framework, we support 11 tools: Date, Current Weather, Historical Weather (Zippenfenig, 2024), Calculator, Wiki Search (Majlis, 2017), Google Search (SerpApi, 2024), Wolfram Alpha (Wolfram Research, 2024), Intra-day Stock Info, Daily Stock Info, Stock Symbol Search (Alpha Vantage), and Python. There were several considerations when choosing these set of tools, namely, we wanted to cover a broad range of use cases from fact retrieval to financial assistant, have some overlap in use cases to encourage various valid trajectories, ensure the tools are general enough to not require specialized knowledge for LLMs to use, and allow for interesting interactions between tools. A detailed breakdown of each tool, including descriptions, parameters, input examples, and output examples are available in Appendix E.

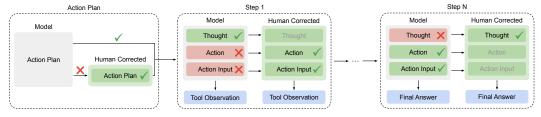


Figure 1: An example annotation path for collecting data that provides tool-call trajectories with human verified-final answers along with step-by-step process supervision labels. Each model generated step (Action Plan and ReAct steps) are first labelled as correct or incorrect. For the components labelled incorrect, a rewrite is made to correct the corresponding component. The annotations and rewrites are made by human annotators for the benchmark. A full annotated trajectory example is available in Appendix D.2.

### 7 3.2 ReAct Format

We chose the ReAct format as it is frequently used for tool use and agentic workflows (Wang et al., 2024b; Mekala et al., 2024; Zhuang et al., 2023). The ReAct format combines reasoning and tool calls by prompting the model to first generate a thought, which contains the rationale behind the following tool call action (Yao et al., 2023). The structured nature of the ReAct format into a thought, action, action input, and observation allows us to collect granular signals at each sub-step, and the relative simplicity of the ReAct format makes it easier to operationalize for annotations.

# 3.3 Prompt Creation

104

In developing the prompts for this dataset, there are two main criteria we desire each prompt to satisfy:

1) the solution to the prompt contains a chain of dependent tool calls to answer and 2) the final answer

107 to the prompt can be programmatically verified. To achieve this, we generate a set of candidate

108 prompts through few-shot prompting which are then refined and validated by human annotators.

The overall process includes 1) manually developing in-context (IC) examples, 2) generating initial prompts, 3) an iterative process of filtering prompts, adding filtered prompts as negative IC examples, and regenerating more prompts, and 4) human refinement. These steps are described in more detail in Appendix B.1

#### 113 3.4 Chat vs. Enterprise Use Cases

In creating the benchmark, we developed two subsets of prompts, coined ToolComp-Enterprise 114 and ToolComp-Chat. ToolComp-Enterprise allows the use of 11 tools and aims to emulate settings in which LLM agents must compose a larger number of expressive APIs together correctly, such 116 as in enterprise settings. The second subset, ToolComp-Chat, is designed to test general purpose 117 chatbots with the minimally sufficient set of tools for information retrieval and processing tasks, 118 namely Google Search and Python. We chose only google search and python execution as these are 119 standard tools found in major chat-bot providers. We only allow the respective tools for each subset 120 during prompt generation, labeling, and evaluation. ToolComp-Enterprise contains 296 examples and 121 ToolComp-Chat contains 197 examples. 122

#### 3.5 Label Creation

123

To create the process supervision labels as well as the final answer for each prompt, we utilize a hybrid human-AI approach, where the language model and human annotators use the same tools to collaborate to get to the final answer. We start by prompting the Policy Model LLM to outline a plan, called Action Plan, on which tools to call and in what order using the prompt in C.1. We have human annotators validate/modify the Action Plan, which is then appended to the sequence before using the LLM to formulate tool calls. We then use the LLM to call tools in the ReAct format, where the specific prompt can be found in C.2.

We asked the annotators to rate if a step is Correct (i.e., the step is a reasonable action towards 131 achieving the final answer) or Incorrect (i.e., the step is nonsensical, incorrect, or is not a reasonable 132 action towards acheiving the final answer). All components of the ReAct Step (Thought, Action, 133 Action Input) must be marked as Correct or Incorrect by the annotator. If the annotator marks a 134 step as Correct, the model is allowed to proceed further and generate the next step. If the annotator 135 deems a step as Incorrect, they must modify the step to make it correct. Once corrected, the model is 136 then prompted to advance to the next step with the human-corrected step as part of its context. This is repeated until the Finish Action is chosen by the LLM and marked as Correct by the annotator 138 or until the annotator corrects an Action step to 'Finish' because we have enough information to 139 answer the question. The overall flow is shown in Figure 1. An example golden trajectory is 140 available in Appendix D.1 and an example annotated trajectory is available in Appendix D.2. We use 141 FireFunction-V1 as the Policy Model LLM (at the time, this was the best open-source tool-use LLM) 142 and humans as the annotators (Fireworks, 2024). 143

With this process, we retrieve, per task, a valid step-by-step chain of tool calls that successfully gets to the final answer along with step-wise correct/incorrect labels and associated rewrites. The correct/incorrect labels and the associated rewrites allow us to assess intermediate reasoning through LLM-as-judge evaluations (described in Section 4.3).

#### 3.6 Quality Control

148

149

150 Any data samples with ambiguous prompts, erroneous process supervision labels, or incorrect final answers are redone. After the initial creation of the benchmark, the authors collaborated with three 151 trusted annotators to perform a final re-review of all samples and make any necessary corrections. 152 As a final quality control step, we evaluate the entire benchmark using GPT-40 (May 2024), GPT-4 153 Turbo, Claude 3.5 Sonnet, and Llama 3.1 405b (OpenAI et al., 2024; Dubey et al., 2024; Anthropic). 154 We identify the set of data samples where all models' answers differed from the ground truth final 155 answers. We then repeated the refinement process on these samples, as they represented the most 156 challenging and/or potentially mislabeled data points. This iterative approach yielded the final version 157 of ToolComp.

To ensure the highest quality of ToolComp, we conduct a thorough manual inspection of all examples.

# 59 4 ToolComp Evaluations

#### 4.1 Evaluation Metric

We have two metrics to evaluate the quality or the correctness of a model's final answers: LLM Grading and Exact Match. For the final answer evaluations in this section (Table 2), we use LLM Grading since it rewards correct answers without penalizing minor formatting issues. Our Exact Match evaluation methodology and the corresponding results are shown in Appendix A.1.

**LLM Grading** By using LLM grading against ground truth answers we opt to be charitable to exact formatting and focus on assessing the tool use capabilities of the model. We intentionally choose not to focus on final answer formatting given that (1) there are existing benchmarks that assess formatting ability (e.g. FOFO (Xia et al., 2024)) and (2) our final answers are quite complex, containing multiple elements, lists which may or may not be sorted, and dictionaries. This approach prompts an LLM Judge to look at the prompt, the ground truth answer, and the model's answer and asks the model to classify it as Incorrect, Correct, or Correct with Bad Formatting. We use GPT-4 Turbo as the de-facto judge for all of our models (OpenAI et al., 2024). The prompt used is shown in Appendix C.5. We consider both Correct and Correct with Bad Formatting as a win (accurate) and Incorrect as a loss (inaccurate).

#### 4.2 Final Answer Evaluations

Table 2: Accuracy and the 95% CIs of all selected models using the final answer and scored using an LLM judge (Dubey et al., 2024; OpenAI et al., 2024; Gemini et al., 2024; Anthropic; Mistral; Cohere). We combined the results of each subset to give an overall score for the entire benchmark. Exact Match results are reported in Appendix A.1 but the rankings do not significantly differ.

<b>Model Family</b>	Model Name	Total (%)	Chat (%)	Enterprise (%)
	GPT-5	$79.81 \pm 5.00$	$76.92 \pm 5.91$	$81.75 \pm 4.40$
	o3	$78.29 \pm 5.12$	$76.14 \pm 5.95$	$79.72 \pm 4.57$
OpenAI	o1	$66.25 \pm 5.92$	$60.41 \pm 6.82$	$70.14 \pm 5.32$
OpenAi	GPT-4o (Aug 2024)	$58.68 \pm 4.39$	$56.85 \pm 6.92$	$59.93 \pm 5.67$
	GPT-4	$45.89 \pm 4.43$	$37.88 \pm 6.78$	$51.39 \pm 5.77$
	GPT-4o Mini	$44.03 \pm 4.41$	$32.83 \pm 6.54$	$51.74 \pm 5.77$
	Claude 4.1 Opus	$75.85 \pm 5.30$	$76.14 \pm 5.95$	$75.67 \pm 4.88$
Anthumia	Claude 4 Sonnet	$75.65 \pm 4.39$	$74.61 \pm 6.91$	$76.35 \pm 5.67$
Anthropic	Claude 3.5 Sonnet	$58.03 \pm 4.39$	$56.06 \pm 6.91$	$59.38 \pm 5.67$
	Claude 3 Opus	$51.03 \pm 4.44$	$48.49 \pm 6.96$	$52.78 \pm 5.77$
	Claude 3 Sonnet	$48.56 \pm 4.44$	$40.4\pm6.84$	$54.17 \pm 5.78$
	Gemini 2.5 Pro	$77.07 \pm 5.21$	$77.15 \pm 5.86$	$77.02 \pm 4.79$
Google	Gemini 1.5 Pro (Aug 2024)	$56.61 \pm 4.41$	$51.27 \pm 6.98$	$60.28 \pm 5.66$
	Gemini 1.5 Pro (May 2024)	$38.43 \pm 4.34$	$35.50 \pm 6.57$	$40.42 \pm 5.68$
Mistral	Mistral Large 2	$46.30 \pm 4.43$	$40.4\pm6.84$	$50.35\pm5.78$
	Llama 4 Scout 17B Instruct	$61.64 \pm 4.44$	$61.42 \pm 6.79$	$61.82 \pm 5.53$
Meta	Llama 3.1 405B Instruct	$46.19 \pm 4.44$	$40.10 \pm 6.84$	$50.35 \pm 5.78$
	Llama 3.1 70B Instruct	$35.74 \pm 4.27$	$33.50 \pm 6.59$	$37.23 \pm 5.60$
	Llama 3.1 8B Instruct	$12.81 \pm 2.98$	$6.090 \pm 3.34$	$17.42 \pm 4.39$
Cohere	Command R+	$26.13 \pm 3.91$	$20.20 \pm 5.59$	$30.21 \pm 5.3$
	Average	55.46	51.10	56.94

The overall scores of the various state-of-the-art tool-use models are shown in Table 2. We combine ToolComp-Chat and ToolComp Enterprise subsets to get an overall score and 95% confidence-intervals (CIs) for the entire benchmark. We use native function calling for all the models and we allow each model to retry up to 3 times if it fails to output a final answer. This is determined by whether there is a parse-able JSON object in the final output with the key "final\_answer". To ensure

Table 3: Accuracy and the 95% CIs (third column) of all of our models on the process supervision labels in ToolComp. We evaluate the model's effectiveness as a pairwise judge in selecting the human-corrected answer versus the model-generated incorrect answer. We show judge accuracy using the ReAct steps (fourth column) and the Action Plan (fifth column).

<b>Model Family</b>	Model Name	Total (%)	ReAct (%)	Action Plan (%)
	GPT-5	$78.64 \pm 1.93$	$77.08 \pm 2.25$	$82.96 \pm 3.53$
	03	$75.58 \pm 1.96$	$72.56 \pm 2.32$	$83.95 \pm 3.49$
OpenAI	o1	$76.92 \pm 1.89$	$78.15 \pm 2.22$	$73.51 \pm 3.55$
OpenAi	GPT-4o (Aug 2024)	$72.61 \pm 2.11$	$72.84 \pm 2.46$	$71.98 \pm 4.13$
	GPT-40 Mini	$63.02 \pm 2.28$	$64.27 \pm 2.64$	$59.56 \pm 4.51$
	GPT-4	$60.02 \pm 2.32$	$55.87 \pm 2.74$	$71.54 \pm 4.15$
	Claude 4.1 Opus	$82.31 \pm 2.19$	$80.25 \pm 2.53$	$88.02 \pm 4.36$
Anthropic	Claude 4 Sonnet	$80.06 \pm 2.23$	$84.50 \pm 2.58$	$78.47 \pm 4.44$
Anunopic	Claude 3.5 Sonnet	$66.46 \pm 2.23$	$67.74 \pm 2.58$	$62.97 \pm 4.44$
	Claude 3 Opus	$64.28 \pm 2.27$	$64.55 \pm 2.64$	$63.52 \pm 4.42$
	Claude 3 Sonnet	$61.10 \pm 2.31$	$62.93 \pm 2.67$	$56.04 \pm 4.56$
	Gemini 2.5 Pro	$83.42\pm2.12$	$80.92\pm2.23$	$90.32 \pm 4.31$
Google	Gemini 1.5 Pro (Aug 2024)		$68.48 \pm 2.56$	$70.88 \pm 4.17$
	Gemini 1.5 Pro (May 2024)	$67.89 \pm 2.21$	$67.72 \pm 2.58$	$68.35 \pm 4.27$
Mistral	Mistral Large 2	$72.67 \pm 2.11$	$73.16 \pm 2.45$	$71.32 \pm 4.16$
	Llama 4 Scout 17B Instruct	$75.45 \pm 2.13$	$76.23 \pm 2.42$	$73.30 \pm 4.37$
Meta	Llama 3.1 405B Instruct	$71.62 \pm 2.13$	$73.87 \pm 2.42$	$65.39 \pm 4.37$
	Llama 3.1 70B Instruct	$70.75 \pm 2.15$	$71.33 \pm 2.50$	$69.12 \pm 4.25$
	Llama 3.1 8B Instruct	$57.63 \pm 2.34$	$59.60 \pm 2.71$	$52.20 \pm 4.56$
Cohere	Command R+	$61.31 \pm 2.30$	$64.91 \pm 2.63$	$51.32 \pm 4.59$
	Average	70.64	70.85	70.32

scores are not indicative of tool or endpoint failures due to rate limiting, we use verbose logging to log all failures and retry any prompt where a tool or model outputs failed due to rate/load limits. In addition, we run error analysis on the types of failures for each model. A description of the error category taxonomy and the breakdown of failure modes for each model can be found in Appendix A.2.

We also show exact match evaluation numbers in Table 4 of Appendix A.1 to ensure that our LLM Judge (OpenAI's o1) isn't biased in favor of outputs from the same model family. Upon inspection of the discrepancies (i.e., examples marked correct by the LLM judge but incorrect under exact match), we find that they are all due to issues with the model's formatting of the final answer despite getting to the correct answer.

#### 4.3 LLM-as-Judge Evaluations

We further evaluate these models using our process supervision labels, aiming to assess each model's effectiveness as a pairwise judge in selecting the human-corrected step over the step generated by the original policy used during annotation. To mitigate position bias, we swap the order of the human-corrected and model-generated steps and conduct two separate predictions for each arrangement. Additionally, models are permitted to indicate a tie. If a model designates a tie at least once, or consistently predicts the same position (before and after swapping) for a given data sample, we classify the outcome as a tie. Mirroring the methodology used in RewardBench (Lambert et al. (2024)), we score losses as 0, ties as 0.5, and wins as 1. We show the results below in Table 3.

# 4.4 Intermediate Reasoning vs. Final Answer

Figure 2 shows the correlation between a model's intermediate reasoning performance and final answer accuracy based on the multi-step tool-use tasks in ToolComp. The standard Pearson correlation

coefficient is r=0.83 with a statistical p-value of 0.000005, which makes the correlation statistically significant under a significance level of 0.05 (Freedman et al., 2007). Intuitively, this suggests that with stronger step-wise performance as assessed by our LLM-as-judge evaluations, we can expect an increased likelihood of reaching the correct final answer. However, the moderate magnitude of the correlation value could be due to additional signals captured by the step-wise reasoning evaluations that are not captured by evaluating final answers. Work done by Havrilla et al. (2024) similarly suggests that there is complementary and non-overlapping information in step-wise and final answer refinement, further highlighting the importance of assessing intermediate reasoning.

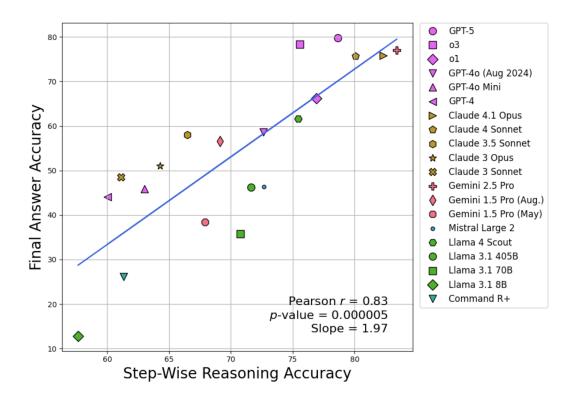


Figure 2: Comparison of step-wise reasoning accuracy (x-axis) and final answer accuracy (y-axis) on ToolComp across 6 different model families.

# Limitations and Biases

## 5.1 Methodological Limitations

**Limited Tool Scope** This work focuses on a restricted set of tools primarily designed for information retrieval and data processing. In contrast, many state-of-the-art systems employ specialized models for diverse tasks such as image generation, translation, and complex reasoning. This limitation raises important questions about how process supervision could scale to more nuanced capabilities when integrating with other specialized models and broader tool ecosystems.

## 5.2 Dataset Construction Biases

**Human Preference Bias in Step Correction** During the step correction process, human annotators naturally gravitated toward tools that were more intuitive or convenient to use. This preference created a systematic skew in the Step-Wise Reasoning data, with certain tools becoming overrepresented. We deliberately preserved this bias as it reflects authentic human (and likely model) tool selection patterns, prioritizing convenience and practical usability over uniform tool distribution.

- Programmatic Verification Constraints Each ToolComp prompt was engineered to have an unambiguous, programmatically verifiable ground truth answer. This design requirement necessitated somewhat artificial output formats that strictly conform to automated evaluation criteria. While this constraint ensures reliable evaluation metrics, it may not capture the natural variability and ambiguity present in real-world task specifications.
- Compositional Focus Limitations Given ToolComp's primary objective of evaluating tool composition abilities, the benchmark systematically excludes several categories of prompts: tasks requiring no tool usage; tasks falling outside the scope of provided tools; and tasks requiring clarifying questions or iterative dialogue. These exclusions ensure focused evaluation of compositional reasoning but limit the dataset's coverage of broader real-world use cases where tool selection and usage patterns may differ significantly.
- Generator Model Bias The use of Firefunction-v1 as the base model for generating initial trajectories introduces potential systematic biases into the dataset. This model's inherent preferences for certain tools, input formats, or reasoning patterns may propagate through the human annotation process, potentially skewing the final dataset distribution in ways that reflect the base model's limitations rather than optimal tool usage patterns.

# **6** Ethics Statement

240

We ensure all prompts in this dataset do not contain any harmful or sensitive material by requiring
annotators to flag any such prompts. The authors of this paper have also manually inspected all the
prompts and tool calls for harmful content. In addition, we applied best practices for code execution,
ensuring that all the code execution is done in a sand-boxed environment for any past and/or future
benchmark evaluations. We also ensured that all tools used have a permissive license for research
purposes, and we plan to open-source both the code for running evaluations and the full benchmark
dataset.

# 248 7 Reproducibility

For the creation of the benchmark, we detail the exact process by which we create the dataset in Section 3. We also detail the exact evaluation method used to evaluate each model in Section 4 and Appendix A.1. We have open sourced both the code for evaluation and the benchmark dataset for the final answer evaluation as well as the intermediate reasoning evaluation.

### 253 References

- Inc. AlphaVantage. Alphavantage stock api. https://www.alphavantage.co/. Accessed: February 2023–September 2024.
- Anthropic. Claude 3: An ai assistant by anthropic. https://www.anthropic.com. Accessed: 2024-02 to 2024-10.
- Wenhu Chen, Xinyi Wang, and William Yang Wang. A dataset for answering time-sensitive questions, 2021. URL https://arxiv.org/abs/2108.06314.
- 260 Cohere. Cohere command r+ model. https://cohere.com. Accessed: 2024-10-01.
- Bhuwan Dhingra, Jeremy R. Cole, Julian Martin Eisenschlos, Daniel Gillick, Jacob Eisenstein, and William W. Cohen. Time-aware language models as temporal knowledge bases. *Transactions of the Association for Computational Linguistics*, 10:257–273, 2022. ISSN 2307-387X. doi: 10.1162/tacl\_a\_00459. URL http://dx.doi.org/10.1162/tacl\_a\_00459.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris

McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruy Choudhary, Dhruy Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghayan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manay Avalani, Manish Bhatt, Maria

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

300

301

302

303

304

305

306 307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, 329 Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle 330 Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, 331 Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, 332 Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, 333 Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia 334 Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro 335 Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, 336 Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, 337 Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan 338 Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara 339 Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh 340 Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, 341 Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan 343 Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, 344 Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe 345 Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, 346 Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, 347 Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, 348 Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, 349 Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, 350 Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, 351 Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd 352 of models, 2024. URL https://arxiv.org/abs/2407.21783. 353

AI Fireworks. Firefunction-v1: Gpt-4 level function calling. https://fireworks.ai/blog/firefunction-v1-gpt-4-level-function-calling, 2024.

David Freedman, Robert Pisani, and Roger Purves. Statistics (international student edition). *Pisani*, *R. Purves*, 4th edn. WW Norton & Company, New York, 2007.

Gemini, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, 358 Damien Vincent, Zhufeng Pan, Shibo Wang, Soroosh Mariooryad, Yifan Ding, Xinyang Geng, Fred 359 Alcober, Roy Frostig, Mark Omernick, Lexi Walker, Cosmin Paduraru, Christina Sorokin, Andrea 360 Tacchetti, Colin Gaffney, Samira Daruki, Olcan Sercinoglu, Zach Gleicher, Juliette Love, Paul 361 Voigtlaender, Rohan Jain, Gabriela Surita, Kareem Mohamed, Rory Blevins, Junwhan Ahn, Tao 362 Zhu, Kornraphop Kawintiranon, Orhan Firat, Yiming Gu, Yujing Zhang, Matthew Rahtz, Manaal 363 Faruqui, Natalie Clay, Justin Gilmer, JD Co-Reyes, Ivo Penchev, Rui Zhu, Nobuyuki Morioka, 364 Kevin Hui, Krishna Haridasan, Victor Campos, Mahdis Mahdieh, Mandy Guo, Samer Hassan, 365 Kevin Kilgour, Arpi Vezer, Heng-Tze Cheng, Raoul de Liedekerke, Siddharth Goyal, Paul Barham, 366 DJ Strouse, Seb Noury, Jonas Adler, Mukund Sundararajan, Sharad Vikram, Dmitry Lepikhin, 367 Michela Paganini, Xavier Garcia, Fan Yang, Dasha Valter, Maja Trebacz, Kiran Vodrahalli, 368 Chulayuth Asawaroengchai, Roman Ring, Norbert Kalb, Livio Baldini Soares, Siddhartha Brahma, 369 David Steiner, Tianhe Yu, Fabian Mentzer, Antoine He, Lucas Gonzalez, Bibo Xu, Raphael Lopez 370 Kaufman, Laurent El Shafey, Junhyuk Oh, Tom Hennigan, George van den Driessche, Seth Odoom, 371 Mario Lucic, Becca Roelofs, Sid Lall, Amit Marathe, Betty Chan, Santiago Ontanon, Luheng He, 372 Denis Teplyashin, Jonathan Lai, Phil Crone, Bogdan Damoc, Lewis Ho, Sebastian Riedel, Karel 373 Lenc, Chih-Kuan Yeh, Aakanksha Chowdhery, Yang Xu, Mehran Kazemi, Ehsan Amid, Anastasia 374 Petrushkina, Kevin Swersky, Ali Khodaei, Gowoon Chen, Chris Larkin, Mario Pinto, Geng Yan, 375 Adria Puigdomenech Badia, Piyush Patil, Steven Hansen, Dave Orr, Sebastien M. R. Arnold, 376 Jordan Grimstad, Andrew Dai, Sholto Douglas, Rishika Sinha, Vikas Yadav, Xi Chen, Elena 377 Gribovskaya, Jacob Austin, Jeffrey Zhao, Kaushal Patel, Paul Komarek, Sophia Austin, Sebastian 378 Borgeaud, Linda Friso, Abhimanyu Goyal, Ben Caine, Kris Cao, Da-Woon Chung, Matthew 379 Lamm, Gabe Barth-Maron, Thais Kagohara, Kate Olszewska, Mia Chen, Kaushik Shivakumar, 380 Rishabh Agarwal, Harshal Godhia, Ravi Rajwar, Javier Snaider, Xerxes Dotiwalla, Yuan Liu, 381 Aditya Barua, Victor Ungureanu, Yuan Zhang, Bat-Orgil Batsaikhan, Mateo Wirth, James Qin, Ivo 382 Danihelka, Tulsee Doshi, Martin Chadwick, Jilin Chen, Sanil Jain, Quoc Le, Arjun Kar, Madhu 383 Gurumurthy, Cheng Li, Ruoxin Sang, Fangyu Liu, Lampros Lamprou, Rich Munoz, Nathan Lintz, 384 Harsh Mehta, Heidi Howard, Malcolm Reynolds, Lora Aroyo, Quan Wang, Lorenzo Blanco, Albin 385

Cassirer, Jordan Griffith, Dipanjan Das, Stephan Lee, Jakub Sygnowski, Zach Fisher, James Besley, Richard Powell, Zafarali Ahmed, Dominik Paulus, David Reitter, Zalan Borsos, Rishabh Joshi, Aedan Pope, Steven Hand, Vittorio Selo, Vihan Jain, Nikhil Sethi, Megha Goel, Takaki Makino, Rhys May, Zhen Yang, Johan Schalkwyk, Christina Butterfield, Anja Hauth, Alex Goldin, Will Hawkins, Evan Senter, Sergey Brin, Oliver Woodman, Marvin Ritter, Eric Noland, Minh Giang, Vijay Bolina, Lisa Lee, Tim Blyth, Ian Mackinnon, Machel Reid, Obaid Sarvana, David Silver, Alexander Chen, Lily Wang, Loren Maggiore, Oscar Chang, Nithya Attaluri, Gregory Thornton, Chung-Cheng Chiu, Oskar Bunyan, Nir Levine, Timothy Chung, Evgenii Eltyshev, Xiance Si, Timothy Lillicrap, Demetra Brady, Vaibhav Aggarwal, Boxi Wu, Yuanzhong Xu, Ross McIlroy, Kartikeya Badola, Paramjit Sandhu, Erica Moreira, Wojciech Stokowiec, Ross Hemsley, Dong Li, Alex Tudor, Pranav Shyam, Elahe Rahimtoroghi, Salem Haykal, Pablo Sprechmann, Xiang Zhou, Diana Mincu, Yujia Li, Ravi Addanki, Kalpesh Krishna, Xiao Wu, Alexandre Frechette, Matan Eyal, Allan Dafoe, Dave Lacey, Jay Whang, Thi Avrahami, Ye Zhang, Emanuel Taropa, Hanzhao Lin, Daniel Toyama, Eliza Rutherford, Motoki Sano, HyunJeong Choe, Alex Tomala, Chalence Safranek-Shrader, Nora Kassner, Mantas Pajarskas, Matt Harvey, Sean Sechrist, Meire Fortunato, Christina Lyu, Gamaleldin Elsayed, Chenkai Kuang, James Lottes, Eric Chu, Chao Jia, Chih-Wei Chen, Peter Humphreys, Kate Baumli, Connie Tao, Rajkumar Samuel, Cicero Nogueira dos Santos, Anders Andreassen, Nemanja Rakićević, Dominik Grewe, Aviral Kumar, Stephanie Winkler, Jonathan Caton, Andrew Brock, Sid Dalmia, Hannah Sheahan, Iain Barr, Yingjie Miao, Paul Natsev, Jacob Devlin, Feryal Behbahani, Flavien Prost, Yanhua Sun, Artiom Myaskovsky, Thanumalayan Sankaranarayana Pillai, Dan Hurt, Angeliki Lazaridou, Xi Xiong, Ce Zheng, Fabio Pardo, Xiaowei Li, Dan Horgan, Joe Stanton, Moran Ambar, Fei Xia, Alejandro Lince, Mingqiu Wang, Basil Mustafa, Albert Webson, Hyo Lee, Rohan Anil, Martin Wicke, Timothy Dozat, Abhishek Sinha, Enrique Piqueras, Elahe Dabir, Shyam Upadhyay, Anudhyan Boral, Lisa Anne Hendricks, Corey Fry, Josip Djolonga, Yi Su, Jake Walker, Jane Labanowski, Ronny Huang, Vedant Misra, Jeremy Chen, RJ Skerry-Ryan, Avi Singh, Shruti Rijhwani, Dian Yu, Alex Castro-Ros, Beer Changpinyo, Romina Datta, Sumit Bagri, Arnar Mar Hrafnkelsson, Marcello Maggioni, Daniel Zheng, Yury Sulsky, Shaobo Hou, Tom Le Paine, Antoine Yang, Jason Riesa, Dominika Rogozinska, Dror Marcus, Dalia El Badawy, Qiao Zhang, Luyu Wang, Helen Miller, Jeremy Greer, Lars Lowe Sjos, Azade Nova, Heiga Zen, Rahma Chaabouni, Mihaela Rosca, Jiepu Jiang, Charlie Chen, Ruibo Liu, Tara Sainath, Maxim Krikun, Alex Polozov, Jean-Baptiste Lespiau, Josh Newlan, Zeyncep Cankara, Soo Kwak, Yunhan Xu, Phil Chen, Andy Coenen, Clemens Meyer, Katerina Tsihlas, Ada Ma, Juraj Gottweis, Jinwei Xing, Chenjie Gu, Jin Miao, Christian Frank, Zeynep Cankara, Sanjay Ganapathy, Ishita Dasgupta, Steph Hughes-Fitt, Heng Chen, David Reid, Keran Rong, Hongmin Fan, Joost van Amersfoort, Vincent Zhuang, Aaron Cohen, Shixiang Shane Gu, Anhad Mohananey, Anastasija Ilic, Taylor Tobin, John Wieting, Anna Bortsova, Phoebe Thacker, Emma Wang, Emily Caveness, Justin Chiu, Eren Sezener, Alex Kaskasoli, Steven Baker, Katie Millican, Mohamed Elhawaty, Kostas Aisopos, Carl Lebsack, Nathan Byrd, Hanjun Dai, Wenhao Jia, Matthew Wiethoff, Elnaz Davoodi, Albert Weston, Lakshman Yagati, Arun Ahuja, Isabel Gao, Golan Pundak, Susan Zhang, Michael Azzam, Khe Chai Sim, Sergi Caelles, James Keeling, Abhanshu Sharma, Andy Swing, YaGuang Li, Chenxi Liu, Carrie Grimes Bostock, Yamini Bansal, Zachary Nado, Ankesh Anand, Josh Lipschultz, Abhijit Karmarkar, Lev Proleev, Abe Ittycheriah, Soheil Hassas Yeganeh, George Polovets, Aleksandra Faust, Jiao Sun, Alban Rrustemi, Pen Li, Rakesh Shivanna, Jeremiah Liu, Chris Welty, Federico Lebron, Anirudh Baddepudi, Sebastian Krause, Emilio Parisotto, Radu Soricut, Zheng Xu, Dawn Bloxwich, Melvin Johnson, Behnam Neyshabur, Justin Mao-Jones, Renshen Wang, Vinay Ramasesh, Zaheer Abbas, Arthur Guez, Constant Segal, Duc Dung Nguyen, James Svensson, Le Hou, Sarah York, Kieran Milan, Sophie Bridgers, Wiktor Gworek, Marco Tagliasacchi, James Lee-Thorp, Michael Chang, Alexey Guseynov, Ale Jakse Hartman, Michael Kwong, Ruizhe Zhao, Sheleem Kashem, Elizabeth Cole, Antoine Miech, Richard Tanburn, Mary Phuong, Filip Pavetic, Sebastien Cevey, Ramona Comanescu, Richard Ives, Sherry Yang, Cosmo Du, Bo Li, Zizhao Zhang, Mariko Iinuma, Clara Huiyi Hu, Aurko Roy, Shaan Bijwadia, Zhenkai Zhu, Danilo Martins, Rachel Saputro, Anita Gergely, Steven Zheng, Dawei Jia, Ioannis Antonoglou, Adam Sadovsky, Shane Gu, Yingying Bi, Alek Andreev, Sina Samangooei, Mina Khan, Tomas Kocisky, Angelos Filos, Chintu Kumar, Colton Bishop, Adams Yu, Sarah Hodkinson, Sid Mittal, Premal Shah, Alexandre Moufarek, Yong Cheng, Adam Bloniarz, Jaehoon Lee, Pedram Peiman, Paul Michel, Stephen Spencer, Vladimir Feinberg, Xuehan Xiong, Nikolay Savinov, Charlotte Smith, Siamak Shakeri, Dustin Tran, Mary Chesus, Bernd Bohnet, George Tucker, Tamara von Glehn, Carrie Muir, Yiran Mao, Hideto Kazawa, Ambrose Slone, Kedar Soparkar, Disha Shrivastava, James Cobon-Kerr, Michael Sharman, Jay

386

387

388

389

390

391

392

393

394

395

396

397

398

400

401

402

403

404

405

406

407

408

409

410

411

412

413

416

417

418

419

420

421

422

423

424

425

426

427

428

430

431

432

433

434

435

436

437

438

439

440

441

442

443

Payagadhi, Carlos Araya, Karolis Misiunas, Nimesh Ghelani, Michael Laskin, David Barker, 445 Qiujia Li, Anton Briukhov, Neil Houlsby, Mia Glaese, Balaji Lakshminarayanan, Nathan Schucher, 446 Yunhao Tang, Eli Collins, Hyeontaek Lim, Fangxiaoyu Feng, Adria Recasens, Guangda Lai, 447 Alberto Magni, Nicola De Cao, Aditya Siddhant, Zoe Ashwood, Jordi Orbay, Mostafa Dehghani, 448 Jenny Brennan, Yifan He, Kelvin Xu, Yang Gao, Carl Saroufim, James Molloy, Xinyi Wu, Seb 449 Arnold, Solomon Chang, Julian Schrittwieser, Elena Buchatskaya, Soroush Radpour, Martin 450 Polacek, Skye Giordano, Ankur Bapna, Simon Tokumine, Vincent Hellendoorn, Thibault Sottiaux, Sarah Cogan, Aliaksei Severyn, Mohammad Saleh, Shantanu Thakoor, Laurent Shefey, Siyuan 452 Qiao, Meenu Gaba, Shuo yiin Chang, Craig Swanson, Biao Zhang, Benjamin Lee, Paul Kishan 453 Rubenstein, Gan Song, Tom Kwiatkowski, Anna Koop, Ajay Kannan, David Kao, Parker Schuh, 454 Axel Stjerngren, Golnaz Ghiasi, Gena Gibson, Luke Vilnis, Ye Yuan, Felipe Tiengo Ferreira, 455 Aishwarya Kamath, Ted Klimenko, Ken Franko, Kefan Xiao, Indro Bhattacharya, Miteyan Patel, 456 Rui Wang, Alex Morris, Robin Strudel, Vivek Sharma, Peter Choy, Sayed Hadi Hashemi, Jessica 457 Landon, Mara Finkelstein, Priya Jhakra, Justin Frye, Megan Barnes, Matthew Mauger, Dennis 458 Daun, Khuslen Baatarsukh, Matthew Tung, Wael Farhan, Henryk Michalewski, Fabio Viola, Felix 459 de Chaumont Quitry, Charline Le Lan, Tom Hudson, Qingze Wang, Felix Fischer, Ivy Zheng, 460 Elspeth White, Anca Dragan, Jean baptiste Alayrac, Eric Ni, Alexander Pritzel, Adam Iwanicki, 461 Michael Isard, Anna Bulanova, Lukas Zilka, Ethan Dyer, Devendra Sachan, Srivatsan Srinivasan, 462 Hannah Muckenhirn, Honglong Cai, Amol Mandhane, Mukarram Tariq, Jack W. Rae, Gary Wang, 463 Kareem Ayoub, Nicholas FitzGerald, Yao Zhao, Woohyun Han, Chris Alberti, Dan Garrette, 464 Kashyap Krishnakumar, Mai Gimenez, Anselm Levskaya, Daniel Sohn, Josip Matak, Inaki Iturrate, 465 Michael B. Chang, Jackie Xiang, Yuan Cao, Nishant Ranka, Geoff Brown, Adrian Hutter, Vahab 466 Mirrokni, Nanxin Chen, Kaisheng Yao, Zoltan Egyed, François Galilee, Tyler Liechty, Praveen 467 Kallakuri, Evan Palmer, Sanjay Ghemawat, Jasmine Liu, David Tao, Chloe Thornton, Tim Green, 468 Mimi Jasarevic, Sharon Lin, Victor Cotruta, Yi-Xuan Tan, Noah Fiedel, Hongkun Yu, Ed Chi, 469 Alexander Neitz, Jens Heitkaemper, Anu Sinha, Denny Zhou, Yi Sun, Charbel Kaed, Brice Hulse, 470 Swaroop Mishra, Maria Georgaki, Sneha Kudugunta, Clement Farabet, Izhak Shafran, Daniel 471 Vlasic, Anton Tsitsulin, Rajagopal Ananthanarayanan, Alen Carin, Guolong Su, Pei Sun, Shashank 472 V, Gabriel Carvajal, Josef Broder, Iulia Comsa, Alena Repina, William Wong, Warren Weilun Chen, Peter Hawkins, Egor Filonov, Lucia Loher, Christoph Hirnschall, Weiyi Wang, Jingchen Ye, Andrea Burns, Hardie Cate, Diana Gage Wright, Federico Piccinini, Lei Zhang, Chu-Cheng Lin, Ionel 475 Gog, Yana Kulizhskaya, Ashwin Sreevatsa, Shuang Song, Luis C. Cobo, Anand Iyer, Chetan Tekur, 476 Guillermo Garrido, Zhuyun Xiao, Rupert Kemp, Huaixiu Steven Zheng, Hui Li, Ananth Agarwal, 477 Christel Ngani, Kati Goshvadi, Rebeca Santamaria-Fernandez, Wojciech Fica, Xinyun Chen, 478 Chris Gorgolewski, Sean Sun, Roopal Garg, Xinyu Ye, S. M. Ali Eslami, Nan Hua, Jon Simon, 479 Pratik Joshi, Yelin Kim, Ian Tenney, Sahitya Potluri, Lam Nguyen Thiet, Quan Yuan, Florian 480 Luisier, Alexandra Chronopoulou, Salvatore Scellato, Praveen Srinivasan, Minmin Chen, Vinod 482 Koverkathu, Valentin Dalibard, Yaming Xu, Brennan Saeta, Keith Anderson, Thibault Sellam, Nick Fernando, Fantine Huot, Junehyuk Jung, Mani Varadarajan, Michael Quinn, Amit Raul, 483 Maigo Le, Ruslan Habalov, Jon Clark, Komal Jalan, Kalesha Bullard, Achintya Singhal, Thang 484 Luong, Boyu Wang, Sujeevan Rajayogam, Julian Eisenschlos, Johnson Jia, Daniel Finchelstein, 485 Alex Yakubovich, Daniel Balle, Michael Fink, Sameer Agarwal, Jing Li, Dj Dvijotham, Shalini 486 Pal, Kai Kang, Jaclyn Konzelmann, Jennifer Beattie, Olivier Dousse, Diane Wu, Remi Crocker, 487 Chen Elkind, Siddhartha Reddy Jonnalagadda, Jong Lee, Dan Holtmann-Rice, Krystal Kallarackal, 488 Rosanne Liu, Denis Vnukov, Neera Vats, Luca Invernizzi, Mohsen Jafari, Huanjie Zhou, Lilly 489 Taylor, Jennifer Prendki, Marcus Wu, Tom Eccles, Tianqi Liu, Kavya Kopparapu, Françoise 490 Beaufays, Christof Angermueller, Andreea Marzoca, Shourya Sarcar, Hilal Dib, Jeff Stanway, 491 Frank Perbet, Nejc Trdin, Rachel Sterneck, Andrey Khorlin, Dinghua Li, Xihui Wu, Sonam 492 Goenka, David Madras, Sasha Goldshtein, Willi Gierke, Tong Zhou, Yaxin Liu, Yannie Liang, 493 Anais White, Yunjie Li, Shreya Singh, Sanaz Bahargam, Mark Epstein, Sujoy Basu, Li Lao, 494 Adnan Ozturel, Carl Crous, Alex Zhai, Han Lu, Zora Tung, Neeraj Gaur, Alanna Walton, Lucas 495 Dixon, Ming Zhang, Amir Globerson, Grant Uy, Andrew Bolt, Olivia Wiles, Milad Nasr, Ilia 496 Shumailov, Marco Selvi, Francesco Piccinno, Ricardo Aguilar, Sara McCarthy, Misha Khalman, 497 Mrinal Shukla, Vlado Galic, John Carpenter, Kevin Villela, Haibin Zhang, Harry Richardson, 498 James Martens, Matko Bosnjak, Shreyas Rammohan Belle, Jeff Seibert, Mahmoud Alnahlawi, 499 Brian McWilliams, Sankalp Singh, Annie Louis, Wen Ding, Dan Popovici, Lenin Simicich, Laura 500 Knight, Pulkit Mehta, Nishesh Gupta, Chongyang Shi, Saaber Fatehi, Jovana Mitrovic, Alex 501 Grills, Joseph Pagadora, Dessie Petrova, Danielle Eisenbud, Zhishuai Zhang, Damion Yates, 502 Bhavishya Mittal, Nilesh Tripuraneni, Yannis Assael, Thomas Brovelli, Prateek Jain, Mihailo 503

451

Velimirovic, Canfer Akbulut, Jiaqi Mu, Wolfgang Macherey, Ravin Kumar, Jun Xu, Haroon 504 Qureshi, Gheorghe Comanici, Jeremy Wiesner, Zhitao Gong, Anton Ruddock, Matthias Bauer, 505 Nick Felt, Anirudh GP, Anurag Arnab, Dustin Zelle, Jonas Rothfuss, Bill Rosgen, Ashish Shenoy, 506 Bryan Seybold, Xinjian Li, Jayaram Mudigonda, Goker Erdogan, Jiawei Xia, Jiri Simsa, Andrea 507 Michi, Yi Yao, Christopher Yew, Steven Kan, Isaac Caswell, Carey Radebaugh, Andre Elisseeff, 508 Pedro Valenzuela, Kay McKinney, Kim Paterson, Albert Cui, Eri Latorre-Chimoto, Solomon Kim, 509 William Zeng, Ken Durden, Priya Ponnapalli, Tiberiu Sosea, Christopher A. Choquette-Choo, 510 James Manyika, Brona Robenek, Harsha Vashisht, Sebastien Pereira, Hoi Lam, Marko Velic, 511 Denese Owusu-Afriyie, Katherine Lee, Tolga Bolukbasi, Alicia Parrish, Shawn Lu, Jane Park, 512 Balaji Venkatraman, Alice Talbert, Lambert Rosique, Yuchung Cheng, Andrei Sozanschi, Adam 513 Paszke, Praveen Kumar, Jessica Austin, Lu Li, Khalid Salama, Wooyeol Kim, Nandita Dukkipati, 514 Anthony Baryshnikov, Christos Kaplanis, XiangHai Sheng, Yuri Chervonyi, Caglar Unlu, Diego 515 de Las Casas, Harry Askham, Kathryn Tunyasuvunakool, Felix Gimeno, Siim Poder, Chester 516 Kwak, Matt Miecnikowski, Vahab Mirrokni, Alek Dimitriev, Aaron Parisi, Dangyi Liu, Tomy Tsai, Toby Shevlane, Christina Kouridi, Drew Garmon, Adrian Goedeckemeyer, Adam R. Brown, 518 Anitha Vijayakumar, Ali Elqursh, Sadegh Jazayeri, Jin Huang, Sara Mc Carthy, Jay Hoover, 519 Lucy Kim, Sandeep Kumar, Wei Chen, Courtney Biles, Garrett Bingham, Evan Rosen, Lisa 520 Wang, Qijun Tan, David Engel, Francesco Pongetti, Dario de Cesare, Dongseong Hwang, Lily 521 Yu, Jennifer Pullman, Srini Narayanan, Kyle Levin, Siddharth Gopal, Megan Li, Asaf Aharoni, 522 Trieu Trinh, Jessica Lo, Norman Casagrande, Roopali Vij, Loic Matthey, Bramandia Ramadhana, 523 Austin Matthews, CJ Carey, Matthew Johnson, Kremena Goranova, Rohin Shah, Shereen Ashraf, 524 Kingshuk Dasgupta, Rasmus Larsen, Yicheng Wang, Manish Reddy Vuyyuru, Chong Jiang, Joana 525 Ijazi, Kazuki Osawa, Celine Smith, Ramya Sree Boppana, Taylan Bilal, Yuma Koizumi, Ying 526 Xu, Yasemin Altun, Nir Shabat, Ben Bariach, Alex Korchemniy, Kiam Choo, Olaf Ronneberger, 527 Chimezie Iwuanyanwu, Shubin Zhao, David Soergel, Cho-Jui Hsieh, Irene Cai, Shariq Iqbal, 528 Martin Sundermeyer, Zhe Chen, Elie Bursztein, Chaitanya Malaviya, Fadi Biadsy, Prakash Shroff, 529 Inderjit Dhillon, Tejasi Latkar, Chris Dyer, Hannah Forbes, Massimo Nicosia, Vitaly Nikolaev, 530 Somer Greene, Marin Georgiev, Pidong Wang, Nina Martin, Hanie Sedghi, John Zhang, Praseem 531 Banzal, Doug Fritz, Vikram Rao, Xuezhi Wang, Jiageng Zhang, Viorica Patraucean, Dayou Du, Igor Mordatch, Ivan Jurin, Lewis Liu, Ayush Dubey, Abhi Mohan, Janek Nowakowski, Vlad-Doru 533 Ion, Nan Wei, Reiko Tojo, Maria Abi Raad, Drew A. Hudson, Vaishakh Keshava, Shubham 534 Agrawal, Kevin Ramirez, Zhichun Wu, Hoang Nguyen, Ji Liu, Madhavi Sewak, Bryce Petrini, 535 DongHyun Choi, Ivan Philips, Ziyue Wang, Ioana Bica, Ankush Garg, Jarek Wilkiewicz, Priyanka 536 Agrawal, Xiaowei Li, Danhao Guo, Emily Xue, Naseer Shaik, Andrew Leach, Sadh MNM Khan, 537 Julia Wiesinger, Sammy Jerome, Abhishek Chakladar, Alek Wenjiao Wang, Tina Ornduff, Folake 538 Abu, Alireza Ghaffarkhah, Marcus Wainwright, Mario Cortes, Frederick Liu, Joshua Maynez, 539 Andreas Terzis, Pouya Samangouei, Riham Mansour, Tomasz Kepa, François-Xavier Aubet, Anton 540 541 Algymr, Dan Banica, Agoston Weisz, Andras Orban, Alexandre Senges, Ewa Andrejczuk, Mark Geller, Niccolo Dal Santo, Valentin Anklin, Majd Al Merey, Martin Baeuml, Trevor Strohman, 542 Junwen Bai, Slav Petrov, Yonghui Wu, Demis Hassabis, Koray Kavukcuoglu, Jeffrey Dean, and 543 Oriol Vinyals. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of 544 context, 2024. URL https://arxiv.org/abs/2403.05530. 545

Alexander Havrilla, Sharath Chandra Raparthy, Christoforos Nalmpantis, Jane Dwivedi-Yu, Maksym Zhuravinskyi, Eric Hambro, and Roberta Raileanu. Glore: When, where, and how to improve llm reasoning via global and local refinements. In *Forty-first International Conference on Machine Learning*, 2024.

Yue Huang, Jiawen Shi, Yuan Li, Chenrui Fan, Siyuan Wu, Qihui Zhang, Yixin Liu, Pan Zhou, Yao
Wan, Neil Zhenqiang Gong, and Lichao Sun. Metatool benchmark for large language models:
Deciding whether to use tools and which to use, 2024. URL https://arxiv.org/abs/2310.
03128.

Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik
 Narasimhan. Swe-bench: Can language models resolve real-world github issues?, 2024. URL
 https://arxiv.org/abs/2310.06770.

Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension, 2017. URL https://arxiv.org/abs/1705.03551.

- Jungo Kasai, Keisuke Sakaguchi, Yoichi Takahashi, Ronan Le Bras, Akari Asai, Xinyan Yu, Dragomir 560 Radey, Noah A. Smith, Yejin Choi, and Kentaro Inui. Realtime qa: What's the answer right now?, 561 2024. URL https://arxiv.org/abs/2207.13332. 562
- Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, 563 Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. Rewardbench: Evaluating reward models 564 for language modeling. arXiv preprint arXiv:2403.13787, 2024. 565
- Minghao Li, Yingxiu Zhao, Bowen Yu, Feifan Song, Hangyu Li, Haiyang Yu, Zhoujun Li, Fei Huang, 566 567 and Yongbin Li. Api-bank: A comprehensive benchmark for tool-augmented llms, 2023. URL https://arxiv.org/abs/2304.08244. 568
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan 569 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step, 2023. URL 570 https://arxiv.org/abs/2305.20050. 571
- Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, 572 Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang. 574 Agentbench: Evaluating Ilms as agents, 2023. URL https://arxiv.org/abs/2308.03688. 575
- Martin Majlis. Wikipedia-api, 2017. URL https://github.com/martin-majlis/ 576 Wikipedia-API/tree/master. 577
- Dheeraj Mekala, Jason Weston, Jack Lanchantin, Roberta Raileanu, Maria Lomeli, Jingbo Shang, and 578 Jane Dwivedi-Yu. Toolverifier: Generalization to new tools via self-verification. arXiv preprint 579 arXiv:2402.14158, 2024. 580
- Grégoire Mialon, Clémentine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom. 581 Gaia: a benchmark for general ai assistants, 2023. URL https://arxiv.org/abs/2311.12983. 582
- Mistral. Mistral large 2 model. https://mistral.ai. Accessed: 2024-10-01. 583
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni 584 Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor 585 Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, 586 Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny 587 Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, 588 Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea 589 Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, 590 Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, 591 Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, 592 Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty 593 Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, 594 Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel 595 Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua 596 Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon 598 Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne 599 Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo 600 Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, 601 Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik 602 Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, 603 Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy 604 Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie 605 Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, 606 Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, 607 Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David 608 Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie 609 Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, 610 Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo 611 Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano,

- Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Payloy, Andrew Peng, 613
- Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, 614
- Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, 615
- Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis 616
- Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted 617
- Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel 618
- Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon 619
- Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, 620
- Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie 621
- Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, 622
- Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun 623
- Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang,
- Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian 625
- Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren
- Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming
- Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao
- 628
- Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL 629
- https://arxiv.org/abs/2303.08774. 630
- Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. Gorilla: Large language model 631 connected with massive apis, 2023. URL https://arxiv.org/abs/2305.15334. 632
- Yun Peng, Shuqing Li, Wenwei Gu, Yichen Li, Wenxuan Wang, Cuiyun Gao, and Michael Lyu. 633 Revisiting, benchmarking and exploring api recommendation: How far are we?, 2021. 634
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru 635
- Tang, Bill Qian, Sihan Zhao, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, 636
- Zhiyuan Liu, and Maosong Sun. Toolllm: Facilitating large language models to master 16000+ 637
- real-world apis, 2023. 638
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, 639 Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to 640
- use tools, 2023. URL https://arxiv.org/abs/2302.04761. 641
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy 642 optimization algorithms, 2017. URL https://arxiv.org/abs/1707.06347. 643
- SerpApi. Serpapi search engine results api. https://serpapi.com/, 2024. Accessed: February 2023-September 2024. 645
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. Fever: a large-scale 646 dataset for fact extraction and verification, 2018. URL https://arxiv.org/abs/1803.05355. 647
- Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, 648 Denny Zhou, Quoc Le, and Thang Luong. Freshllms: Refreshing large language models with 649 search engine augmentation, 2023. URL https://arxiv.org/abs/2310.03214. 650
- Peiyi Wang, Lei Li, Zhihong Shao, R. X. Xu, Damai Dai, Yifei Li, Deli Chen, Y. Wu, and Zhifang 651 Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations, 2024a. 652 URL https://arxiv.org/abs/2312.08935. 653
- Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji. 654 Executable code actions elicit better llm agents, 2024b. URL https://arxiv.org/abs/2402. 655 656 01030.
- Inc. Wolfram Research. Mathematica, Version 14.1, 2024. URL https://www.wolfram.com/ 657 mathematica. Champaign, IL, 2024. 658
- Congying Xia, Chen Xing, Jiangshu Du, Xinyi Yang, Yihao Feng, Ran Xu, Wenpeng Yin, and 659 Caiming Xiong. Fofo: A benchmark to evaluate llms' format-following capability, 2024. URL 660
- https://arxiv.org/abs/2402.18667.

- Qiantong Xu, Fenglu Hong, Bo Li, Changran Hu, Zhengyu Chen, and Jian Zhang. On the tool
   manipulation capability of open-source large language models, 2023.
- Fanjia Yan, Huanzhi Mao, Charlie Cheng-Jie Ji, Tianjun Zhang, Shishir G. Patil, Ion Stoica, and
   Joseph E. Gonzalez. Berkeley function calling leaderboard. https://gorilla.cs.berkeley.
   edu/blogs/8\_berkeley\_function\_calling\_leaderboard.html, 2024.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov,
   and Christopher D. Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question
   answering, 2018. URL https://arxiv.org/abs/1809.09600.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
   React: Synergizing reasoning and acting in language models, 2023. URL https://arxiv.org/abs/2210.03629.
- Michael J. Q. Zhang and Eunsol Choi. Situatedqa: Incorporating extra-linguistic contexts into qa, 2021. URL https://arxiv.org/abs/2109.06157.
- Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and Chao Zhang. Toolqa: A dataset for llm
   question answering with external tools. *Advances in Neural Information Processing Systems*, 36:
   50117–50143, 2023.
- Patrick Zippenfenig. Open-Meteo.com Weather API, 2024. URL https://github.com/open-meteo/open-meteo.

# **A ToolComp Extended Evaluations**

This appendix presents comprehensive supplementary evaluations that provide deeper insights into model performance and failure modes. We include exact match grading analysis (A.1) and detailed error categorization for each evaluated model (A.2 and A.3). Note that the frontier model evaluations presented here reflect the state-of-the-art as of August 2024, while the main text incorporates more recent frontier models released through August 2025. These extended evaluation methodologies offer model developers actionable frameworks for conducting thorough assessments of their systems' tool composition capabilities and identifying specific areas for improvement.

## 688 A.1 Exact Match

This paradigm aims to assess both the tool use capabilities and the instruction/format following capabilities of the model. Formatting is particularly important when we want to use the LLM to automate a backend process. This paradigm programmatically evaluates unsorted lists (eg. prompt asks for a list of all states in the US), sorted lists (eg. prompt asks for a list of all states in the US in alphabetical order), numbers (eg. prompt asks for the areas of Texas in square miles) and strings (eg. prompt asks for the name of the football team that won the Superbowl in 2016)

Unsorted lists are sorted and exact matched (set match gets rid of duplicates) Sorted lists are exact matched Number are checked if they are within a tolerance param (the tolerance param is to account for variance among different sources online) String are stripped, lower cased, and exact matched

Table 4: Model Family Performance Comparison: Accuracy and 95% Confidence Intervals

Model Family	Model Name	Total Accuracy (%)
OpenAI	GPT-40 (Aug 2024) GPT-40 (May 2024) GPT-4 Turbo Preview GPT-4 GPT-40 Mini	$43.52 \pm 4.43$ $40.60 \pm 4.38$ $40.11 \pm 4.39$ $38.45 \pm 4.34$ $34.70 \pm 4.25$
Anthropic	Claude 3.5 Sonnet Claude 3 Opus Claude 3 Sonnet	$42.92 \pm 4.42  36.96 \pm 4.43  33.58 \pm 4.21$
Google	Gemini 1.5 Pro (August 27, 2024) Gemini 1.5 Pro (May 2024)	$43.22 \pm 4.43$ $27.36 \pm 3.98$
Mistral	Mistral Large 2	$33.63 \pm 4.21$
Meta	Llama 3.1 405B Instruct* Llama 3.1 70B Instruct* Llama 3.1 8B Instruct*	$33.10 \pm 4.20$ $26.19 \pm 3.93$ $11.75 \pm 2.88$
Cohere	Command R+	$0.00 \pm 0.00$

# A.2 Final Answer Failure Analysis

699

701

702

In order to better understand the reasons behind each model's failures, we come up with an Error Taxonomy and use GPT-4 Turbo to categorize the reasoning behind each failure. We note that the error categories are not mutually exclusive. We inspect the individual failure cases predicted by GPT-4 Turbo and find that it is reasonably accurate. The different categories and their definitions are shown in Table 5 and the error counts for each model is shown in Figure 3.

Table 5: Common Error Category Taxonomy.

Category	Description
Final Answer Missing Information	The model's trajectory got to the final answer however the final answer fails to answer all parts of the prompt.
Called Incorrect Tool	The model called irrelevant tools that lead it down the wrong direction.
Incorrect Tool Call Formatting	The model tried to call the relevant tool but consistently used the wrong formatting for the input arguments (e.g., wrong input format, didn't include a required argument). You can tell this is occurring if the tool call's result is an error message.
Terminated Early Unexpectedly	The model stopped short of reaching the final answer even though it should have kept proceeding. It is unclear why the model stopped early.
Hallucinated Information	The model either didn't call the relevant tool and just made up information or it called the relevant tool but didn't use its outputs in the next tool call or final answer properly (made up information afterwards).
Misunderstood Tool Info	The model called the relevant tool but misunderstood the information it gave back.
Repeatedly Calling Same Tool	The model called the same tool with the same arguments multiple times (even though it didn't have any errors) and didn't use the returned info to proceed to the next step or the final answer.
Action Plan Flawed	The Action Plan provided to the model in the user query was fundamentally flawed.
Miscellaneous	The reason for the error doesn't fit into any of the above categories.

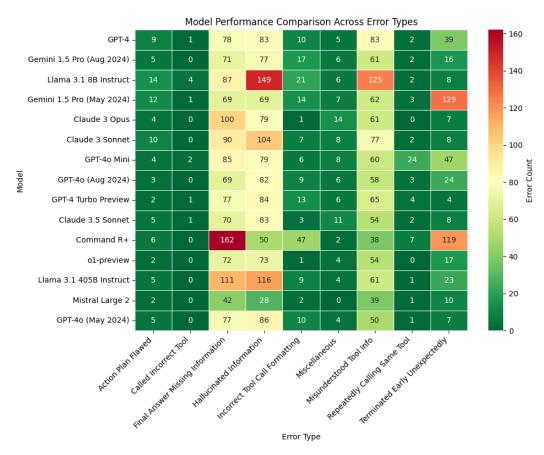


Figure 3: Breakdown of the various error categories in our taxonomy for each model (on the ToolComp-Enterprise).

#### 704 A.3 Intermediate Reasoning Failure Analysis

707

712

713

714

715

716

In this appendix section, we conduct a thorough failure analysis for the intermediate reasoning evaluations shown in Table 3.

# A.3.1 ReAct-Step-Error-Based Failure Trends in Models

Figures 4 and 5 shows the count for type of mistake between the human corrected substep and the original incorrect substep whenever the model fails to pick the more appropriate trajectory (see Figure 1 for an overview on the annotation process). We define the failure cases in terms of which subset of the ReAct step needed correction. We end up with 5 different cases:

- Case 1: Thought Correct, Action Correct, Action Input Incorrect
- Case 2: Thought Incorrect, Action Incorrect, Action Input Incorrect
- Case 3: Thought Incorrect, Action Correct, Action Input Correct
- Case 4: Thought Incorrect, Action Correct, Action Input Incorrect
- Case 5: Thought Correct, Action Incorrect, Action Input Incorrect

Together, these figures highlight what types of errors are most common during a lapse in reasoning when picking the best next course of action or invoking a tool correctly. In particular, we notice that models often fail in reasoning about the better course of action when the deciding factor is in picking the better Action Input with all else equal.

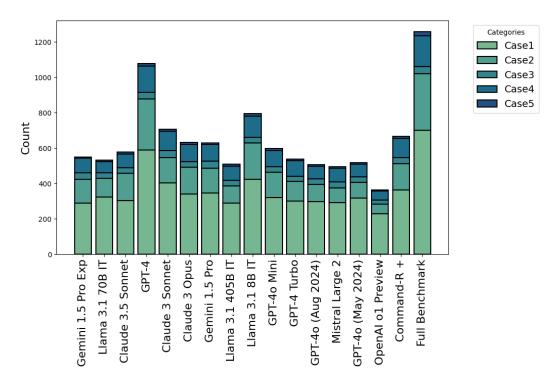


Figure 4: Histogram showing the LLM as judge evaluation failure counts for each model, which is further categorized by subset of the ReAct step that needed correction. Full Benchmark denotes the counts for the entire ToolComp benchmark. Recall from 4.3, we have 3 outcomes for LLM judge evaluation: win, tie, or loss. Here we count a failure as either a tie or a loss outcome.

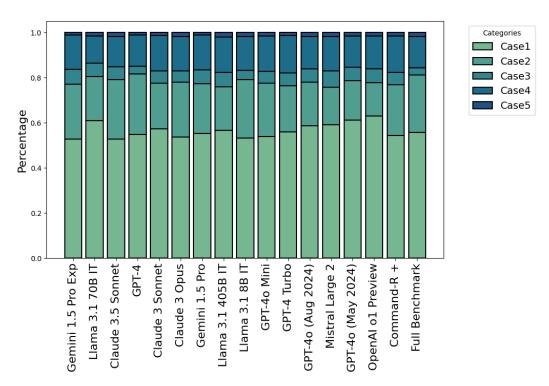


Figure 5: Density of the error-type between correct and incorrect step for the LLM as judge evaluation failures for each model. Full Benchmark denotes the distribution for the entire ToolComp benchmark.

#### A.3.2 Position-Based Error Trends in Models

Figures 6 and 7 shows the count and percentage of the relative positions where each respective model failed to chose the better step when serving as an LLM judge choosing between two steps. In order to calculate the position, we divide the step number at which the decision is taking place by the total number of steps in the trajectory and multiply by 100. Hence, the position of a step will be a number between 0 and 100. We bin these position values by increments of 20. Overall, these figures illustrate that most, if not, all of the models struggle when judging steps towards the middle-end (position values between 60 and 80) of the trajectory. Intuitively this makes sense because this is likely where models have to compose the observations of previous tools into the input for the next tool call, which requires more nuanced and sophisticated reasoning.

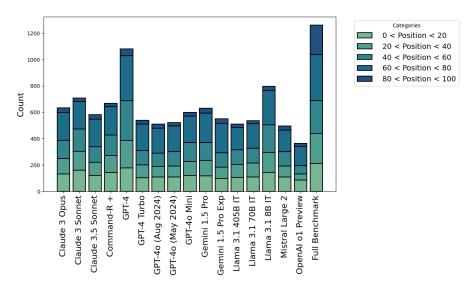


Figure 6: Histogram showing the LLM as judge evaluation failure counts for each model, which is further categorized by the position of the decision step. Full Benchmark denotes the counts for the entire ToolComp benchmark.

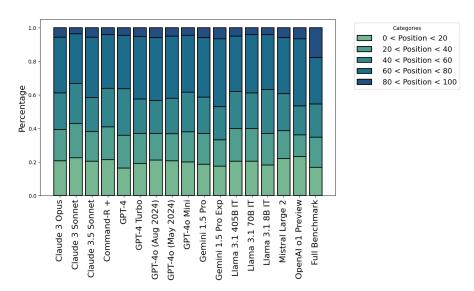


Figure 7: Density of the position of the LLM as judge evaluation failures for each model. Full Benchmark denotes the distribution for the entire ToolComp benchmark.

# 731 B ToolComp Details

- In this appendix section, we provide further details regarding benchmark creation steps such as prompt creation (B.1, B.2, B.3). We also provide additional benchmark metadata revolving different characteristics and statistics about the benchmark (B.4).
- 735 B.1 Prompt Creation Details
- Step 1: Develop In-Context Examples We crafted high-quality in-context (IC) examples with supporting reasoning, which we call 'processes', to guide the prompt generation. These processes are Chain of Thought reasonings that describe the process by which we came up with the prompt. One of the IC Prompts and a corresponding CoT is shown in Appendix B.2
- Step 2: Generate Initial Prompts Using the IC examples, we generated synthetic prompts, ensuring diversity by selecting random subsets of IC examples. Each subset used distinct in-context prompts and randomly sampled tools from its set of available tools. The seed prompt used in this step in Appendix B.3.
- Step 3: Filtering We manually inspected each prompt to ensure they were reasonable, interesting, and challenging, labeling them as Good, Too Simple, or Nonsensical with justifications for each classification. These labeled examples served as IC inputs for GPT-4 Turbo (OpenAI et al., 2024) to classify additional prompts. We iteratively review the outputs, make necessary edits, and add more IC examples. Through three iterations, the filtered prompts were of high quality, exhibiting only minor mistakes.
- Step 4: Human Refinement After filtering, annotators reviewed the finals prompts to resolve any issues related to complexity, clarity and ambiguity. We gave clear instructions on ambiguity (only one possible correct answer) and complexity (requires two or more tool calls to answer), instructing our annotators to ensure the prompt has only one correct answer that is complex, challenging and requires the use of tools.

### 55 B.2 In Context Example

### 756 Prompt

I wanna know if eating meat is correlated with heart issues, find the annual per capita consumption of meat in (kg/person) and also the per capita heart attack rates (in heart attacks/person) for every country. Then run a linear regression with y as heart attack rates and x as meat consumption, return the Pearson's correlation as well as the slope of the fit line.

# 758 Process

757

I will first start by creating a prompt that requires the use of google search. I want to make this prompt about investigating whether the amount of meat you consume is correlated to heart disease. In order to make sure there is only one possible answer, I will ask to find the per capita consumption of meat (in kg/person) and heart attacks rates (heart attacks per person) in all countries. This standardizes the actual data that needs to be pulled and specifies the units to ensure there is only one possible answer. I will then ask for a linear regression using that data since it requires a python interpreter. Since linear regression is deterministic when the data is fixed and the data required to fit the linear regression is well defined, I can ask to output its parameters and ensure there is only one possible answer that can be returned. This ensures that the good prompt is clear, unambiguous and has an answer that is easy to verify through an exact string match while also requiring a chain of dependent tool calls (google search call, then python interpreter call) to solve.

# 760 B.3 Seed Prompt

761

I want you to act as a Prompt Writer.

Please adhere to the following instructions:

- Write a prompt that requires the use of all of the tools.
- The prompt should require a chain of dependent tools calls who's outputs influence the inputs of the next tool invocation.
- The prompt should be appropriate for someone in {grade}.
- Please do not specify the tools to be used in the prompt. We want the assistant to figure out on it's own what tools to call so it should not be specified in the prompt itself. No phrases like "Use the ... tool" should be in the written prompt.
- The prompt should be a couple sentences.
- Make sure the prompt has only one possible answer that is concrete and easily verifiable. We want to be able to check the final answer using exact match.
- Make sure the answer is not in the prompt.
- Place [STOP] at the end of the prompt.

Examples:

{examples}

[BEGIN ALLOWED TOOLS]

{tools}

[END ALLOWED TOOLS]

762

763

### B.4 Benchmark Metadata

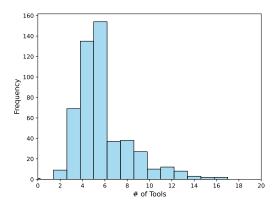


Figure 8: About 85% of prompts in ToolComp require at least 3 tool calls to solve, indicating that they have a decent amount of complexity and difficulty. Furthermore, 20% of prompts still require 7 or more tool calls to solve. This indicates that an agent being evaluated on this benchmark requires high context length, sophisticated reasoning over long context, and advanced tool calling capabilities in order to process long tool chains, formulate a high level plan, and understand the outputs of each tool call to proceed to the next step and subsequently achieve a high score.

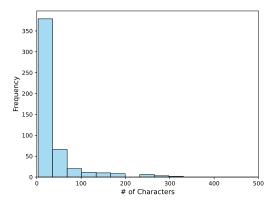


Figure 9: Due to the nature of ToolComp needing to have answers that are easily verifiable, we choose to create prompts that have numbers and short strings to match. However, there are still some examples of prompts that require long structured outputs such as dictionaries, tuples and lists. These test the agent's ability to follow complex queries that involve returning long outputs such as lists or dictionaries of city names, temperatures, altitudes, etc.

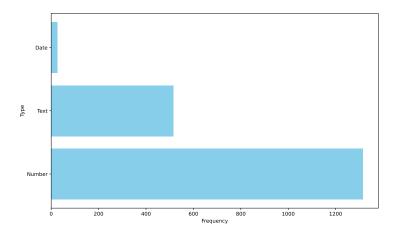


Figure 10: We show the distribution of the following primitive data types: number, string and date. We care most about evaluation of compositional tool use and reasoning rather than aesthetic output structuring and formatting. This is why the benchmark's labels are predominantly numeric while containing a significant fraction of string outputs. In many cases, strings and names are intermediary outputs, but we most often ask for numerical final answers to make the answer easier to unambiguously verify.

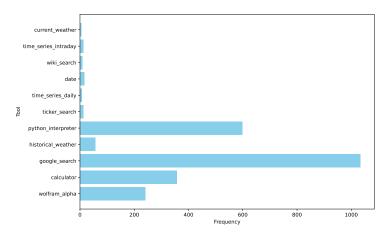


Figure 11: The distribution of tools called in our human supervised tool call chains. The heavy bias towards Google and Python are due to ToolComp Chat only allowing these tools as well them being generally applicable for a wide range of tasks (web retrieval and information processing).

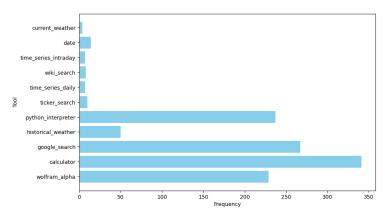


Figure 12: The distribution of tools called in our human supervised tool call chains for just the ToolComp Enterprise subset.

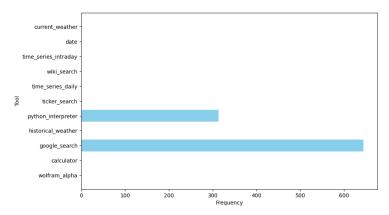


Figure 13: The distribution of tools called in our human supervised tool call chains for just the ToolComp Chat subset.

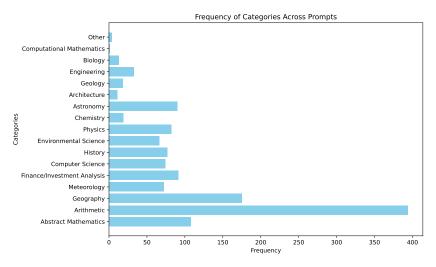


Figure 14: Here, we show the various topics our prompts address. Many prompts require arithmetic operations and mathematical reasoning along with a somewhat uniform distribution of multiple disciplines ranging from Geography, Finance, History, Physics, Chemistry, Astronomy, Architecture etc. The topics are not mutually exclusive since many of these prompts span multiple domains and require multiple tools, multiple sources of knowledge and diverse forms of reasoning.

# 64 C Tool-Use Prompts

In this section, we summarize all of the prompts that were used during the creation of the benchmark, evaluation of the benchmark, and creation of the synthetic training data. For the creation of the benchmark, we state the "Action Plan Prompt" for the Policy Model in Section C.1 and the "Tool Call Prompt" for the Policy Model in Section C.2. For the evaluation of the benchmark, we state the LLM grading prompt and the in-context examples used to aid grading in Section C.5. Lastly, for the creation of the synthetic training data, we use the same policy model prompts for the action plan and tool call, and we additionally include the "Action Plan Prompt" for the Critic Model in Section C.3 and the "Tool Call Prompt" for the Critic Model in Section C.4.

# 773 C.1 Action Plan Prompt (Policy Model)

You are a helpful action planner with access to functions. Please use the tools to provide information accurate up to current date: {current\_date}

FUNCTIONS: {func\_spec}

Question: {question}

Given the tools available to you above, please formulate an action plan to answer the question in a bulleted list for each step. Refrain from using any specific tool calls in your action plan, instead focus on the high-level steps you would take to answer the question and the name of the tool you would use and how you would use it. Refrain from trying to answer the question directly in the action plan.

765

766

767

768

769

770

# 775 C.2 ReAct Tool Call Prompt (Policy Model)

#### SYSTEM:

776

You are a helpful assistant with access to functions, each function will be regarded as an action. Your job is to take relevant and necessary actions to get to the final answer to a user question. Please use the actions to provide information accurate up to current date and time: {current\_date}. The user will provide you a question and a high level action plan. Your job is to execute on the action plan to answer the question. It's okay to slightly deviate from the action plan if you think it's necessary.

FUNCTIONS: {func\_spec}

Please stick to the following format:

Thought: \(\langle\) your reasoning/thought on why/how to use an action\(\rangle\)

Action: (the action to take, should be one of {func\_list})

Action Input: (the input to the action (should be in JSON format with the required fields))

**End Action** 

If you believe that you have obtained enough information (which can be judged from the history observations) to answer the question, please call:

Thought: I have enough information to answer the question

Action: finish

Action Input: {"answer": [your answer string]}}

**End Action** 

For your final answer (the finish action input), make sure you answer the full question. Additionally, we want to make sure the final answers/outputs in the finish action input are returned in the order that they are given in a list format so we can verify them with an exact string match. For eg. if the prompt asks for a city name, its temperature and a list of names of all the NBA teams whose home stadium is within a 400 mile radius, you would output ['San Francisco', 78, ['Los Angeles Lakers', 'Golden State Warriors']].

If the prompt asks for a special sorting of the list, make sure to output wrap the list in {{}} and if doesn't require any special sorting wrap it in [] like you normally would. So if the prompt instead asked to list the names of all the NBA teams whose home stadium is within a 400 mile radius in alphabetical order, you would output [San Francisco, 78, {{Golden State Warriors, Los Angeles Lakers}}].

Only output the final answer with no additional text or natural language. Give dates in YYYY-MM-DD format, temperatures in celcius, prices in dollars, lengths in meters, area in meters<sup>2</sup>, volume in  $m^3$  and angles in degrees if the prompt doesn't specify what format/units to output the answer in.

Given a user provided question and action plan, as well as your previous actions and observations, take your next action.

USER:

Question: {question}

Action Plan: {action\_plan}

ASSISTANT:

{history\_of\_react\_steps}

# 778 C.3 Action Plan Prompt (Critic Model)

You are an expert planner of tool calls. Your job is to critique the action plan of an assistant.

The following information is shown to the assistant in order to devise an action plan:

[Start of the message]

You are a helpful assistant with access to functions. Please use the tools to provide information accurate up to current date and time: {current\_date}.

FUNCTIONS: {func\_spec}

Question: {question}

Given the question and the tools available to you above, please formulate an action plan to answer the question in a bulleted list for each step.

Refrain from using any specific tool calls in your action plan, instead focus on the high-level steps you would take to answer the question and the name of the tool you would use and how you would use it. Refrain from trying to answer the question directly in the action plan.

[End of the message]

Given the set of functions and the question, please critique the action plan provided by the assistant.

First, determine if the action plan is correct or incorrect. To do so, provide a reasoning and then label the action plan as correct or incorrect. In order to determine if the action plan needs revision, consider the following:

- Is the action plan reasonable given the set of functions available?
- Is the action plan clear and concise?
- Is the action plan missing any steps?

Please err on the side of giving the assistant the benefit of the doubt, and only critique the action plan if it is clearly incorrect.

If the action plan is incorrect, provide a revised action plan that you believe would be correct.

Furthermore, your output should follow the format:

Reasoning: \( \text{ your reasoning for the correctness or incorrectness of the action plan } \)

Label: \( \text{correct/incorrect} \)

Revised Action Plan: (your revised action plan or empty if no revision needed)

Here is the action plan provided by the assistant:

{action\_plan}

Please provide your critique of the action plan.

# 780 C.4 ReAct Tool Call Prompt (Critic Model)

You are an expert judge of tool calls. Your job is to critique each of the ReAct steps of an assistant.

The following information is shown to the assistant in order to devise a ReAct step.

[Start of the message]

781

You are a helpful assistant with access to functions. Use them if required. Please use the tools to provide information accurate up to current date and time: {current\_date}.

FUNCTIONS: {func\_spec}

Please stick to the following format:

Thought: you should always think about what to do Action: the action to take, should be one of {func\_list}

Action Input: the input to the action

**End Action** 

If you believe that you have obtained enough information (which can be judged from the history observations) to answer the question, please call:

Thought: I have enough information to answer the question

Action: finish

Action Input: "answer": [your answer string]

**End Action** 

Question: {question}

[End of the message]

Given the set of functions, question, action plan and history of past actions, critique the Thought, Action, and Action Input step. Assume the action plan and history of past actions are optimal. To assess the thought step, if the step is roughly reasonable and the action and action input step are correlated with the thought step, then the thought step is correct. Please give the assistant the benefit of the doubt and be lenient in your assessment.

To assess the action step, let's assume that the Assistant cannot complete simple functionalities such as simple arithmetic, converting units, or utilizing simple facts without the use of tools. If the action specifies a reasonable function to use, then the action step is correct.

To assess the action input step, if the input is reasonable and the action is correct, then the action input step is correct.

If any of the steps are incorrect, label them as incorrect in the Labels section.

For the Revised ReAct Step section, provide the correct step that the assistant should have taken. If the assistant's step is correct, provide the assistant's step as the revised step. If the assistant's step is incorrect, provide the correct step that the assistant should have taken. As a general rule of thumb, if your revised step is different from the assistant's step, then the assistant's step is incorrect, and if your revised step is the same as the assistant's step, then the assistant's step is correct.

As an important reminder, for your final answer (the finish action input), we want to make sure the final answers/outputs in the finish action input are returned in the order that they

are given in a list format so we can verify them with an exact string match. For eg. if the prompt asks for a city name, its temperature and a list of names of all the NBA teams whose home stadium is within a 400 mile radius, you would output ['San Francisco', 78, ['Los Angeles Lakers', 'Golden State Warriors']]. If the prompt asks for a special sorting of the list, make sure to output wrap the list in {{}} and if doesn't require any special sorting wrap it in [] like you normally would. So if the prompt instead asked to list the names of all the NBA teams whose home stadium is within a 400 mile radius in alphabetical order, you would output [San Francisco, 78, {{Golden State Warriors, Los Angeles Lakers}}].

Only output the final answer with no additional text or natural language or units. Give dates in YYYY-MM-DD format, temperatures in Celcius, prices in dollars, lengths in meters, area in meters<sup>2</sup>, volume in  $m^3$  and angles in degrees if the prompt doesn't specify what format/units to output the answer in.

As a reminder, you should not use an external information that is not provided in the prompt or by a tool call. As a simple example, you may know a ticker symbol already for a company, but you should not use it unless you have called the ticker\_search or a similar function (e.g. google\_search, wiki\_search, etc.) to retrieve that information.

Your output should follow the format:

```
[Start of format]
```

Reasoning: \( \text{ your reasoning for the correctness or incorrectness of each step } \)

Labels:  $[\langle correct/incorrect \rangle, \langle correct/incorrect \rangle]$  (in the order of Thought, Action, Action Input)

Revised ReAct Step:

```
Thought: \langle your revised thought or assistant's thought if correct \rangle
```

Action: \( \text{ your revised action or assistant's action if correct } \)

Action Input: \( \) your revised action input or assistant's action input if correct \( \)

End Action

[End of format]

Here is the action plan:

{action\_plan}

Here is the history of past actions. If there are no past actions yet, this will be empty:

{history}

Here is the latest ReAct step provided by the assistant:

Thought: {thought}
Action: {action}

Action Input: {action input}

**End Action** 

Observation: {observation}

Please provide your critique of the latest ReAct step provided by the assistant.

# 784 C.5 LLM Grading Prompt

#### 85 C.5.1 Main Prompt

786

You are an expert test grader. You have been given a student answer ('Student Answer:') to grade. You have also been the correct answer ('Correct Answer:') and the original question ('Question:'). Each correct answer is a list of strings.

# {In-Context Examples}

The possible grades are

**INCORRECT**: 'Student Answer:' is different from 'Correct Answer:'

- numbers are completely different
- lists are completely different
- 'Question:' asks for special sorting of a list but the list in 'Student Answer:' is sorted differently than 'Correct Answer:'
- strings are completely different or information present in the string is completely different

**CORRECT BUT BAD FORMATTING**: 'Student Answer:' has the same info as 'Correct Answer:' but is formatted differently.

- 'Student Answer:' includes natural language or additional text
- numbers are formatted differently but they are close to one another ('Student Answer:' is within
- lists are wrapped differently than the correct answer but contains the same information and sorted the same way as 'Correct Answer:' if asked 'Question:' asks for a special sorting
- Strings are the same but may be formatted differently

**CORRECT**: The student answer has the same info as 'Correct Answer:' and is also formatted the same as 'Correct Answer:'

- numbers are close to one another ('Student Answer:' is within 10% of the correct answer)
- if 'Question:' asks for a special sorting of the list the 'Student Answer:' list is sort the same as 'Correct Answer:'
- lists are wrapped the same
- Strings are identical

Remember you are assuming the correct answer provided is correct, your job is is only to compare the correct answer to the student answer and grade it based on the above criteria. Do not try to determine the correct answer yourself. Make sure to include a reasoning and final grade in the format:

Reasoning: \( \text{reasoning} \) Final Grade: \( \text{INCORRECT / CORRECT BUT BAD FORMATTING / CORRECT } \) [ENDOFGRADE]

Now do this for the following user provided question, student answer and correct answer.

# 788 C.5.2 In-Context Examples (Ordering)

We want to make sure the values in the student answer are returned in the order that they are asked in 'Question:'.

789

For example, if 'Question:' asks for a city name, its temperature and a list of names of all the NBA teams whose home stadium is within a 400 mile radius, and 'Correct Answer:' is ['San Francisco', 78, ['Golden State Warriors', 'Los Angeles Lakers']] we would want 'Student Answer:' can be ['San Francisco', 78, ['Los Angeles Lakers', 'Golden State Warriors']].

# Examples:

**Question**: Find the name of the city known for its famous tourist attraction Alcatraz, also give it's current temperature and a list of names of all the NBA teams whose home stadium is within a 400 mile radius

**Correct Answer**: ['San Francisco', 78, ['Golden State Warriors', 'Los Angeles Lakers'] **Student Answer**: ['San Francisco', 74, ['Los Angeles Lakers', 'Golden State Warriors']] **Reasoning**: The Student Answer is correct because it identifies the same city , the temperature is within 10% of the Correct Answer and the same team names are present in the list.

Final Grade: CORRECT

**Question**: Find the name of the city known for its famous tourist attraction Alcatraz, also give it's current temperature and a list of names of all the NBA teams whose home stadium is within a 400 mile radius

**Correct Answer**: ['San Francisco', 78, ['Golden State Warriors', 'Los Angeles Lakers'] **Student Answer**: The city name is San Francisco, its temperature is 80 degrees and the Los Angeles Lakers and the Golden State Warriors are two NBA teams whose home stadium is within a 400 mile radius

**Reasoning**: Although the Student Answer is correct (identifies the same city, the temperature is within 10% of the Correct Answer and the same team names are present), it's not formatted the same and contains extra text and natural language.

Final Grade: CORRECT BUT BAD FORMATTING

**Question**: Find the name of the city known for its famous tourist attraction Alcatraz, also give it's current temperature and a list of names of all the NBA teams whose home stadium is within a 400 mile radius

**Correct Answer**: ['San Francisco', 78, ['Golden State Warriors', 'Los Angeles Lakers'] **Student Answer**: ['San Francisco', -15, ['Los Angeles Lakers', 'Golden State Warriors']] **Reasoning**: The Student Answer is incorrect because although identifies the same city and the same team names are present in the list, the temperature is well outside of 10% of the Correct Answer.

Final Grade: INCORRECT

790

# C.5.3 In-Context Examples (Sorting)

If 'Question:' asks for a special sorting of the list, make sure 'Student Answer:' is sorted the same as 'Correct Answer:'. So if 'Question:' instead asked to list the names of all the NBA teams whose home stadium is within a 400 mile radius in alphabetical order, we would want 'Student Answer:' to contain ['San Francisco', 78, ['Golden State Warriors', 'Los Angeles Lakers']].

Examples:

**Question**: Find the name of the city known for its famous tourist attraction Alcatraz, also give it's current temperature and a list of names of all the NBA teams whose home stadium is within a 400 mile radius in alphabetical order

Correct Answer: ['San Francisco', 78, ['Golden State Warriors', 'Los Angeles Lakers']

**Student Answer**: ['SF', 75, ['Golden State Warriors', 'Los Angeles Lakers']]

**Reasoning**: The Student Answer is correct because it identifies the same city (SF is a commonly known short form for San Francisco), the temperature is within 10% of the Correct Answer and the same team names are present in the list and sorted the same given the Question asks for a special sorting (alphabetical order).

Final Grade: CORRECT

**Question**: Find the name of the city known for its famous tourist attraction Alcatraz, also give it's current temperature and a list of names of all the NBA teams whose home stadium is within a 400 mile radius in alphabetical order

**Correct Answer**: ['San Francisco', 78, ['Golden State Warriors', 'Los Angeles Lakers'] **Student Answer**: The city name is San Francisco, its temperature is 80 degrees and the Golden State Warriors and the Los Angeles Lakers are two NBA teams whose home stadium is within a 400 mile radius (in alphabetical order)

**Reasoning**: Although the Student Answer is correct (identifies the same city, the temperature is within 10% of the Correct Answer and the same team names are present in the same order), it's not formatted the same and contains extra text and natural language.

Final Grade: CORRECT BUT BAD FORMATTING

**Question**: Find the name of the city known for its famous tourist attraction Alcatraz, also give it's current temperature and a list of names of all the NBA teams whose home stadium is within a 400 mile radius in alphabetical order

Correct Answer: ['San Francisco', 78, ['Golden State Warriors', 'Los Angeles Lakers'] **Student Answer**: ['San Francisco', 79, ['Los Angeles Lakers', 'Golden State Warriors']] **Reasoning**: The Student Answer is incorrect because although it identifies the same city and the temperature is within of 10% of the Correct Answer, the list of teams name is in a differently ordered even though 'Question:' explicitly asks for a special sorting.

Final Grade: INCORRECT

# 794 D ToolComp Examples

#### 5 D.1 Example Corrected Trajectory

**Question:** Suppose you have a 2 inch  $\times$  2 inch box. How many of these boxes can you fit into the biggest U-Haul truck (biggest truck as of 2024)?

796

#### Action Plan

- 1. Use the google\_search tool to find the name of the biggest U-Haul truck and find its dimensions.
- 2. Use python\_interpreter to calculate the total number of boxes that would fit given the dimensions.
- 3. Return the answer.

798 799

# ReAct Step 1

- **Thought:** I need to find the dimensions of the biggest U-Haul truck.
- Action: google\_search
- Action Input: {"query": "biggest U-Haul truck 2024"}
- **Observation:** {(search results omitted for brevity)}

800

# 801 ReAct Step 2

- **Thought:** I now have the dimensions of the biggest U-Haul truck, so I can use Python to calculate how many boxes can fit inside.
- Action: python\_interpreter
- Action Input:

• Observation: null}

```
{"result": "Total number of boxes: 356132", "error":
```

802

## ReAct Step 3

- **Thought:** I have the total number of boxes that can fit inside the truck, so I should return that as the answer.
- Action: finish
- Action Input: {"answer": "| 356132 |"}
- Observation: 356132

# D.2 Example Annotated Trajectory

**Question:** Using the 2023 birth-rate, death-rate, and population of Japan, assuming no one leaves or enters, what will be the projected population in 10 years?

806 807

#### Action Plan

#### Model

1. Use the python\_interpreter tool to calculate the population growth rate using the birth-rate, death-rate, and population of Japan.

2. Use the python\_interpreter tool to calculate the population in 10 years using the population growth rate and the current population of Japan.

- 3. Use the google\_search tool to search for the current birth-rate, death-rate, and population of Japan.
- 4. ... (truncated since the model goes in circles)

(Label: X)

808

# (Label. A

# ReAct Step 1

## Model

- Thought: I need to use the google\_search tool to calculate the population growth rate using the birth-rate, death-rate, and population of Japan in 2023. (Label: ✓)
- Action: google\_search (Label: ✓)
- Action Input: { "query": "Japan 2023 population growth rate using birth-rate, death-rate, and population" } (Label: )

• Observation:

[ {"source": "nippon.com",
"date": "Mar 13, 2024",
"title": "Number of Births
in Japan Reaches New Low in
2023", "snippet": "In 2023,
Japan had more than twice as
many deaths as births, leading
to a natural population decline
(excluding migration) of over
800000.", "highlight": "had
more than twice as many deaths
as births"}, ... truncated]

## Human

- 1. Use the google\_search tool to calculate the population growth rate using the birthrate, death-rate, and population of Japan in 2023.
- 2. Use the python\_interpreter tool to calculate the population in 10 years using the population growth rate and the population of Japan in 2023.
- 3. Use the finish action to return the final answer.

#### Human

- Thought: N/A
- Action: N/A
- Action Input: { "query": "Japan 2023 birth-rate per 1000"}
- Observation:

[ {"source":
 "en.wikipedia.org", "date":
 null, "title": "Demographics
 of Japan", "snippet": "6.6
 births/1,000 population
 (2023 est.) Death rate,
 11.7 deaths/1,000 population
 (2023 est.) Life expectancy,
 84.83 years. 2022 male, 81
 years.", "highlight": "6.6
 births/1,000 population"}, ...
 truncated]

#### 811 ReAct Step 2

812

#### Model

- Thought: I need to use the google\_search tool to calculate the population growth rate using the birth-rate, death-rate, and population of Japan in 2023. (Label: X)
- Action: google\_search (Label: ✓)
- Action Input: { "query": "Japan 2023 death-rate per 1000" } (Label:

• Observation:

[ {"source":
 "data.worldbank.org", "date":
 null, "title": "Death rate,
 crude (per 1000 people) Japan", "snippet": "Death rate,
 crude (per 1000 people) - Japan
 from The World Bank: Data.",
 "highlight": "Death rate | per
 1000 | Japan" }, ... truncated]

#### Human

- **Thought**: According to wikipedia the 2023 estimates for the birthrate 6.6 births/1000 population and 11.7 deaths/1000 population. I need to find the 2023 population for Japan.
- Action: N/A
- Action Input: { "query": "Japan 2023 population"}
- Observation:

[{ "source": "statista.com",
 "date": "Jul 04, 2024",
 "title": "Japan: total
 population 2019-2029",
 "snippet": "In 2023, the total
 population of Japan amounted
 to around 124.62 million
 inhabitants. See the figures
 for the population of South
 Korea for ...", "highlight":
 "124.62 million inhabitants" },
 ... truncated]

#### 813 ReAct Step 3

## Model

• **Thought**: I need to use the python interpreter tool to calculate the population in 10 years using the population growth rate and the population of Japan in 2023.

(Label: ✓)

- Action: python\_interpreter (Label: ✓)
- Action Input:

814

import numpy as np
population = 124352000
population\_growth\_rate = 0.0012
population\_in\_10\_years
= population \* (1 +
population\_growth\_rate)\*\*10
print(population\_in\_10\_years)
(Label: X)

• **Observation:** 125852307.84945872

#### Human

Thought: N/AAction: N/AAction Input:

population = 124620000birth\_rate\_per\_thousand = 6.6 death\_rate\_per\_thousand = 11.7 total\_births\_extrapolated = birth\_rate\_per\_thousand \* (population / 1000) total\_deaths\_extrapolated = death\_rate\_per\_thousand \* (population / 1000) population\_growth\_rate = (total\_births\_extrapolated - total\_deaths\_extrapolated) / population population\_in\_10\_years = population \* (1 + population\_growth\_rate) \*\*10 print(population\_in\_10\_years)

• **Observation**: 118408275.35965854

# ReAct Step 4

816

# Model

- Thought: I have enough information to answer the question (Label: ✓)
- Action: finish (Label: ✓)

• Action Input: {"answer": ["118,408,275"]} (Label: ✗)

• **Observation:** ["118, 408, 275"]

# Human

• Thought: N/A

• Action: N/A • Action Input:

{"answer": ["118408275"]}

• **Observation**: ["118408275"]

# 17 E Tools

- In this appendix section, we provide a detail breakdown of each of the 11 tools, providing a description,
- the parameters, an input example and a corresponding output example.
- 820 E.1 Date
- Description: Returns the current date (e.g., January 1, 2024).
- 822 Input Example:

```
823
824 {}
```

# 826 Output Example:

# 833 Parameters:

```
834
§38 []
```

## 837 E.2 Calculator

838 **Description:** Calculates expressions including basic arithmetic and brackets.

# Input Example:

839

852

```
840
841 {
842    "operation": "2*32-4+456+(1+2)+3+(1/2*3+3+(1+2))"
843 }
```

# 845 Output Example:

```
846
847
848 "error": "",
849 "result": "529.5"
}
```

## **Parameters:**

#### E.3 Current Weather

Description: Retrieves current daily averages for temperature, rainfall, and hours of precipitation for 864 a specified city and country. Does not return historical data.

# **Input Example:**

866

```
868
       "city_name": "London",
869
870
       "country_code": "GB"
87<u>1</u>
    }
```

# **Output Example:**

```
873
874
875
      "error": ""
876
      "result": [
877
878
        {
           "date": "2024-03-25 00:00:00",
879
           "temperature (F)": "47.78615",
880
           "total rain (mm)": "1.4000001",
881
           "total snowfall (mm)": "0.0",
882
883
           "precipitation hours (hours)": "4.0"
884
885
           "date": "2024-03-26 00:00:00",
886
           "temperature (F)": "48.374897",
887
           "total rain (mm)": "8.2",
888
           "total snowfall (mm)": "0.0"
889
           "precipitation hours (hours)": "11.0"
890
        }
891
892
           "date": "2024-03-27 00:00:00",
893
           "temperature (F)": "47.217274",
894
           "total rain (mm)": "2.399999",
895
           "total snowfall (mm)": "0.0",
896
           "precipitation hours (hours)": "4.0"
897
898
      ]
899
    }
989
```

# **Parameters:**

```
903
    Γ
904
905
      {
        "Input Name": "city_name",
906
        "Type": "String",
907
        "Description": "The name of the city."
908
      },
909
910
911
        "Input Name": "country_code",
        "Type": "Two Alphabet - Number",
912
        "Description": "The country code (ISO 3166-2). The list can be
913
            found here: https://en.wikipedia.org/wiki/ISO_3166-2"
914
915
    ]
916
```

#### E.4 Historical Weather

Description: Retrieves daily averages for temperature and precipitation starting from the 1940s for a given city. Note: 5-day data delay, meaning you cannot get current weather data for the last 5 days.

## 921 Input Example:

```
922

923

924 "city_name": "London",

925 "country_code": "GB",

926 "start_date": "2023-03-09",

927 "end_date": "2023-03-21"

928 }
```

# Output Example:

930

```
931
932
      "error": "",
933
      "result": [
934
935
           "date": "2024-03-09 00:00:00",
936
           "temperature (F)": "48.102356"
937
           "total rain (mm)": "0.4"
938
           "total snowfall (mm)": "0.0"
939
           "precipitation hours (hours)": "2.0"
940
        },
941
942
943
           "date": "2024-03-23 00:00:00",
944
           "temperature (F)": "43.373596"
945
           "total rain (mm)": "1.0999999"
946
           "total snowfall (mm)": "0.42000002"
947
           "precipitation hours (hours)": "3.0"
948
        }
949
      ]
950
   }
952
```

#### **Parameters:**

```
954
955
    Γ
956
      {
        "Input Name": "city_name",
957
        "Type": "String",
958
        "Description": "The name of the city."
959
      },
960
961
        "Input Name": "country_code",
962
        "Type": "Two Alphabet - Number",
963
964
        "Description": "The country code (ISO 3166-2). The list can be
965
            found here https://en.wikipedia.org/wiki/ISO_3166-2"
      },
966
967
        "Input Name": "start_date",
968
        "Type": "Date Format",
969
        "Description": "The start date in YYYY-MM-DD format"
970
976
972
973
        "Input Name": "end_date",
        "Type": "Date Format",
974
         "Description": "The start date in YYYY-MM-DD format"
975
976
   ]
9<del>7</del>7
```

#### E.5 Wiki Search

980 **Description:** Searches Wikipedia and returns a summary of the top pages matching the query.

# **Input Example:**

981

988

```
982

983

984  "query": "covid-19",

985  "num_results": "1"

986 }
```

## **Output Example:**

```
989
990
       "error": ""
991
       "result": [
992
993
994
           "title": "COVID-19",
           "summary": "Coronavirus disease 2019 (COVID-19) is a contagious
995
              disease caused by the coronavirus SARS-CoV-2. The first
996
              known case was identified in Wuhan, China, in December 2019.
997
                Most scientists believe the SARS-CoV-2 virus entered into
998
999
              human populations through natural zoonosis, similar to the
              SARS-CoV-1 and MERS-CoV outbreaks, and consistent with other
1000
                pandemics in human history. Social and environmental
1001
1002
              factors including climate change, natural ecosystem
              destruction and wildlife trade increased the likelihood of
1003
1004
              such zoonotic spillover. The disease quickly spread
              worldwide, resulting in the \texttt{COVID-19} pandemic. The \texttt{symptoms}
1005
              of COVID-19 are variable but often include fever, fatigue,
1006
              cough, breathing difficulties, loss of smell, and loss of
1007
1008
              taste. Symptoms may begin one to fourteen days after
              exposure to the virus. At least a third of people who are
1009
              infected do not develop noticeable symptoms. Of those who
1010
1011
              develop symptoms noticeable enough to be classified as
              patients, most (81%) develop mild to moderate symptoms (up
1012
              to mild pneumonia), ... truncated"
1013
1014
1015
      ]
    }
1819
```

#### **Parameters:**

```
1019
     [
1020
1021
         "Input Name": "query",
1022
         "Type": "String"
1023
         "Description": "The search query."
1024
1025
1026
         "Input Name": "num_results (Optional)",
1027
         "Type": "Integer",
1028
         "Description": "Number of search results to return."
1029
1030
    ]
1832
```

#### E.6 Intraday Stock Info

1034 **Description:** Provides intraday time series data for specified equities.

# **Input Example:**

1035

1042

```
1036
1037 {
1038     "symbol": "AAPL",
1039     "interval": "60min"
}
```

## **Output Example:**

```
1043
1044
       "error": "",
1045
1046
       "result": [
1047
         {
            "timestamp": "2024-07-16 19:00:00",
1048
            "open_market_value": "234.6520",
1049
            "high_market_value": "234.7200",
1050
            "low_market_value": "234.2200",
105$
1052
            "close_market_value": "234.3200",
            "volume": "38722"
1053
         },
1054
1055
            "timestamp": "2024-07-16 18:00:00",
1056
           "open_market_value": "234.6220",
1057
            "high_market_value": "234.7500",
1058
            "low_market_value": "234.5050",
1059
            "close_market_value": "234.7000",
1060
            "volume": "24098"
106$
         },
1062
1063
         {
1064
           "timestamp": "2024-07-08 16:00:00",
1065
            "open_market_value": "227.8100",
1066
            "high_market_value": "227.8800",
1067
            "low_market_value": "226.0630";
1068
            "close_market_value": "227.6400",
1069
            "volume": "14364524"
1070
1078
       ]
1072
    }
1879
```

### **Parameters:**

```
1076
1077
1078
         "Input Name": "symbol",
1079
         "Type": "String",
1080
1081
         "Description": "The ticker symbol of the equity."
       },
1082
1083
         "Input Name": "interval",
1084
         "Type": "String",
1085
1086
         "Description": "Data point interval (1min, 5min, etc.)."
1087
       },
1088
         "Input Name": "month (optional)",
1089
         "Type": "String",
1090
         "Description": "You can use the month parameter (in YYYY-MM format
1091
1092
             ) to query a specific month in history."
       }
1093
    ]
1893
```

# E.7 Daily Stock Info

Description: Returns daily time series data for specified equities.

## **Input Example:**

1097

1098

```
1099
1100
        "symbol": "AAPL",
1101
1102
        "number_of_days": 5
     }
1183
```

## **Output Example:**

```
1105
1106
1107
       "error": "",
1108
       "result": [
1109
         {
1110
1115
           "timestamp": "2024-07-16",
           "open_market_value": "235.0000",
1112
            "high_market_value": "236.2700",
1113
            "low_market_value": "232.3300",
1114
            "close_market_value": "234.8200",
1115
            "volume": "43234278"
1116
         },
1117
1118
            "timestamp": "2024-07-15",
1119
           "open_market_value": "236.4800",
1120
           "high_market_value": "237.2300",
1125
            "low_market_value": "233.0900",
1122
            "close_market_value": "234.4000",
1123
            "volume": "62631252"
1124
1125
         },
1126
         {
1127
           "timestamp": "2024-07-10",
1128
           "open_market_value": "229.3000",
1129
            "high_market_value": "233.0800",
1130
            "low_market_value": "229.2500",
1135
            "close_market_value": "232.9800",
1132
            "volume": "62627687"
1133
1134
       ]
1135
    }
1139
```

# **Parameters:**

```
1139
     Ε
1140
       {
1141
         "Input Name": "symbol",
1142
         "Type": "String",
1143
         "Description": "The ticker symbol of the equity."
1144
       },
1145
1146
         "Input Name": "number_of_days",
1147
         "Type": "Integer",
1148
         "Description": "The number of days before today to return data for
1149
1150
       }
1151
    ]
1153
```

# 1154 E.8 Stock Symbol Search

1155 **Description:** Searches for stock tickers based on provided keywords.

# Input Example:

1156

1162

```
1157
1158 {
"keywords": "tesla"
1169 }
```

#### **Output Example:**

```
1163
1164
       "error": "",
1165
       "result": [
1166
1167
            "symbol": "TSLA",
1168
            "name": "Tesla Inc",
1169
            "type": "Equity",
1170
            "region": "United States",
117$
            "market_open": "09:30",
1172
            "market_close": "16:00"
1179
            "timezone": "UTC-04",
1174
            "currency": "USD",
1175
            "match_score": "0.8889"
1176
1177
1178
            "symbol": "TLO.DEX",
1179
            "name": "Tesla Inc",
1180
            "type": "Equity", "region": "XETRA"
1181
1182
            "market_open": "08:00"
1183
            "market_close": "20:00",
1184
            "timezone": "UTC+02",
1185
            "currency": "EUR",
1186
            "match_score": "0.7143"
1187
         },
1188
1189
          {
1190
            "symbol": "TLO1.FRK",
1191
            "name": "TESLA INC. CDR DL-001",
1192
            "type": "Equity",
1198
            "region": "Frankfurt",
1194
            "market_open": "08:00"
1195
            "market_close": "20:00",
1196
            "timezone": "UTC+02",
1197
            "currency": "EUR",
1198
1199
            "match_score": "0.3846"
1200
       ]
1201
    }
1203
```

#### **Parameters:**

```
1205
     Ε
1206
       {
1207
         "Input Name": "keywords",
1208
         "Type": "String",
1209
         "Description": "Keywords to search, , e.g., company name, to
1210
             retrieve the ticker symbol for"
1211
1212
    ]
1213
```

# E.9 Python

1215

1223

1244

1250

1257

1216 **Description:** Runs a python interpreter on a code snippet.

# Input Example:

```
1218
1219
1220 "code": "print(4 + 5)"
1221 }
```

## **Output Example:**

```
1224
1225
1226 "result": "9",
1227 "error": ""
1228 }
```

#### 1230 Parameters:

# 1241 E.10 Wolfram Alpha

Description: Accesses Wolfram Alpha to generate outputs from the Knowledgebase for computations and data queries. Wolfram Alpha excels at complex number-crunching, computation and calculations.

# **Input Example:**

```
1245
1246 {
1247 "query": "what is Ronaldo's age?"
1248 }
```

#### **Output Example:**

```
1251
1252 {
1253    "error": "",
1254    "result": "47 years 5 months 13 days"
1255 }
```

# **Parameters:**

```
Ε
1259
1260
         "Input Name": "query",
1261
         "Type": "String"
1262
         "Description": "The query to perform computations/searches on.
1263
             When unsure of your query search, try searching yourself on
1264
1265
             the website!"
       }
1266
    ]
1268
```

# E.11 Google Search

Description: Performs a Google search and returns snippet results, without linked page details
Google is often used for popular culture, location-awareness and crowdsourcing.

#### Input Example:

1269

1272

1279

```
1273

1274

1275

1276

1276

1277

| "query": "What is the capital of France?",

"location": "Paris"

| 1277

}
```

# **Output Example:**

```
1280
    {
1281
       "error": "",
1282
       "result": [
1283
1284
           "source": "en.wikipedia.org",
1285
           "date": "None",
1286
           "title": "Paris",
1287
           "snippet": "Paris is the capital and largest city of France.
1288
               With an official estimated population of 2,102,650 residents
1289
1290
                as of 1 January 2023 in an area of more than ...",
           "highlight": "Paris"
1299
         },
1292
1293
           "source": "home.adelphi.edu",
1294
           "date": "None",
1295
           "title": "Paris facts: the capital of France in history",
1296
           "snippet": "Paris facts: Paris, the capital of France. Paris is
1297
               the capital of France, the largest country of Europe with 55
1298
               0 000 km2 (65 millions inhabitants).",
1299
           "highlight": "Paris"
1300
         },
1301
1302
         {
1303
           "source": "britannica.com",
1304
           "date": "None",
1305
           "title": "France | History, Maps, Flag, Population, Cities,
1306
               Capital, & ...",
1307
           "snippet": "Get a special academic rate on Britannica Premium.
1308
               The capital and by far the most important city of France is
1309
               Paris, one of the world's preeminent cultural ...",
1310
           "highlight": "Paris"
1314
1312
         },
      ]
1313
    }
1314
```

#### **Parameters:**

```
1317
1318
1319
       {
         "Input Name": "query",
1320
         "Type": "String"
1324
         "Description": "The search query."
1322
       },
1328
1324
         "Input Name": "location (Optional)",
1325
         "Type": "String",
1326
          "Description": "The geographical location for the search (optional
1327
1328
       }
1329
     ]
1339
```