TOOLCOMP: A MULTI-TOOL REASONING & PROCESS SUPERVISION BENCHMARK

Anonymous authors

Paper under double-blind review

ABSTRACT

Despite recent advances in AI, the development of systems capable of executing complex, multi-step reasoning tasks involving multiple tools remains a significant challenge. Current benchmarks fall short in capturing the real-world complexity of tool-use reasoning, where verifying the correctness of not only the final answer but also the intermediate steps is important for evaluation, development, and identifying failures during inference time. To bridge this gap, we introduce ToolComp, a comprehensive benchmark designed to evaluate multi-step tool-use reasoning. ToolComp is developed through a collaboration between models and human annotators, featuring human-edited/verified prompts, final answers, and process supervision labels, allowing for the evaluation of both final outcomes and intermediate reasoning. Evaluation across six different model families demonstrates the challenging nature of our dataset, with the majority of models achieving less than 50% accuracy. Additionally, we generate synthetic training data to compare the performance of outcome-supervised reward models (ORMs) with process-supervised reward models (PRMs) to assess their ability to improve complex tool-use reasoning as evaluated by ToolComp. Our results show that PRMs generalize significantly better than ORMs, achieving a 19% and 11% improvement in rank@1 accuracy for ranking base and fine-tuned model trajectories, respectively. These findings highlight the critical role of process supervision in both the evaluation and training of AI models, paving the way for more robust and capable systems in complex, multi-step tool-use tasks. *

031 032

033

000

001

002 003 004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

1 INTRODUCTION

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

041 However, evaluating the effectiveness of these tool-use capabilities remains a pressing challenge. 042 While there have been notable efforts in developing benchmarks for tool-use capability, these often 043 assess isolated instances of tool use, focusing on whether the model can invoke the correct tool 044 at the right time (Huang et al., 2024; Zhuang et al., 2023; Peng et al., 2021). Additionally, while 045 benchmarks for multi-step tool usage exist, most focus only on scoring the correctness of the final 046 answer (Mialon et al., 2023), despite that the complex nature of multi-step reasoning often requires 047 the evaluation for partial correctness or step-wise correctness of the reasoning trajectories. This can 048 be valuable for both understanding model failure modes and developing systems that can improve upon these intermediate reasoning flaws. 049

To address these shortcomings, we introduce ToolComp, a benchmark comprising 485 complex, 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

⁰⁵³

^{*}Code, ToolComp benchmark, and synthetic training data will be made publicly available.

providing human verification, ToolComp offers a rigorous assessment of a model's ability to per form 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.

Moreover, in light of recent works demonstrating how process supervision significantly improve reasoning in language models (Lightman et al., 2023; Wang et al., 2024a), we explore the best methods for improving agentic tool-use reasoning by conducting an initial comparative analysis between process-supervised reward models (PRMs) and outcome-supervised reward models (ORMs) on ToolComp. Our results demonstrate that process-supervised models outperform outcome-based approaches, underscoring the importance of training models with and evaluating against process supervision signals.

In order to avoid early contamination, we plan to open source the entire benchmark when either 1) any model scores over 85% on ToolComp or 2) at the end of 2025, whichever comes earlier.

- 068 1.1 CONTRIBUTIONS AND KEY TAKEAWAYS
- 070 Our key contributions and takeaways are summarized as follows:
 - **Introduction of ToolComp** We introduce ToolComp, a multi-tool reasoning and process supervision benchmark with 485 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 1731 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 16 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 o1-preview has the best performance, achieving 61.83% against the human-verified final answers and 80.18% against the process supervision labels (Section 4 and Section A).
 - Process-Supervision Outperforms Outcome-Supervision Our analysis shows that process-supervised reward models outperform outcome-based reward models by 19% in rank@1 accuracy on base model generations and by 11% in rank@1 accuracy on fine-tuned model generations (Section 5 and Section B).
 - Release of Synthetic Training Data and Evaluation Pipeline We provide the complete synthetic training dataset, which includes 11K prompts, 11K preferred and dis-preferred full trajectories for ORM training, and 15K preferred and dis-preferred steps for PRM training. We also release the evaluation pipeline for ToolComp. (Section 5).
 - 2 RELATED WORKS

071

073

075

076

077 078

079

081

082

084

085

090

092

Benchmarks for Complex Tool Use Planning With rising interest in tool-augmented LLMs 094 (Schick et al., 2023; Patil et al., 2023; Qin et al., 2023), several benchmarks have been introduced 095 to assess their abilities. Earlier benchmarks were designed to assess a model's ability to do proper 096 retrieval, execution, and extraction of one tool call for specific tasks such as general question an-097 swering (Yang et al., 2018; Joshi et al., 2017), fact verification (Thorne et al., 2018), or answering 098 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 100 multiple tool calls to answer more complex queries. More recent benchmarks aimed at evaluating 101 multiple tool calls are often placed within or dependent on state-full systems (such as a code-base 102 and/or a dynamic database) (Yan et al., 2024; Jimenez et al., 2024; Liu et al., 2023). Although 103 these types of benchmarks assess a language model's ability to chain together multiple tool calls, 104 the evaluation may penalize general-purpose language models that are not familiar with the given 105 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 106 against another reference state-of-the-art model (Qin et al., 2023), both of which lack the rigour of 107 human-verified ground truth final answers. Closest to our work, the GAIA benchmark is a collection

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	 Image: A second s	1	1	1	1	 Image: A second s
Multi-Tools Scenario	1	1	 Image: A second s	1	×	×
Multi-Step Reasoning	1	 Image: A second s	 Image: A second s	1	1	1
Step-Wise Labels	✓	×	×	×	×	×
Verified Final Answer	✓	 Image: A second s	×	×	×	√ *
Number of Tools	11	23	3	3451	53	8

120 121 122

1 1

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 di verse 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).

135 In this work, through a hybrid human-AI annotation workflow, we generate per-step process su-136 pervision labels, which uniquely enable us to rigorously evaluate a model's intermediate reasoning 137 capability. Table 1 provides a comparative overview of popular tool-use benchmarks, including 138 our work, ToolComp. In addition, we investigate how to best apply process supervision signals to 139 improve multi-step tool-use reasoning, which introduces novel design challenges compared to its application in mathematical reasoning. For instance, the granularity of supervision becomes a key 140 consideration, where we must decide between supervising the entire ReAct (Yao et al., 2023) pro-141 cess or its subcomponents. These design choices alongside comparisons with outcome supervision 142 are explored in detail in Section 5. 143

144 145

3 TOOLCOMP

146 147

3.1 Tools

148 For the creation of this benchmark and evaluation framework, we support 11 tools: Date, Current 149 Weather, Historical Weather (Zippenfenig, 2024), Calculator, Wiki Search (Majlis, 2017), Google 150 Search (SerpApi, 2024), Wolfram Alpha (Wolfram Research, 2024), Intra-day Stock Info, Daily 151 Stock Info, Stock Symbol Search (AlphaVantage), and Python. There were several considerations 152 when choosing these set of tools, namely, we wanted to cover a broad range of use cases from fact 153 retrieval to financial assistant, have some overlap in use cases to encourage various valid trajectories, 154 ensure the tools are general enough to not require specialized knowledge for LLMs to use, and allow 155 for interesting interactions between tools. A detailed breakdown of each tool, including descriptions, 156 parameters, input examples, and output examples are available in Appendix G.

157 158

159

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



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 (or by a state-of-the-art LM for generating synthetic training data as further described in Section 5.1). A full annotated trajectory example is available in Appendix F.2.

176 177 178

179

180 181

170

171

172

173

174

175

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.

182 3.3 PROMPT CREATION

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 to the prompt can be programmatically verified. To achieve this, we generate a set of candidate 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 C.1

191 192

193

3.4 CHAT VS. ENTERPRISE USE CASES

In creating the benchmark, we developed two subsets of prompts, coined ToolComp-Enterprise and 194 ToolComp-Chat. ToolComp-Enterprise allows the use of 11 tools and aims to emulate settings 195 in which LLM agents must compose a larger number of expressive APIs together correctly, such 196 as in enterprise settings. The second subset, ToolComp-Chat, is designed to test general purpose 197 chatbots with the minimally sufficient set of tools for information retrieval and processing tasks, namely Google Search and Python. We chose only google search and python execution as these are 199 standard tools found in major chat-bot providers. We only allow the respective tools for each subset 200 during prompt generation, labeling, and evaluation. ToolComp-Enterprise contains 287 examples 201 and ToolComp-Chat contains 198 examples.

- 202
- 203 3.5 LABEL CREATION

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 E.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 E.2.

We asked the annotators to rate if a step is Correct (i.e., the step is a reasonable action towards achieving the final answer) or Incorrect (i.e., the step is nonsensical, incorrect, or is not a reasonable action towards acheiving the final answer). All components of the ReAct Step (Thought, Action, Action Input) must be marked as Correct or Incorrect by the annotator. If the annotator marks a step as Correct, the model is allowed to proceed further and generate the next step. If the annotator 216 deems a step as Incorrect, they must modify the step to make it correct. Once corrected, the model 217 is then prompted to advance to the next step with the human-corrected step as part of its context. 218 This is repeated until the Finish Action is chosen by the LLM and marked as Correct by the anno-219 tator or until the annotator corrects an Action step to 'Finish' because we have enough information 220 to answer the question. The overall flow is shown in Figure 1. An example golden trajectory is available in Appendix F.1 and an example annotated trajectory is available in Appendix F.2. We 221 use FireFunction-V1 as the Policy Model LLM (at the time, this was the best open-source tool-use 222 LLM) and humans as the annotators (Fireworks, 2024). 223

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).

228 229

230

3.6 QUALITY CONTROL

To ensure the highest quality of ToolComp, we conduct a thorough manual inspection of all examples. 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 trusted annotators to perform a final re-review of all samples and make any necessary corrections.

As a final quality control step, we evaluate the entire benchmark using GPT-40 (May 2024), GPT-4
Turbo, Claude 3.5 Sonnet, and Llama 3.1 405b (OpenAI et al., 2024; Dubey et al., 2024; Anthropic).
We identify the set of data samples where all models' answers differed from the ground truth final answers. We then repeated the refinement process on these samples, as they represented the most challenging and/or potentially mislabeled data points. This iterative approach yielded the final version of ToolComp.

241 242

4 TOOLCOMP EVALUATIONS

243 244 245

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.

251 LLM Grading By using LLM grading against ground truth answers we opt to be charitable to 252 exact formatting and focus on assessing the tool use capabilities of the model. We intentionally 253 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, 254 containing multiple elements, lists which may or may not be sorted, and dictionaries. This approach 255 prompts an LLM Judge to look at the prompt, the ground truth answer, and the model's answer and 256 asks the model to classify it as Incorrect, Correct, or Correct with Bad Formatting. We use GPT-4 257 Turbo as the de-facto judge for all of our models (OpenAI et al., 2024). The prompt used is shown 258 in Appendix E.5. We consider both Correct and Correct with Bad Formatting as a win (accurate) 259 and Incorrect as a loss (inaccurate).

260 261 262

4.2 FINAL ANSWER EVALUATIONS

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% confidenceintervals (CIs) for the entire benchmark. We use native function calling for all the models except for o1-preview. Since the o1-preview API does not accept system prompts nor allows for native function calling, we prepend our ReAct System Instruction (Appendix E.2) to the user query. Additionally, 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 scores are not indicative of tool or endpoint failures due to rate limiting, we use verbose logging to

300

301

302

303

270 Table 2: Accuracy and the 95% CIs of all selected models using the final answer and scored using 271 an LLM judge (Dubey et al., 2024; OpenAI et al., 2024; Gemini et al., 2024; Anthropic; Mistral; 272 Cohere). We combined the results of each subset to give an overall score for the entire benchmark. 273 Exact Match results are reported in Appendix A.1 but the rankings do not significantly differ, with 274 the top 5 and bottom 4 models remaining the same. *Llama models sometimes fail to call tools/terminate early or call tools in the wrong format. Using constrained decoding and other techniques to 275 guarantee structured outputs can improve their performance. **Since o1-preview does not support 276 native function calling via API, we prompt the model to formulate tool calls in the ReAct format. 277

Model Family	Model Name	Total (%)	Chat (%)	Enterprise (%)
	o1-preview**	61.83 ± 4.34	55.1 ± 6.96	66.43 ± 5.47
	GPT-40 (Aug 2024)	58.68 ± 4.39	56.85 ± 6.92	59.93 ± 5.67
Open A I	GPT-40 (May 2024)	58.44 ± 4.38	49.5 ± 6.96	64.58 ± 5.52
OpenAi	GPT-4 Turbo Preview	57.61 ± 4.39	53.03 ± 6.95	60.76 ± 5.64
	GPT-4	45.89 ± 4.43	37.88 ± 6.78	51.39 ± 5.77
	GPT-40 Mini	44.03 ± 4.41	32.83 ± 6.54	51.74 ± 5.77
	Claude 3.5 Sonnet	58.03 ± 4.39	56.06 ± 6.91	59.38 ± 5.67
Anthropic	Claude 3 Opus	51.03 ± 4.44	48.49 ± 6.96	52.78 ± 5.77
_	Claude 3 Sonnet	48.56 ± 4.44	40.4 ± 6.84	54.17 ± 5.78
Casala	Gemini 1.5 Pro (Aug 2024)	56.61 ± 4.41	51.27 ± 6.98	60.28 ± 5.66
Google	Gemini 1.5 Pro (May 2024)	38.43 ± 4.34	35.5 ± 6.57	40.42 ± 5.68
Mistral	Mistral Large 2	46.30 ± 4.43	40.4 ± 6.84	50.35 ± 5.78
	Llama 3.1 405B Instruct*	46.19 ± 4.44	40.1 ± 6.84	50.35 ± 5.78
Meta	Llama 3.1 70B Instruct*	35.74 ± 4.27	33.5 ± 6.59	37.23 ± 5.6
	Llama 3.1 8B Instruct*	12.81 ± 2.98	6.09 ± 3.34	17.42 ± 4.39
Cohere	Command R+	26.13 ± 3.91	20.2 ± 5.59	30.21 ± 5.3
	Average	$\textbf{46.64} \pm \textbf{4.27}$	$\textbf{41.08} \pm \textbf{6.50}$	$\textbf{50.46} \pm \textbf{5.58}$

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).

Model Family	Model Name	Total (%)	ReAct (%)	Action Plan (%)
	o1-preview	80.19 ± 1.89	79.62 ± 2.22	81.76 ± 3.55
	GPT-40 (Aug 2024)	72.61 ± 2.11	72.84 ± 2.46	71.98 ± 4.13
OpenAI	GPT-40 (May 2024)	71.24 ± 2.14	71.37 ± 2.49	70.88 ± 4.17
OpenAi	GPT-4 Turbo Preview	70.66 ± 2.15	70.18 ± 2.52	71.98 ± 4.13
	GPT-40 Mini	63.02 ± 2.28	64.27 ± 2.64	59.56 ± 4.51
	GPT-4	60.02 ± 2.32	55.87 ± 2.74	71.54 ± 4.15
	Claude 3.5 Sonnet	66.46 ± 2.23	67.74 ± 2.58	62.97 ± 4.44
Anthropic	Claude 3 Opus	64.28 ± 2.27	64.55 ± 2.64	63.52 ± 4.42
-	Claude 3 Sonnet	61.10 ± 2.31	62.93 ± 2.67	56.04 ± 4.56
Caral	Gemini 1.5 Pro (Aug 2024)	69.11 ± 2.19	68.48 ± 2.56	70.88 ± 4.17
Google	Gemini 1.5 Pro (May 2024)	67.89 ± 2.21	67.72 ± 2.58	68.35 ± 4.27
Mistral	Mistral Large 2	72.67 ± 2.11	73.16 ± 2.45	71.32 ± 4.16
	Llama 3.1 405B Instruct	71.62 ± 2.13	73.87 ± 2.42	65.39 ± 4.37
Meta	Llama 3.1 70B Instruct	70.75 ± 2.15	71.33 ± 2.50	69.12 ± 4.25
	Llama 3.1 8B Instruct	57.63 ± 2.34	59.60 ± 2.71	52.20 ± 4.56
Cohere	Command R+	61.31 ± 2.30	64.91 ± 2.63	51.32 ± 4.59
	Average	$\textbf{67.54} \pm \textbf{2.20}$	$\textbf{68.03} \pm \textbf{2.55}$	$\textbf{66.18} \pm \textbf{4.28}$

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 6 of Appendix A.1 to ensure that our LLM Judge (GPT-4 Turbo) 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.

333 334

335

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 humancorrected 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.

343 344 345

356

357

359

360

361

362

364

366

367

368

369 370

372

373

374 375

4.4 INTERMEDIATE REASONING VS. FINAL ANSWER

346 Figure 2 shows the correlation between a model's intermediate reasoning performance and final an-347 swer accuracy based on the multi-step tool-use tasks in ToolComp. The standard Pearson correlation 348 coefficient is r = 0.63 with a statistical p-value of 0.0084, which makes the correlation statistically 349 significant under a significance level of 0.05 (Freedman et al., 2007). Intuitively, this suggests that 350 with stronger step-wise performance as assessed by our LLM-as-judge evaluations, we can expect 351 an increased likelihood of reaching the correct final answer. However, the moderate magnitude of 352 the correlation value could be due to additional signals captured by the step-wise reasoning evalua-353 tions 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 354 refinement, further highlighting the importance of assessing intermediate reasoning. 355





Scale of Training Data (%)

379 380 **Base Model Generations** SFT Model Generations 381 100 100 Process-Supervised RM Process-Supervised RM 382 Outcome-Supervised RM Outcome-Supervised RM 80 80 Rank@1 Accuracy (%) (%) 384 Average Pass@30 Rank@1 Accuracy 385 60 60 Average Pass@30 386 387 40 40 388 389 20 20 390 Pass@

PROCESS SUPERVISION VS. OUTCOME SUPERVISION

Figure 3: A comparison of outcome-supervised and process-supervised reward models across various scales of training data (10%, 25%, 50%, 100%), evaluated by their ability to pick out the best answer out of 30 tool call trajectories. The 95% confidence intervals captures the variance of 500 random samples of 30 completions out of 50 completions. We plot both the performance on generations from Llama-3.1-8b-Instruct (left) and Llama-3.1-8b-Instruct fine-tuned on all the preferred trajectories (right) (Dubey et al., 2024). The plot also shows the Pass@1 given by greedy sampling and the average Pass@30 accuracies for the respective generating models.

100

0

10

Scale of Training Data (%)

100

401

391

392

378

5

С

10

Recent works have demonstrated the power of process supervision signals in domains such as mathematical reasoning (Lightman et al., 2023; Wang et al., 2024a). Despite this, the application of
process-supervised signals towards tool-augmented LLMs remains under-explored. By focusing on
the process rather than just the outcome, process-supervised models could provide more granular
feedback during multi-step tool-use, leading to faster convergence, especially in complex applications where each step is associated with dynamic feedback from the environment. In this section,
we provide a preliminary analysis to assess the value of process supervision for multi-step tool-use
by training PRM and ORM models using additionally generated synthetic data.

409 410

411

5.1 EXPERIMENT DESIGN

Training Dataset For the generation of synthetic prompts and the process/outcome supervision labels, we mirror the strategies outlined in Section 3.3 and Section 3.5, respectively. We use Llama-3.1-8b-Instruct as the Policy Model generator and Chat GPT-40 as the Critic Model generator (OpenAI et al., 2024; Dubey et al., 2024). A detailed accounting of the training data generation, including generation parameters, dataset sizes, and methodology can be found in Appendix D.1. Furthermore, the construction of the outcome-supervised, process-supervised reward modelling and supervised fine-tuning datasets are detailed in Appendix D.2.

419

Reward Model Training Objective We equip the base model, Llama-3.1-8b-Instruct (Dubey 420 et al., 2024), with a linear layer to serve as the reward head. For ORM training, the reward model 421 places a probability of correctness of the whole preferred and dis-preferred trajectory. For PRM 422 training, we experiment with 4 different levels of process supervision. The first axis of variation 423 experiments with including/excluding the "Observation" of the ReAct step, which is the output ob-424 served from the tool call. The second axis of variation ablates the granularity of the process signals, 425 i.e., rewarding the whole ReAct step or rewarding each substep of the ReAct step (recall that a ReAct 426 step contains Thought, Action, and Action Input, the correctness of which is determined/modified 427 by the Critic Model). This leads us to 4 total supervision methods: Full ReAct step with or without 428 Observation, Sub ReAct Step with or without Observation. Under each type of supervision, the 429 PRM places a probability of correctness of the preferred step and dis-preferred step. The training objective is given by the average Binary Cross Entropy loss, assessing the probability of correctness 430 with the corresponding label. A more detailed explanation of the training objective is in Appendix 431 D.3 and training implementation along with hyper-parameters is in Appendix D.5.

432 **Evaluation** To evaluate, we use our ORM and PRM models to select the best response out of a 433 set of candidates. Specifically, we use a generator model to produce 50 completions per problem in 434 ToolComp. We experiment with two different generator model in these experiments: base Llama-435 3.1-8b-Instruct and a supervised fine-tuned Llama-3.1-8b-Instruct model, which is trained on all of 436 the full preferred trajectories in the synthetic training data. After collecting these completions, we use the trained ORM and PRM to rank these completions, returning the best answer according to 437 this ranking. We then assess the rank@1 accuracy by judging the correctness of said highest scoring 438 completion against the ground truth final answer per problem. For the ORM, we simply use the 439 reward score it places on each completion. For the PRM, in order to compress per-step scores into 440 one single score for the entire chain, we experiment with different aggregation functions, namely: 441 min, max, average, and product. 442

To account for variance in the rankings, we perform 500 permutations of 30 random samples from the 50 total completion per problem, and calculate the average rank@1 accuracy. Moreover, we vary the dataset sizes at 10, 25, 50, and 100 percent of the full dataset in order to assess how performance of each method scales with increasing training data.

5.2 Results

447

448

458

449 **PRM outperforms ORM in selecting the best trajectory** In Figure 3, we observe that for both 450 the base model generations and the SFT model generations, our best PRM model (see Table 4) 451 outperforms the ORM in rank@1 accuracy. On the base model generations, the PRM achieves an 452 accuracy of 42.65% compared to the ORM accuracy of 23.89%. Moreover, towards enhancing an 453 already fine-tuned model, we see that the PRM is able to push rank@1 accuracy to 60.25%, nearly 454 matching the generator's pass@30 performance. Overall, these results suggest that PRMs are able 455 to efficiently translate these per tool call step-level signals to provide superior performance than utilizing just outcome-level signals. In Appendix B, we also find the PRM scales better than the 456 ORM on increasing prompt complexity. 457

Full step with observation is the best PRM supervision method Table 4 shows the performance 459 of the PRM using different supervision methods trained on the full-scale dataset. We see that pro-460 viding the model signals about whether an entire step, including the observation, was correct or 461 incorrect led to the best performance. These findings suggest that 1) intermediate information from 462 the environment (tool call outputs) provide valuable signal to PRMs as these dynamic signals pro-463 vide additional insight in determining the correctness of a trajectory, and 2) there is a balance to 464 be struck when determining the granularity of supervision during training. Intuitively, providing 465 less granular, full-step signals gives the model more freedom to learn and identify what makes one 466 step better than another, without being constrained by potentially noisy or overly detailed sub-step labels. This allows the model to generalize more effectively, rather than being limited by specific, 467 finer-grained supervision. 468

Method	Rank@1 Accuracy (%)
Full Step with Observation	60.25
Full Step without Observation	55.43
Sub Step with Observation	49.48
Sub Step without Observation	45.70

Table 4: Comparison of the average rank@1 accuracy across different methods of PRM supervised training, which are all trained on the full-scale training dataset.

477 478

PRM scales better than ORM with increasing data Part of the design of our experiment is to measure the scaling performance of PRMs versus ORMs as we incorporate more training data. In Figure 3, we find that the PRM rank@1 performance consistently outperforms the ORM at all training data scales in ranking both the base model completions as well as the SFT model completions. Interestingly, we observe greater performance scaling for the base model completions, emphasizing the importance of utilizing process level signals as the PRM is able to still pick out the best trajectory amongst lots of low-quality trajectories. This ability could also serve to pick out high-quality trajectories for further training the base model for multi-step tool-use reasoning.

487 Table 5: Comparison of the performance of dif-488 ferent aggregation methods used to combine step-489 wise level PRM scores. Results here use the PRM 490 model trained on all data with the Full Step with 491 Observation supervision method. 492

Method	rank@1 Accuracy (%)
Max	60.25
Min	23.68
Average	25.74
Product	23.06



Figure 4: Distribution of the position of the maximum scoring step, normalized by the length of the trajectory, for the rank@1 selected trajectories.

Max is the best PRM aggregation function Since the PRM provides a score for each step in the trajectory, an important design decision is how to combine step-level scores into a single score 504 that can be used for ranking. Table 5 clearly demonstrates that max is by far the best aggregation function for scoring trajectories. By examining many trajectories and their step-level scores, we 505 find that min and average both heavily penalize any wrong steps that are taken, even if the model 506 eventually recovers and gets to the correct final answer. When using product as the aggregation function, the final aggregation results in low magnitudes that are biased towards shorter trajectories. Max is a better aggregation function because it avoids these aforementioned pitfalls and tends to 509 favor the later steps (as shown by Figure 4) which are a better proxy for a successful trajectory. 510

511 512

507

508

486

493 494

495

496

497

498

LIMITATIONS AND FUTURE DIRECTIONS 6

513 514 515

In this study, we focus solely on applying outcome and process supervision to the reward model. 516 Although fine-tuning the policy model with supervision from the reward model using reinforcement 517 learning (RL) is a logical next step, we leave this for future work and focus instead on the contributed 518 dataset and the value of process supervision even without RL.

519 A notable limitation of our work is the reliance on synthetic data to scale the policy. We hypothe-520 size that incorporating human-generated data to expand the training set could enhance the tool-use 521 capabilities beyond the performance of the state-of-the-art critic model, which was used to label our 522 synthetic dataset. 523

Additionally, the restricted set of tools used in this work, primarily focused on information retrieval and data processing, presents another limitation. In contrast, a common approach in the field involves employing specialized models for various tasks such as image generation and translation. This opens up further questions regarding how process supervision could facilitate the scaling of more nuanced capabilities when integrating with other specialized models.

527 528 529

524

525

526

530

ETHICS STATEMENT 7

531 532

533 In the creation of this dataset and any further analysis, we abide by the ICLR Code of Ethics. We 534 ensure all prompts in this dataset do not contain any harmful or sensitive material by requiring 535 annotators to flag any such prompts. The authors of this paper have also manually inspected all the 536 prompts and tool calls for harmful content. In addition, we applied best practices for code execution, 537 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 538 purposes, and we plan to open-source both the code for running evaluations and the full benchmark dataset.

540 8 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 D.4
and Appendix A.1. Moreover, for the training of the ORM, PRM and SFT models, we detail the
exact process (including methodology, hyperparmeters, and additional settings) in Section 5.1 and
Appendix D. We plan to open source both the code for evaluation and the benchmark dataset.

References

547 548

549

552

553

554

556

- 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.
- 558 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.
- 564 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha 565 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony 566 Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, 567 Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, 568 Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, 569 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny 570 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, 571 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael 572 Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Ander-573 son, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah 574 Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan 575 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-576 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy 577 Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, 578 Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Al-579 wala, 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 580 Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Man-582 nat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, 583 Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, 584 Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur 585 Celebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhar-586 gava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, 588 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 592 Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta,

594 Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petro-595 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, 596 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 597 598 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 600 Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew 601 Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita 602 Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh 603 Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De 604 Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Bran-605 don Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina 606 Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, 607 Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, 608 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 Ar-610 caute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco 611 Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella 612 Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory 613 Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, 614 Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Gold-615 man, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, 616 James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer 617 Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe 618 Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie 619 Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun 620 Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, 621 Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian 622 Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, 623 Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Ke-624 neally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel 625 Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mo-626 hammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navy-627 ata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, 628 Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, 629 Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, 630 Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, 631 Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, 632 Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Sa-633 tadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lind-634 say, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang 635 Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen 636 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, 637 Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, 638 Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Tim-639 othy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, 640 Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu 641 Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, 642 Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef 644 Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models, 2024. 645 URL https://arxiv.org/abs/2407.21783. 646

649

650

651

652

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.

653 Gemini, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, 654 Damien Vincent, Zhufeng Pan, Shibo Wang, Soroosh Mariooryad, Yifan Ding, Xinyang Geng, Fred Alcober, Roy Frostig, Mark Omernick, Lexi Walker, Cosmin Paduraru, Christina Sorokin, 655 Andrea Tacchetti, Colin Gaffney, Samira Daruki, Olcan Sercinoglu, Zach Gleicher, Juliette Love, 656 Paul Voigtlaender, Rohan Jain, Gabriela Surita, Kareem Mohamed, Rory Blevins, Junwhan Ahn, 657 Tao Zhu, Kornraphop Kawintiranon, Orhan Firat, Yiming Gu, Yujing Zhang, Matthew Rahtz, 658 Manaal Faruqui, Natalie Clay, Justin Gilmer, JD Co-Reyes, Ivo Penchev, Rui Zhu, Nobuyuki 659 Morioka, Kevin Hui, Krishna Haridasan, Victor Campos, Mahdis Mahdieh, Mandy Guo, Samer Hassan, Kevin Kilgour, Arpi Vezer, Heng-Tze Cheng, Raoul de Liedekerke, Siddharth Goyal, 661 Paul Barham, DJ Strouse, Seb Noury, Jonas Adler, Mukund Sundararajan, Sharad Vikram, Dmitry 662 Lepikhin, Michela Paganini, Xavier Garcia, Fan Yang, Dasha Valter, Maja Trebacz, Kiran Vo-663 drahalli, Chulayuth Asawaroengchai, Roman Ring, Norbert Kalb, Livio Baldini Soares, Siddhartha Brahma, David Steiner, Tianhe Yu, Fabian Mentzer, Antoine He, Lucas Gonzalez, Bibo 665 Xu, Raphael Lopez Kaufman, Laurent El Shafey, Junhyuk Oh, Tom Hennigan, George van den Driessche, Seth Odoom, Mario Lucic, Becca Roelofs, Sid Lall, Amit Marathe, Betty Chan, Santiago Ontanon, Luheng He, Denis Teplyashin, Jonathan Lai, Phil Crone, Bogdan Damoc, Lewis 667 Ho, Sebastian Riedel, Karel Lenc, Chih-Kuan Yeh, Aakanksha Chowdhery, Yang Xu, Mehran 668 Kazemi, Ehsan Amid, Anastasia Petrushkina, Kevin Swersky, Ali Khodaei, Gowoon Chen, Chris 669 Larkin, Mario Pinto, Geng Yan, Adria Puigdomenech Badia, Piyush Patil, Steven Hansen, Dave 670 Orr, Sebastien M. R. Arnold, Jordan Grimstad, Andrew Dai, Sholto Douglas, Rishika Sinha, Vikas 671 Yadav, Xi Chen, Elena Gribovskaya, Jacob Austin, Jeffrey Zhao, Kaushal Patel, Paul Komarek, 672 Sophia Austin, Sebastian Borgeaud, Linda Friso, Abhimanyu Goyal, Ben Caine, Kris Cao, Da-673 Woon Chung, Matthew Lamm, Gabe Barth-Maron, Thais Kagohara, Kate Olszewska, Mia Chen, 674 Kaushik Shivakumar, Rishabh Agarwal, Harshal Godhia, Ravi Rajwar, Javier Snaider, Xerxes 675 Dotiwalla, Yuan Liu, Aditya Barua, Victor Ungureanu, Yuan Zhang, Bat-Orgil Batsaikhan, Ma-676 teo Wirth, James Qin, Ivo Danihelka, Tulsee Doshi, Martin Chadwick, Jilin Chen, Sanil Jain, Quoc Le, Arjun Kar, Madhu Gurumurthy, Cheng Li, Ruoxin Sang, Fangyu Liu, Lampros Lam-677 prou, Rich Munoz, Nathan Lintz, Harsh Mehta, Heidi Howard, Malcolm Reynolds, Lora Aroyo, 678 Quan Wang, Lorenzo Blanco, Albin Cassirer, Jordan Griffith, Dipanjan Das, Stephan Lee, Jakub 679 Sygnowski, Zach Fisher, James Besley, Richard Powell, Zafarali Ahmed, Dominik Paulus, David 680 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 Butter-682 field, 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 684 Reid, Obaid Sarvana, David Silver, Alexander Chen, Lily Wang, Loren Maggiore, Oscar Chang, 685 Nithya Attaluri, Gregory Thornton, Chung-Cheng Chiu, Oskar Bunyan, Nir Levine, Timothy 686 Chung, Evgenii Eltyshev, Xiance Si, Timothy Lillicrap, Demetra Brady, Vaibhav Aggarwal, Boxi 687 Wu, Yuanzhong Xu, Ross McIlroy, Kartikeva Badola, Paramjit Sandhu, Erica Moreira, Wojciech Stokowiec, Ross Hemsley, Dong Li, Alex Tudor, Pranav Shyam, Elahe Rahimtoroghi, Salem 688 Haykal, Pablo Sprechmann, Xiang Zhou, Diana Mincu, Yujia Li, Ravi Addanki, Kalpesh Krishna, 689 Xiao Wu, Alexandre Frechette, Matan Eyal, Allan Dafoe, Dave Lacey, Jay Whang, Thi Avrahami, 690 Ye Zhang, Emanuel Taropa, Hanzhao Lin, Daniel Toyama, Eliza Rutherford, Motoki Sano, Hyun-691 Jeong Choe, Alex Tomala, Chalence Safranek-Shrader, Nora Kassner, Mantas Pajarskas, Matt 692 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 697 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 699 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,

702 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, 704 Tom Le Paine, Antoine Yang, Jason Riesa, Dominika Rogozinska, Dror Marcus, Dalia El Badawy, 705 Qiao Zhang, Luyu Wang, Helen Miller, Jeremy Greer, Lars Lowe Sjos, Azade Nova, Heiga Zen, 706 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, 708 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 710 Amersfoort, Vincent Zhuang, Aaron Cohen, Shixiang Shane Gu, Anhad Mohananey, Anastasija 711 Ilic, Taylor Tobin, John Wieting, Anna Bortsova, Phoebe Thacker, Emma Wang, Emily Caveness, 712 Justin Chiu, Eren Sezener, Alex Kaskasoli, Steven Baker, Katie Millican, Mohamed Elhawaty, 713 Kostas Aisopos, Carl Lebsack, Nathan Byrd, Hanjun Dai, Wenhao Jia, Matthew Wiethoff, El-714 naz Davoodi, Albert Weston, Lakshman Yagati, Arun Ahuja, Isabel Gao, Golan Pundak, Su-715 san Zhang, Michael Azzam, Khe Chai Sim, Sergi Caelles, James Keeling, Abhanshu Sharma, 716 Andy Swing, YaGuang Li, Chenxi Liu, Carrie Grimes Bostock, Yamini Bansal, Zachary Nado, Ankesh Anand, Josh Lipschultz, Abhijit Karmarkar, Lev Proleev, Abe Ittycheriah, Soheil Has-717 sas Yeganeh, George Polovets, Aleksandra Faust, Jiao Sun, Alban Rrustemi, Pen Li, Rakesh 718 Shivanna, Jeremiah Liu, Chris Welty, Federico Lebron, Anirudh Baddepudi, Sebastian Krause, 719 Emilio Parisotto, Radu Soricut, Zheng Xu, Dawn Bloxwich, Melvin Johnson, Behnam Neyshabur, 720 Justin Mao-Jones, Renshen Wang, Vinay Ramasesh, Zaheer Abbas, Arthur Guez, Constant Segal, 721 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 723 Hartman, Michael Kwong, Ruizhe Zhao, Sheleem Kashem, Elizabeth Cole, Antoine Miech, 724 Richard Tanburn, Mary Phuong, Filip Pavetic, Sebastien Cevey, Ramona Comanescu, Richard 725 Ives, Sherry Yang, Cosmo Du, Bo Li, Zizhao Zhang, Mariko Iinuma, Clara Huiyi Hu, Aurko Roy, 726 Shaan Bijwadia, Zhenkai Zhu, Danilo Martins, Rachel Saputro, Anita Gergely, Steven Zheng, 727 Dawei Jia, Ioannis Antonoglou, Adam Sadovsky, Shane Gu, Yingying Bi, Alek Andreev, Sina 728 Samangooei, Mina Khan, Tomas Kocisky, Angelos Filos, Chintu Kumar, Colton Bishop, Adams Yu, Sarah Hodkinson, Sid Mittal, Premal Shah, Alexandre Moufarek, Yong Cheng, Adam Blo-729 niarz, Jaehoon Lee, Pedram Pejman, Paul Michel, Stephen Spencer, Vladimir Feinberg, Xuehan 730 Xiong, Nikolay Savinov, Charlotte Smith, Siamak Shakeri, Dustin Tran, Mary Chesus, Bernd 731 Bohnet, George Tucker, Tamara von Glehn, Carrie Muir, Yiran Mao, Hideto Kazawa, Ambrose 732 Slone, Kedar Soparkar, Disha Shrivastava, James Cobon-Kerr, Michael Sharman, Jay Pavagadhi, 733 Carlos Araya, Karolis Misiunas, Nimesh Ghelani, Michael Laskin, David Barker, Qiujia Li, An-734 ton Briukhov, Neil Houlsby, Mia Glaese, Balaji Lakshminarayanan, Nathan Schucher, Yunhao 735 Tang, Eli Collins, Hyeontaek Lim, Fangxiaoyu Feng, Adria Recasens, Guangda Lai, Alberto 736 Magni, Nicola De Cao, Aditya Siddhant, Zoe Ashwood, Jordi Orbay, Mostafa Dehghani, Jenny Brennan, Yifan He, Kelvin Xu, Yang Gao, Carl Saroufim, James Molloy, Xinyi Wu, Seb Arnold, Solomon Chang, Julian Schrittwieser, Elena Buchatskaya, Soroush Radpour, Martin Polacek, 739 Skye Giordano, Ankur Bapna, Simon Tokumine, Vincent Hellendoorn, Thibault Sottiaux, Sarah Cogan, Aliaksei Severyn, Mohammad Saleh, Shantanu Thakoor, Laurent Shefey, Siyuan Qiao, 740 Meenu Gaba, Shuo yiin Chang, Craig Swanson, Biao Zhang, Benjamin Lee, Paul Kishan Ruben-741 stein, Gan Song, Tom Kwiatkowski, Anna Koop, Ajay Kannan, David Kao, Parker Schuh, Axel 742 Stjerngren, Golnaz Ghiasi, Gena Gibson, Luke Vilnis, Ye Yuan, Felipe Tiengo Ferreira, Aish-743 warya Kamath, Ted Klimenko, Ken Franko, Kefan Xiao, Indro Bhattacharya, Miteyan Patel, Rui 744 Wang, Alex Morris, Robin Strudel, Vivek Sharma, Peter Choy, Sayed Hadi Hashemi, Jessica 745 Landon, Mara Finkelstein, Priya Jhakra, Justin Frye, Megan Barnes, Matthew Mauger, Dennis 746 Daun, Khuslen Baatarsukh, Matthew Tung, Wael Farhan, Henryk Michalewski, Fabio Viola, Fe-747 lix de Chaumont Quitry, Charline Le Lan, Tom Hudson, Qingze Wang, Felix Fischer, Ivy Zheng, 748 Elspeth White, Anca Dragan, Jean baptiste Alayrac, Eric Ni, Alexander Pritzel, Adam Iwan-749 icki, Michael Isard, Anna Bulanova, Lukas Zilka, Ethan Dyer, Devendra Sachan, Srivatsan Srini-750 vasan, Hannah Muckenhirn, Honglong Cai, Amol Mandhane, Mukarram Tariq, Jack W. Rae, Gary Wang, Kareem Ayoub, Nicholas FitzGerald, Yao Zhao, Woohyun Han, Chris Alberti, Dan Gar-751 rette, Kashyap Krishnakumar, Mai Gimenez, Anselm Levskaya, Daniel Sohn, Josip Matak, Inaki 752 Iturrate, Michael B. Chang, Jackie Xiang, Yuan Cao, Nishant Ranka, Geoff Brown, Adrian Hutter, Vahab Mirrokni, Nanxin Chen, Kaisheng Yao, Zoltan Egyed, Francois Galilee, Tyler Liechty, 754 Praveen Kallakuri, Evan Palmer, Sanjay Ghemawat, Jasmine Liu, David Tao, Chloe Thornton, 755 Tim Green, Mimi Jasarevic, Sharon Lin, Victor Cotruta, Yi-Xuan Tan, Noah Fiedel, Hongkun

756 Yu, Ed Chi, Alexander Neitz, Jens Heitkaemper, Anu Sinha, Denny Zhou, Yi Sun, Charbel Kaed, Brice Hulse, Swaroop Mishra, Maria Georgaki, Sneha Kudugunta, Clement Farabet, Izhak 758 Shafran, Daniel Vlasic, Anton Tsitsulin, Rajagopal Ananthanarayanan, Alen Carin, Guolong Su, 759 Pei Sun, Shashank V, Gabriel Carvajal, Josef Broder, Iulia Comsa, Alena Repina, William Wong, 760 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 761 Zhang, Chu-Cheng Lin, Ionel Gog, Yana Kulizhskaya, Ashwin Sreevatsa, Shuang Song, Luis C. 762 Cobo, Anand Iyer, Chetan Tekur, Guillermo Garrido, Zhuyun Xiao, Rupert Kemp, Huaixiu Steven Zheng, Hui Li, Ananth Agarwal, Christel Ngani, Kati Goshvadi, Rebeca Santamaria-Fernandez, 764 Wojciech Fica, Xinyun Chen, Chris Gorgolewski, Sean Sun, Roopal Garg, Xinyu Ye, S. M. Ali 765 Eslami, Nan Hua, Jon Simon, Pratik Joshi, Yelin Kim, Ian Tenney, Sahitya Potluri, Lam Nguyen 766 Thiet, Quan Yuan, Florian Luisier, Alexandra Chronopoulou, Salvatore Scellato, Praveen Srini-767 vasan, Minmin Chen, Vinod Koverkathu, Valentin Dalibard, Yaming Xu, Brennan Saeta, Keith 768 Anderson, Thibault Sellam, Nick Fernando, Fantine Huot, Junehyuk Jung, Mani Varadarajan, 769 Michael Quinn, Amit Raul, Maigo Le, Ruslan Habalov, Jon Clark, Komal Jalan, Kalesha Bullard, 770 Achintya Singhal, Thang Luong, Boyu Wang, Sujeevan Rajayogam, Julian Eisenschlos, Johnson Jia, Daniel Finchelstein, Alex Yakubovich, Daniel Balle, Michael Fink, Sameer Agarwal, Jing Li, 771 Dj Dvijotham, Shalini Pal, Kai Kang, Jaclyn Konzelmann, Jennifer Beattie, Olivier Dousse, Diane Wu, Remi Crocker, Chen Elkind, Siddhartha Reddy Jonnalagadda, Jong Lee, Dan Holtmann-Rice, Krystal Kallarackal, Rosanne Liu, Denis Vnukov, Neera Vats, Luca Invernizzi, Mohsen 774 Jafari, Huanjie Zhou, Lilly Taylor, Jennifer Prendki, Marcus Wu, Tom Eccles, Tianqi Liu, Kavya 775 Kopparapu, Francoise Beaufays, Christof Angermueller, Andreea Marzoca, Shourya Sarcar, Hi-776 lal Dib, Jeff Stanway, Frank Perbet, Nejc Trdin, Rachel Sterneck, Andrey Khorlin, Dinghua Li, 777 Xihui Wu, Sonam Goenka, David Madras, Sasha Goldshtein, Willi Gierke, Tong Zhou, Yaxin 778 Liu, Yannie Liang, Anais White, Yunjie Li, Shreya Singh, Sanaz Bahargam, Mark Epstein, Sujoy 779 Basu, Li Lao, Adnan Ozturel, Carl Crous, Alex Zhai, Han Lu, Zora Tung, Neeraj Gaur, Alanna 780 Walton, Lucas Dixon, Ming Zhang, Amir Globerson, Grant Uy, Andrew Bolt, Olivia Wiles, Mi-781 lad Nasr, Ilia Shumailov, Marco Selvi, Francesco Piccinno, Ricardo Aguilar, Sara McCarthy, 782 Misha Khalman, Mrinal Shukla, Vlado Galic, John Carpenter, Kevin Villela, Haibin Zhang, Harry Richardson, James Martens, Matko Bosnjak, Shreyas Rammohan Belle, Jeff Seibert, Mah-783 moud Alnahlawi, Brian McWilliams, Sankalp Singh, Annie Louis, Wen Ding, Dan Popovici, 784 Lenin Simicich, Laura Knight, Pulkit Mehta, Nishesh Gupta, Chongyang Shi, Saaber Fatehi, Jo-785 vana Mitrovic, Alex Grills, Joseph Pagadora, Dessie Petrova, Danielle Eisenbud, Zhishuai Zhang, 786 Damion Yates, Bhavishya Mittal, Nilesh Tripuraneni, Yannis Assael, Thomas Brovelli, Prateek 787 Jain, Mihajlo Velimirovic, Canfer Akbulut, Jiaqi Mu, Wolfgang Macherey, Ravin Kumar, Jun Xu, 788 Haroon Qureshi, Gheorghe Comanici, Jeremy Wiesner, Zhitao Gong, Anton Ruddock, Matthias 789 Bauer, Nick Felt, Anirudh GP, Anurag Arnab, Dustin Zelle, Jonas Rothfuss, Bill Rosgen, Ashish 790 Shenoy, Bryan Seybold, Xinjian Li, Jayaram Mudigonda, Goker Erdogan, Jiawei Xia, Jiri Simsa, Andrea Michi, Yi Yao, Christopher Yew, Steven Kan, Isaac Caswell, Carey Radebaugh, Andre Elisseeff, Pedro Valenzuela, Kay McKinney, Kim Paterson, Albert Cui, Eri Latorre-Chimoto, 793 Solomon Kim, William Zeng, Ken Durden, Priya Ponnapalli, Tiberiu Sosea, Christopher A. Choquette-Choo, James Manyika, Brona Robenek, Harsha Vashisht, Sebastien Pereira, Hoi Lam, 794 Marko Velic, Denese Owusu-Afriyie, Katherine Lee, Tolga Bolukbasi, Alicia Parrish, Shawn Lu, Jane Park, Balaji Venkatraman, Alice Talbert, Lambert Rosique, Yuchung Cheng, Andrei 796 Sozanschi, Adam Paszke, Praveen Kumar, Jessica Austin, Lu Li, Khalid Salama, Wooyeol Kim, Nandita Dukkipati, Anthony Baryshnikov, Christos Kaplanis, XiangHai Sheng, Yuri Chervonyi, 798 Caglar Unlu, Diego de Las Casas, Harry Askham, Kathryn Tunyasuvunakool, Felix Gimeno, Siim 799 Poder, Chester Kwak, Matt Miecnikowski, Vahab Mirrokni, Alek Dimitriev, Aaron Parisi, Dan-800 gyi Liu, Tomy Tsai, Toby Shevlane, Christina Kouridi, Drew Garmon, Adrian Goedeckemeyer, 801 Adam R. Brown, Anitha Vijayakumar, Ali Elqursh, Sadegh Jazayeri, Jin Huang, Sara Mc Carthy, 802 Jay Hoover, Lucy Kim, Sandeep Kumar, Wei Chen, Courtney Biles, Garrett Bingham, Evan 803 Rosen, Lisa Wang, Qijun Tan, David Engel, Francesco Pongetti, Dario de Cesare, Dongseong Hwang, Lily Yu, Jennifer Pullman, Srini Narayanan, Kyle Levin, Siddharth Gopal, Megan Li, 804 Asaf Aharoni, Trieu Trinh, Jessica Lo, Norman Casagrande, Roopali Vij, Loic Matthey, Bramandia Ramadhana, Austin Matthews, CJ Carey, Matthew Johnson, Kremena Goranova, Rohin Shah, 806 Shereen Ashraf, Kingshuk Dasgupta, Rasmus Larsen, Yicheng Wang, Manish Reddy Vuyyuru, Chong Jiang, Joana Ijazi, Kazuki Osawa, Celine Smith, Ramya Sree Boppana, Taylan Bilal, Yuma 808 Koizumi, Ying Xu, Yasemin Altun, Nir Shabat, Ben Bariach, Alex Korchemniy, Kiam Choo, Olaf Ronneberger, Chimezie Iwuanyanwu, Shubin Zhao, David Soergel, Cho-Jui Hsieh, Irene Cai,

832 833

834

835

836

837

838

839

840

847

848

849

850

810 Shariq Iqbal, Martin Sundermeyer, Zhe Chen, Elie Bursztein, Chaitanya Malaviya, Fadi Biadsy, 811 Prakash Shroff, Inderjit Dhillon, Tejasi Latkar, Chris Dyer, Hannah Forbes, Massimo Nicosia, 812 Vitaly Nikolaev, Somer Greene, Marin Georgiev, Pidong Wang, Nina Martin, Hanie Sedghi, John 813 Zhang, Praseem Banzal, Doug Fritz, Vikram Rao, Xuezhi Wang, Jiageng Zhang, Viorica Pa-814 traucean, Dayou Du, Igor Mordatch, Ivan Jurin, Lewis Liu, Ayush Dubey, Abhi Mohan, Janek Nowakowski, Vlad-Doru Ion, Nan Wei, Reiko Tojo, Maria Abi Raad, Drew A. Hudson, Vaishakh 815 Keshava, Shubham Agrawal, Kevin Ramirez, Zhichun Wu, Hoang Nguyen, Ji Liu, Madhavi Se-816 wak, Bryce Petrini, DongHyun Choi, Ivan Philips, Ziyue Wang, Ioana Bica, Ankush Garg, Jarek 817 Wilkiewicz, Priyanka Agrawal, Xiaowei Li, Danhao Guo, Emily Xue, Naseer Shaik, Andrew 818 Leach, Sadh MNM Khan, Julia Wiesinger, Sammy Jerome, Abhishek Chakladar, Alek Wenjiao 819 Wang, Tina Ornduff, Folake Abu, Alireza Ghaffarkhah, Marcus Wainwright, Mario Cortes, Fred-820 erick Liu, Joshua Maynez, Andreas Terzis, Pouya Samangouei, Riham Mansour, Tomasz Kepa, 821 François-Xavier Aubet, Anton Algymr, Dan Banica, Agoston Weisz, Andras Orban, Alexandre 822 Senges, Ewa Andrejczuk, Mark Geller, Niccolo Dal Santo, Valentin Anklin, Majd Al Merey, 823 Martin Baeuml, Trevor Strohman, Junwen Bai, Slav Petrov, Yonghui Wu, Demis Hassabis, Koray 824 Kavukcuoglu, Jeffrey Dean, and Oriol Vinyals. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context, 2024. URL https://arxiv.org/abs/2403.05530. 825

- 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.
 - Jian Hu, Xibin Wu, Weixun Wang, Xianyu, Dehao Zhang, and Yu Cao. Openrlhf: An easy-to-use, scalable and high-performance rlhf framework, 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 Radev, Noah A. Smith, Yejin Choi, and Kentaro Inui. Realtime qa: What's the answer right now?, 2024. URL https://arxiv.org/abs/2207.13332.
 - Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. Rewardbench: Evaluating reward models for language modeling. arXiv preprint arXiv:2403.13787, 2024.
- Minghao Li, Yingxiu Zhao, Bowen Yu, Feifan Song, Hangyu Li, Haiyang Yu, Zhoujun Li, Fei
 Huang, and Yongbin Li. Api-bank: A comprehensive benchmark for tool-augmented llms, 2023.
 URL https://arxiv.org/abs/2304.08244.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step, 2023. URL https://arxiv.org/abs/2305.20050.
- Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, 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. Agentbench: Evaluating Ilms as agents, 2023. URL https://arxiv.org/abs/ 2308.03688.
- 863 Martin Majlis. Wikipedia-api, 2017. URL https://github.com/martin-majlis/ Wikipedia-API/tree/master.

- B64
 B65
 B66
 B66
 Dheeraj Mekala, Jason Weston, Jack Lanchantin, Roberta Raileanu, Maria Lomeli, Jingbo Shang, and Jane Dwivedi-Yu. Toolverifier: Generalization to new tools via self-verification. *arXiv* preprint arXiv:2402.14158, 2024.
- Grégoire Mialon, Clémentine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas
 Scialom. Gaia: a benchmark for general ai assistants, 2023. URL https://arxiv.org/ abs/2311.12983.
- 871 Mistral. Mistral large 2 model. https://mistral.ai. Accessed: 2024-10-01.
- 872 OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Floren-873 cia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red 874 Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher 875 Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, 877 Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, 878 Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey 879 Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, 882 Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gib-883 son, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hal-885 lacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook 889 Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel 890 Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen 891 Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel 892 Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, 893 Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv 894 Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, 895 Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, 896 Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Ra-897 jeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel 899 Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe 900 de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, 901 Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, 902 Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra 903 Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, 904 Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Sel-905 sam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, 906 Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, 907 Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, 908 Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vi-909 jayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan 910 Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, 911 Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Work-912 man, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming 913 Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao 914 Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL 915 https://arxiv.org/abs/2303.08774. 916
- 917 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Ed-

918 919 920	ward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library, 2019.
921 922 923	Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. Gorilla: Large language model connected with massive apis, 2023. URL https://arxiv.org/abs/2305.15334.
924 925 926	Yun Peng, Shuqing Li, Wenwei Gu, Yichen Li, Wenxuan Wang, Cuiyun Gao, and Michael Lyu. Revisiting, benchmarking and exploring api recommendation: How far are we?, 2021.
927 928 929 930	Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. Toolllm: Facilitating large language models to master 16000+ real-world apis, 2023.
931 932 933	Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools, 2023. URL https://arxiv.org/abs/2302.04761.
934 935 936	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL https://arxiv.org/abs/1707.06347.
937 938 939	SerpApi. Serpapi - search engine results api. https://serpapi.com/, 2024. Accessed: February 2023-September 2024.
940 941 942	James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. Fever: a large- scale dataset for fact extraction and verification, 2018. URL https://arxiv.org/abs/ 1803.05355.
943 944 945	Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, and Thang Luong. Freshllms: Refreshing large language models with search engine augmentation, 2023. URL https://arxiv.org/abs/2310.03214.
946 947 948 949	Peiyi Wang, Lei Li, Zhihong Shao, R. X. Xu, Damai Dai, Yifei Li, Deli Chen, Y. Wu, and Zhifang Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations, 2024a. URL https://arxiv.org/abs/2312.08935.
950 951 952	Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji. Exe- cutable code actions elicit better llm agents, 2024b. URL https://arxiv.org/abs/2402. 01030.
953 954 955	<pre>Inc. Wolfram Research. Mathematica, Version 14.1, 2024. URL https://www.wolfram. com/mathematica. Champaign, IL, 2024.</pre>
956 957 958	Congying Xia, Chen Xing, Jiangshu Du, Xinyi Yang, Yihao Feng, Ran Xu, Wenpeng Yin, and Caiming Xiong. Fofo: A benchmark to evaluate llms' format-following capability, 2024. URL https://arxiv.org/abs/2402.18667.
959 960 961	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.
962 963 964 965	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.
966 967 968 969	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.
970 971	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.

1026 A TOOLCOMP EXTENDED EVALUATIONS

In this appendix section, we include additional evaluations, namely the exact match grading (A.1) and error analysis for each model (A.2 and A.3).

1031 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 6: Model Family Performance Comparison: Accuracy and 95% Confidence Intervals

Model Family	Model Name	Total Accuracy (%)
	o1-preview	38.92 ± 4.36
	GPT-40 (Aug 2024)	43.52 ± 4.43
OpenAI	GPT-40 (May 2024)	40.60 ± 4.38
OpenAl	GPT-4 Turbo Preview	40.11 ± 4.39
	GPT-4	38.45 ± 4.34
	GPT-40 Mini	34.70 ± 4.25
	Claude 3.5 Sonnet	42.92 ± 4.42
Anthropic	Claude 3 Opus	36.96 ± 4.43
-	Claude 3 Sonnet	33.58 ± 4.21
Consta	Gemini 1.5 Pro (August 27, 2024)	43.22 ± 4.43
Google	Gemini 1.5 Pro (May 2024)	27.36 ± 3.98
Mistral	Mistral Large 2	33.63 ± 4.21
	Llama 3.1 405B Instruct*	33.10 ± 4.20
Meta	Llama 3.1 70B Instruct*	26.19 ± 3.93
	Llama 3.1 8B Instruct*	11.75 ± 2.88
Cohere	Command R+	0.00 ± 0.00

A.2 FINAL ANSWER FAILURE ANALYSIS

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 7 and the error counts for each model is shown in Figure 5.

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 consis- tently used the wrong formatting for the input argu- ments (e.g., wrong input format, didn't include a re- quired 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 un- clear 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 argu- ments 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.

Table 7: Common Error Category Taxonomy.

1134	Model Performance Comparison Across Error Types										
1135	GPT-4 -	9	1	78	83	10	5	83	2	39	- 160
1130	Gemini 1.5 Pro (Aug 2024) -	5	0	71	77	17	6	61	2	16	
1138	Llama 3.1 8B Instruct -	14	4	87	149	21	6	125	2	8	- 140
1139	Gemini 1.5 Pro (May 2024) -	12	1	69	69	14	7	62	3	129	
1140	Claude 3 Opus -	4	0	100	79	1	14	61	0	7	- 120
1141	Claude 2 Connet			100	104	-	-	77			
1142	Claude 3 Sonnet -	10	0	90	104	/	•	11	2	•	- 100
1143	GPT-4o Mini -	4	2	85	79	6	8	60	24	47	ount
1144	ው GPT-4o (Aug 2024) -	3	0	69	82	9	6	58	3	24	-80 Ŭ
1146	GPT-4 Turbo Preview -	2	1	77	84	13	6	65	4	4	ш
1147	Claude 3.5 Sonnet -	5	1	70	83	3	11	54	2	8	- 60
1148	Command R+ -	6	0	162	50	47	2	38	7	119	
1149	o1-preview -	2	0	72	73	1	4	54	0	17	- 40
1150	Llama 3.1 405B Instruct -	5	0	111	116	9	4	61	1	23	
1151	Mistral Large 2 -	2	0	42	28	2	0	39	1	10	- 20
1153	GPT-4o (May 2024) -	5	0	77	86	10	4	50	1	7	
1154		~	~	0	<u>`</u>	^'	6		~		- 0
1155	8	awec	×100.	ation.	nation	attinus	neous co	ol mile	e 100.	rediv	
1156	Plan	Incorre	Infor	AInfor	allForth	Miscello	, 00d 10	ngsan	Inexpe		
1157	Action	alled .	AISSING	inateu	00 ¹ Ca.	64.	nderst	Callin	arity		
1158		Swer	, Hallin	rect		MISI	. ated	ated	•		
1159		nalan		Inco'			Repe	Termin			
1160	4					Error Type	•	ž			
1161											

Figure 5: Breakdown of the various error categories in our taxonomy for each model (on the ToolComp-Enterprise).

1166 A.3 INTERMEDIATE REASONING FAILURE ANALYSIS

1167 1168 In this appendix section, we conduct a thorough failure analysis for the intermediate reasoning evaluations shown in Table 3.

1170 1171

A.3.1 REACT-STEP-ERROR-BASED FAILURE TRENDS IN MODELS

Figures 6 and 7 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:

1176 1177

1178 1179

1180

1181

1182

- 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 6: 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 7: 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.

1242 A.3.2 POSITION-BASED ERROR TRENDS IN MODELS

Figures 8 and 9 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 8: 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 9: Density of the position of the LLM as judge evaluation failures for each model. Full Benchmark denotes the distribution for the entire ToolComp benchmark.

¹²⁹⁶ B PERFORMANCE SCALING ON INCREASING COMPLEXITY

In this appendix section, we compare how process supervised reward model performance scales with more complex tool use prompts.

Categorizing Prompt Complexity We group ToolComp prompts into three categories of complexity — Easy, Medium, and Hard — based on the number of tool calls required in the human-verified trajectory (see Figure 1 for an overview on the annotation process) to answer the prompt:

- Easy: Prompts solved in 1–4 steps. (199 total prompts)
- Medium: Prompts solved in 5–8 steps. (210 total prompts)
- Hard: Prompts solved in 9–12 steps. (62 total prompts)

While there can be multiple valid trajectories of different lengths for the same prompt, we considerthis categorization a reasonable proxy for complexity.



Figure 10: A comparison of outcome-supervised and process-supervised reward models across various scales of training data (10%, 25%, 50%, 100%) on different complexity of prompts. We use the same 50 completions from the fine-tuned Llama-3.1-8b-Instruct generator in the experiments in Section 5 and we consider all 50 completions when ranking the trajectories.

PRM performance scaling is the greater for more complex prompts. Figure 10 compares the rank@1 performance scaling of ORM and PRM across the different prompt complexity. PRM consistently demonstrates better scalability for harder prompts, with the largest performance gains over ORM observed in the Hard category. This highlights that PRMs are particularly effective in handling more complex queries requiring sophisticated reasoning across multiple tool call steps.

1350 C TOOLCOMP DETAILS

In this appendix section, we provide further details regarding benchmark creation steps such as prompt creation (C.1, C.2, C.3). We also provide additional benchmark metadata revolving different characteristics and statistics about the benchmark (C.4).

1356 1357 C.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 C.2

1362 1363

1374

1380 1381

1382

1384 1385

1386

1387

1388

1389 1390

1391 1392

1393

1394

1395

1399

1400

1401

1402

1403

1352

1353

1354

1355

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 C.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.

C.2 IN CONTEXT EXAMPLE

1383 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.

Process

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.

1404 C.3 SEED PROMPT

ſ	I want you to act as a Prompt Writer.							
	Please adhere to the following instructions:							
	Thease deficite 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:							
	(avamalas)							
	{examples}							
	[BEGIN ALLOWED TOOLS]							
	{tools}							
	[END ALLOWED TOOLS]							
С	.4 BENCHMARK METADATA							
	140 -							
	120 -							
	≥ ¹⁰⁰							
	^{сс} ₆₀ -							
	40 -							
	20 -							
	# of Tools							
C:	gure 11: About 85% of prompts in ToolComp require at least 3 tool calls to solve indicating th							
1'1 71	guie 11. About 65 % of prompts in Tool comp require at least 5 tool cans to solve, indicating in							

Figure 11: 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.

- 1456
- 1457



Figure 13: 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.

- 1509
- 1510
- 1511



Figure 16: The distribution of tools called in our human supervised tool call chains for just the ToolComp Chat subset.



Figure 17: 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.

- 1588
- 1589 1590

D PROCESS SUPERVISION VS. OUTCOME SUPERVISION TRAINING DETAILS

1591 D.1 DETAILED SYNTHETIC TRAINING DATA GENERATION

Synthetic Prompt Generation For the generation of synthetic prompts, we mirror the strategies outlined in Section 3.3 with the notable exception of the Human Refinement step due to the high level of associated cost. Instead, we replace this step with final answer consistency across different model families. Using the following models – GPT-40 (May 2024), GPT- 4 Turbo, Claude 3.5 Sonnet, and Llama 3.1 70b – we generate full trajectories and only keep the prompts for which every model arrives at the same final answer. From empirical evaluation, this serves as a good proxy for unambiguous and sensible prompts. Table 8 notes the initial amount of prompts generated and the final number of prompts that are final answer consistent across the model families.

Table 8: Count of training data through the different stages of generation.

Туре	Count
Initial Prompts	75K
Final Answer Consistent Prompts	17369
Trajectories with Final Answers	13628
Trajectories with Correctly Formatted Final Answers	11654

1610 **Synthetic Training Data** Suppose we have a LLM that acts as a policy, coined the Policy Model, 1611 and another LLM that acts as a judge, coined the Critic Model. Given a query, q, we first use 1612 the Policy Model to generate an action plan, a. Then, conditioned on the action plan, we prompt 1613 the Policy Model to generate a full chain trajectory $\{t_1, \ldots, t_N\}$, where each t_i is ReAct step that 1614 invokes a tool call or invokes the finish action. The prompt to generate the action plan and each 1615 tool call is given in E.1 and E.2, respectively. We bound N by 15, allowing at most 15 tool call in 1616 trajectory. If the model reaches a final answer, then t_N is the finish action. We then use the Critic Model to critique the action plan a and react steps t_i . The prompt for the Critic Model to critique 1617 the action plan and the react steps is given in E.3 and E.4, respectively. In the case that the critique 1618 model finds a fault with a step, it then proposes a corrected step. The corrected step is then used to 1619 continue the chain. Using the latest correct step, the Policy Model will then be invoked to generate

1620 a new ReAct step and the Critic Model will critique the step iteratively until either the Policy Model 1621 reaches a final answer and the Critic Model agrees or the Critic Model proposes a final answer when 1622 correcting a step. All Policy Model full trajectories that did not reach a final answer was discarded 1623 and all Critic Model trajectories that did not reach a final answer was also discarded. For the Policy 1624 Model, we used Llama-3.1-8b-Instruct, and for the Critic Model, we used GPT-40. The process described here mimics the same process we used to create ToolComp described in Figure 1, where 1625 human annotators is replaced by GPT-40 and the policy is Llama-3.1-8b-Instruct. Table 9 shows the 1626 different generation parameters used for the Policy Model and the Critic Model. From Table 8, the 1627 total number of trajectories with correctly formatted final answers is 10342. 1628

1629 1630

1633

1635

1637

1639

Table 9: Policy Model and Critic Model Generation Parameters

Parameter	Policy Model	Critic Model
Temperature	0.5	0
Max New Tokens	1024	4096
Num Retries Per Step	3	3
Stop	"End Action", "End Action\n", "\nEnd Action"	"End Action"

1640

1641 D.2 CONSTRUCTION OF TRAINING DATASETS

1642 **Reward Model Datasets** The preference ranking dataset for outcome reward modelling (ORM) is 1643 composed of the dis-preferred trajectory being the original Policy Model full trajectory and preferred 1644 trajectory being the fully corrected Critic Model trajectory (assuming at least one edit was made in 1645 the duration of the Critic Model trajectory). The step-by-step preference ranking dataset for the 1646 process reward modelling (PRM) is comprised of every Critic Model correction, we take the dis-1647 preferred step to be the original Policy Model's full ReAct step and the preferred step to be the 1648 Critic Model's full ReAct corrected step. The trajectory history is the most correct chain leading 1649 up to the given step. In total, as in Table 10, the ORM dataset size is same as the total number of 1650 trajectories with correctly formatted final answers (11654) and the PRM dataset size is equivalant to the number of times the Policy and Critic models had disputes and corrections (14932). We note 1651 that the ORM samples are full trajectories, whereas the PRM samples are simply the steps that were 1652 correct. 1653

 1654
 1655
 1656
 1656
 SFT Dataset The SFT dataset is comprised of the most corrected trajectory taken by simply following all Critic Model correct steps where applicable.

1657 1658	Table 10:	Training dataset sizes.
1659	Dataset	Number of Samples
1661	ORM	11654
1662 1663	SFT	14932 11654

D.3 DETAILED TRAINING OBJECTIVES

Outcome Reward Model The base model is equipped with a linear layer to serve as the reward head, which outputs a scalar reward value per token. We place a special RM token at the end of the preferred and dis-preferred completion, followed by an EOS token. The training objective is then given by a Binary Cross Entropy loss on the reward value output on the special RM token. We use Llama-3.1-8b-Instruct as the base model.

1672

1664

1666

Process Reward Model Similarly, the base model is equipped with a linear layer to serve as the reward head, which outputs a singular scalar value. We test 4 different levels of supervision. For the

Full Step with Observation and Full Step without Observation, a singular special RM token is placed at the end of the observation step and at the end of the action input step, respectively. Moreover, for Sub Step with Observation and Sub Step without Observation, 4 special RM tokens are placed after each ReAct step (including the observation) and 3 special RM tokens are placed after each ReAct step (excluding the observation), respectively. For all variation, the training objective is given by the average Binary Cross Entropy loss across the reward value output at each of the special RM tokens and the corresponding label. We use Llama-3.1-8b-Instruct as the base model.

1681 1682

D.4 PRM vs. ORM EVALUATION

1683 **Ranking Trajectories** We use a generator to generate 50 completions per problem in ToolComp. 1684 Then for each completion, we discard all completions that did not reach a final answer. We then use 1685 an ORM and PRM to rank the completions per problem. For the ORM, we simply take the sigmoid 1686 of the reward score it places on the completions as that is the probability the model assigns that the 1687 trajectory is correct. For PRM, given a trajectory, we collect a list of probability scores for each 1688 step. To combine the list of probability scores into a single number, we experiment with a couple of 1689 aggregation functions, namely: min, max, average, and product. In order to account for variance in the ranking scores, we perform 500 permutations of the completion per problem, leaving 20 out at random each permutation. We consider the best-of-30 accuracy, which takes the best trajectory as ranked by the corresponding method and evaluates the correctness of that single trajectory.

Training Dataset Scales In order to assess the generalizability of PRM vs. ORM with increasing scales of data, we vary the dataset sizes to be 10, 25, 50, and 100 percent of the full dataset. At each dataset scale, we train the models for 3 epochs and then take the best performing model on a held out synthetically generated validation set.

Base Model vs. SFT Model Generations To account for a variety in the quality of Tool-Use trajectories to rank and to assess the performance gains beyond SFT, in addition to using Llama-3.1-8b-Instruct model as a generator, we also use a Supervise Fine-Tuned Llama-3.1-8b-Instruct. The model is trained on the full set of golden trajectories given by following all the Critic Model's corrections.

1703

1693

1704 D.5 TRAINING IMPLEMENTATION AND HYPER-PARAMETERS

Implementation We implement the reward model and SFT training using the OpenRLHF library (Hu et al., 2024) in combination with PyTorch (Paszke et al., 2019). We made the following modification to the OpenRLHF Library:

1709 1710

1711

1712

• The RewardDataset object optionally took indices that signified where to place the special RM tokens to implement the ORM training and the 4 different PRM supervision training

- We update the BCE loss function implementation to additional consider more than one RM classification token id for the PRM training. We simply average the BCE loss for the multiple settings.
- 1713 1714 1715

1716Training Parameters and SelectionWe use Llama-3.1-8b-Instruct weights as initialization for1717both the ORM and the PRM. For each of the ORM and PRM training across the different dataset1718scales, we keep a constant batch size of 64 with 3 epochs and experiment with 5 different learning1719rates: 1e-5, 1e-6, 9e-7, 6e-7, and 1e-7. We then pick the best respective model using the best1720performance on a held-out (10%) validation set. Generally speaking, the best learning rate for PRM1721training is 9e-7 and for ORM training is 1e-6 across dataset scales. For the SFT training we1722experiment with 3 different learning rates 1e-5, 1e-6, and 1e-7 and 5 epochs. Table 11 and Table172412 summarizes hyper-parameters and training settings.

1723

1724

1725

729		
730	hyper-parameter and settings	value
731	anochs	2
732	context length	16384
33	batch size	64
34	deepspeed zero stage	2
35	seed	42
36	loss function	BCE
37		10-5 $10-6$ $00-7$
38	learning rates	$6e_7 1e_7$
39		0e=7, 1e=7
40		
41		
42	Table 12: Training set	tings for SFT.
43		
44	hyper-parameter and settings	value
45	epochs	5
46	context length	16384
17	batch size	64
18	deepspeed zero stage	3
9	seed	42
50	loss function	CE
51	learning rates	1e-5, 1e-6, 1e-7
52	e	
53		

Table 11: Training settings for ORM and PRM.

1755 E TOOL-USE PROMPTS

1757

1766 1767

1768 1769 1770

1771

1772 1773

1774 1775

1776

1728

1758 In this section, we summarize all of the prompts that were used during the creation of the benchmark, 1759 evaluation of the benchmark, and creation of the synthetic training data. For the creation of the 1760 benchmark, we state the "Action Plan Prompt" for the Policy Model in Section E.1 and the "Tool Call Prompt" for the Policy Model in Section E.2. For the evaluation of the benchmark, we state the 1761 LLM grading prompt and the in-context examples used to aid grading in Section E.5. Lastly, for the 1762 creation of the synthetic training data, we use the same policy model prompts for the action plan and 1763 tool call, and we additionally include the "Action Plan Prompt" for the Critic Model in Section E.3 1764 and the "Tool Call Prompt" for the Critic Model in Section E.4. 1765

E.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.

1782 E.2 REACT TOOL CALL PROMPT (POLICY MODEL)

1784	
1785	SYSTEM:
1786	
1787	You are a helpful assistant with access to functions, each function will be regarded as an
1788	action. Your job is to take relevant and necessary actions to get to the final answer to a user
1789	question. Please use the actions to provide information accurate up to current date and time:
1790	{current_date}. The user will provide you a question and a high level action plan. Your job
1791	is to execute on the action plan to answer the question. It's okay to slightly deviate from the
1792	action plan if you unlik it's necessary.
1793	FUNCTIONS: {func spec}
1794	
1795	Please stick to the following format:
1796	
1797	Thought: \langle your reasoning/thought on why/how to use an action \rangle
1798	Action: (the action to take, should be one of $\{func_list\}$)
1799	Action Input: (the input to the action (should be in JSON format with the required fields)) End Action
1800	
1801	If you believe that you have obtained enough information (which can be judged from the
1802	history observations) to answer the question, please call:
1803	
1804	Thought: I have enough information to answer the question
1805	Action: finish
1000	Action Input: {"answer": [your answer string]}}
1808	End Action
1809	For your final answer (the finish action input) make sure you answer the full question
1810	Additionally we want to make sure the final answers/outputs in the finish action input are
1811	returned in the order that they are given in a list format so we can verify them with an exact
1812	string match. For eg. if the prompt asks for a city name, its temperature and a list of names
1813	of all the NBA teams whose home stadium is within a 400 mile radius, you would output
1814	['San Francisco', 78, ['Los Angeles Lakers', 'Golden State Warriors']].
1815	
1816	If the prompt asks for a special sorting of the list, make sure to output wrap the list in {{}}
1817	prompt instead asked to list the names of all the NBA teams whose home stadium is within
1818	a 400 mile radius in alphabetical order you would output [San Francisco 78 {{Golden
1819	State Warriors, Los Angeles Lakers}}].
1820	
1821	Only output the final answer with no additional text or natural language. Give dates in
1822	YYYY-MM-DD format, temperatures in celcius, prices in dollars, lengths in meters, area in
1823	meters ² , volume in m^3 and angles in degrees if the prompt doesn't specify what format/units
1824	to output the answer in.
1825	observations take your next action
1826	observations, take your next action.
1827	USER:
1828	
1829	Question: {question}
1830	
1831	Action Plan: {action_plan}
1832	A SSISTA NT.
1833	ASSISTANT.
1834	{history_of_react_steps}
1835	

1000		
1030	F 3	ACTION PLAN PROMPT (CRITIC MODEL)
	L .J	Merion I EAN I Rown I (CRITIC MODEL)

837
838
839
840
841
842
843
844
845
846
847
848
840
850
951
051
852
954
855
956
050
001
000
009
000
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887

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}.
UNCTIONS: {func_spec}
Question: {question}
Given the question and the tools available to you above, please formulate an action plan to nswer 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 teps you would take to answer the question and the name of the tool you would use and yow 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?
lease err on the side of giving the assistant the benefit of the doubt, and only critique the ction plan if it is clearly incorrect. f the action plan is incorrect, provide a revised action plan that you believe would be correct
Furthermore, your output should follow the format:
Reasoning: \langle your reasoning for the correctness or incorrectness of the action plan \rangle
Label: (correct/incorrect)
Revised Action Plan: \langle your revised action plan or empty if no revision needed \rangle
Here is the action plan provided by the assistant:
action_plan}
Please provide your critique of the action plan.

1890 E.4 REACT TOOL CALL PROMPT (CRITIC MODEL) 1891

1892	
1893	You are an expert judge of tool calls. Your job is to critique each of the ReAct steps of an
1894	assistant
1895	
1896	The following information is shown to the assistant in order to devise a ReAct step.
1897	
1898	[Start of the message]
1899	
1900	You are a helpful assistant with access to functions. Use them if required. Please use the
1901	tools to provide information accurate up to current date and time: {current_date}.
1902	
1903	FUNCTIONS: {func_spec}
1904	Please stick to the following format:
1905	riedse suck to the following format.
1906	Thought: you should always think about what to do
1907	Action: the action to take, should be one of {func_list}
1908	Action Input: the input to the action
1909	End Action
1910	
1911	If you believe that you have obtained enough information (which can be judged from the
1912	history observations) to answer the question, please call:
1913	
1914	Thought: I have enough information to answer the question
1915	Action: Innish Action Input: "answer": [your answer string]
1916	End Action
1917	
1918	Ouestion: {auestion}
1919	
1920	[End of the message]
1921	
1922	Given the set of functions, question, action plan and history of past actions, critique the
1923	Thought, Action, and Action Input step. Assume the action plan and history of past
1924	actions are optimal. To assess the thought step, if the step is roughly reasonable and the
1925	action and action input step are correlated with the thought step, then the thought step
1926	is correct. Please give the assistant the benefit of the doubt and be rement in your assessment.
1927	To assess the action step, let's assume that the Assistant cannot complete simple function-
1928	alities such as simple arithmetic, converting units, or utilizing simple facts without the
1929	use of tools. If the action specifies a reasonable function to use, then the action step is correct.
1930	
1931	To assess the action input step, if the input is reasonable and the action is correct, then the
1932	action input step is correct.
1933	
1934	If any of the steps are incorrect, label them as incorrect in the Labels section.
1935	For the Deviced De Act Stan section, movide the second star that the second start is 111
1936	For the Revised ReAct Step section, provide the assistant's step as the revised step. If the
1937	assistant's step is incorrect, provide the correct step that the assistant should have taken. As
1938	a general rule of thumb if your revised step is different from the assistant's step, then the
1939	assistant's step is incorrect, and if your revised step is the same as the assistant's step, then the
1940	the assistant's step is correct.
1941	
1942	As an important reminder, for your final answer (the finish action input), we want to make
1943	sure the final answers/outputs in the finish action input are returned in the order that they

1944 are given in a list format so we can verify them with an exact string match. For eg. if the 1945 prompt asks for a city name, its temperature and a list of names of all the NBA teams whose 1946 home stadium is within a 400 mile radius, you would output ['San Francisco', 78, ['Los 1947 Angeles Lakers', 'Golden State Warriors']]. If the prompt asks for a special sorting of the 1948 list, make sure to output wrap the list in $\{\}\}$ and if doesn't require any special sorting 1949 wrap it in [] like you normally would. So if the prompt instead asked to list the names 1950 of all the NBA teams whose home stadium is within a 400 mile radius in alphabetical or-1951 der, you would output [San Francisco, 78, {{Golden State Warriors, Los Angeles Lakers}}]. 1952 1953 Only output the final answer with no additional text or natural language or units. Give dates 1954 in YYYY-MM-DD format, temperatures in Celcius, prices in dollars, lengths in meters, area in meters², volume in m^3 and angles in degrees if the prompt doesn't specify what format/units to output the answer in. 1957 As a reminder, you should not use an external information that is not provided in the 1958 prompt or by a tool call. As a simple example, you may know a ticker symbol already for 1959 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. 1961 Your output should follow the format: 1963 1964 [Start of format] 1965 1966 Reasoning: \langle your reasoning for the correctness or incorrectness of each step \rangle 1967 Labels: [$\langle correct/incorrect \rangle$, $\langle correct/incorrect \rangle$] (in the order of 1968 Thought, Action, Action Input) 1969 1970 Revised ReAct Step: 1971 1972 Thought: $\langle \text{ your revised thought or assistant's thought if correct } \rangle$ Action: \langle your revised action or assistant's action if correct \rangle 1974 Action Input: \langle your revised action input or assistant's action input if correct \rangle 1975 End Action 1976 [End of format] 1977 1978 Here is the action plan: {action_plan} 1981 1982 Here is the history of past actions. If there are no past actions yet, this will be empty: 1983 1984 {history} Here is the latest ReAct step provided by the assistant: 1987 1988 Thought: {thought} Action: {action} 1989 Action Input: {action_input} End Action Observation: {observation} 1993 Please provide your critique of the latest ReAct step provided by the assistant.

1998 E.5 LLM GRADING PROMPT

E.5.1 MAIN PROMPT

2000

2002

2003 2004 2005 2006 You are an expert test grader. You have been given a student answer ('Student Answer:') to 2007 grade. You have also been the correct answer ('Correct Answer:') and the original question 2008 ('Question:'). Each correct answer is a list of strings. 2009 {In-Context Examples} 2011 The possible grades are 2012 2013 **INCORRECT:** 'Student Answer:' is different from 'Correct Answer:' 2014 2015 numbers are completely different 2016 lists are completely different 2017 • 'Question:' asks for special sorting of a list but the list in 'Student Answer:' is 2018 sorted differently than 'Correct Answer:' 2019 • strings are completely different or information present in the string is completely different 2021 CORRECT BUT BAD FORMATTING: 'Student Answer:' has the same info as 'Correct 2023 Answer:' but is formatted differently. 2024 2025 'Student Answer:' includes natural language or additional text 2026 • numbers are formatted differently but they are close to one another ('Student An-2027 swer:' is within 2028 • lists are wrapped differently than the correct answer but contains the same infor-2029 mation and sorted the same way as 'Correct Answer:' if asked 'Question:' asks for 2030 a special sorting 2031 • Strings are the same but may be formatted differently 2032 **CORRECT**: The student answer has the same info as 'Correct Answer:' and is also formatted the same as 'Correct Answer:' 2035 • numbers are close to one another ('Student Answer:' is within 10% of the correct 2036 answer) 2037 2038 • if 'Question:' asks for a special sorting of the list the 'Student Answer:' list is sort the same as 'Correct Answer:' 2039 2040 lists are wrapped the same 2041 · Strings are identical 2042 2043 Remember you are assuming the correct answer provided is correct, your job is is only to 2044 compare the correct answer to the student answer and grade it based on the above criteria. 2045 Do not try to determine the correct answer yourself. Make sure to include a reasoning and 2046 final grade in the format: 2047 Reasoning: (reasoning) Final Grade: (INCORRECT / CORRECT BUT BAD FORMAT-2049 TING / CORRECT > [ENDOFGRADE] 2050 2051 Now do this for the following user provided question, student answer and correct answer.

2052 E.5.2 IN-CONTEXT EXAMPLES (ORDERING)

_	
W	want to make sure the values in the student answer are returned in the order that they
are	e want to make sure the values in the student answer are returned in the order that they
ure	doked in Question.
Fo	r 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
'St	udent Answer:' can be ['San Francisco', 78, ['Los Angeles Lakers', 'Golden State
Wa	arriors']].
E.	amplace
Ľл	ampies.
Or	estion : Find the name of the city known for its famous tourist attraction Alcatraz, also
giv	e it's current temperature and a list of names of all the NBA teams whose home stadium
is v	within a 400 mile radius
Co	rrect Answer: ['San Francisco', 78, ['Golden State Warriors', 'Los Angeles Lakers']
Stu	ident Answer: ['San Francisco', 74, ['Los Angeles Lakers', 'Golden State Warriors']]
Re	asoning: The Student Answer is correct because it identifies the same city, the
ten	nperature is within 10% of the Correct Answer and the same team names are present in
.ne Fi	nsi. nal Grade: CORRECT
. 11	
Or	estion: Find the name of the city known for its famous tourist attraction Alcatraz, also
giv	e it's current temperature and a list of names of all the NBA teams whose home stadium
is v	within a 400 mile radius
Co	rrect Answer: ['San Francisco', 78, ['Golden State Warriors', 'Los Angeles Lakers']
Stı	Ident Answer: The city name is San Francisco, its temperature is 80 degrees and the Los
An	geles Lakers and the Golden State Warriors are two NBA teams whose home stadium is
Wil	thin a 400 mile radius
Ke	asoning: Although the Student Answer is correct (identifies the same city, the temper-
atu for	me is within 10% of the Context Answer and the same team names are present), it's not matted the same and contains extra text and natural language
Fir	al Grade: CORRECT BUT BAD FORMATTING
Qu	estion: Find the name of the city known for its famous tourist attraction Alcatraz, also
giv	re it's current temperature and a list of names of all the NBA teams whose home stadium
is v	within a 400 mile radius
Co	rrect Answer: ['San Francisco', 78, ['Golden State Warriors', 'Los Angeles Lakers']
Sti D	Ident Answer: [San Francisco', -15, [Los Angeles Lakers', Golden State Warriors']]
Ke	asoning: The Student Answer is incorrect because although identifies the same city and same team names are present in the list, the temperature is well outside of 100% of the
Co	rect Answer
Fi	al Grade: INCORRECT
.3	IN-CONTEXT EXAMPLES (SORTING)
If	Ouestion:' 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).
	Fillar 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
j	is within a 400 mile radius in alphabetical order
(Correct Answer: ['San Francisco', 78, ['Golden State Warriors', 'Los Angeles Lakers']
5	Student Answer: The city name is San Francisco, its temperature is 80 degrees and
t	he Golden State Warriors and the Los Angeles Lakers are two NBA teams whose home
5	stadium is within a 400 mile radius (in alphabetical order)
]	Reasoning : Although the Student Answer is correct (identifies the same city, the tempera-
i	ture is within 10% of the Correct Answer and the same team names are present in the same
(Since Crode: CORRECT BUT RAD FORMATTING
-	FINAL GLAUE. CORRECT BUT DAD FORMATTING
1	Ouestion : 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

2160 F TOOLCOMP EXAMPLES 2161

F.1 EXAMPLE GOLDEN TRAJECTORY

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

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.

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)}

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:

```
import math
  box_length = 2
2
  box_width = 2
3
  box_height = 2
4
   # convert truck dimensions into inches
5
  truck_length = (26 * 12) + 5
6
  truck_width = (7 * 12) + 8
7
  truck_height = (8 * 12) + 3
8
9
  num_boxes = math.floor(truck_length // box_length) * \
10
               math.floor(truck_width // box_width) * \
               math.floor(truck_height // box_height)
  print(f'Total number of boxes: {num_boxes}')
12
 • Observation:
                    {"result":
                                 "Total number of boxes:
                                                              356132",
   "error":
             null}
```

ReAct Step 3

• Thought: I have the total number that as the answer.	er of boxes that can fit inside the truck, so I should return
• Action: finish	
• Action Input: {"answer":	" <u>356132</u> "}
Observation: 356132	

F.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?

Action Plan

2216

2217

2218 2219

2220

2222

2223

2224

2225

2226

2227

2228

2229

2230

2231

2232

2233

2234 2235

2236

2237 2238

2239 2240

2241

2242

2243

2244

2245

2246

2247

2248

2249

2250

2251

2252

2253

2254

2255

2256

2257

2258

2259

2260

2261

2262

2263

2264

2265

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)

ReAct Step 1

Model

- **Thought**: I need to use the google_search tool to calculate the population growth rate using the birthrate, 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: X)
```

• 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 birth-rate, 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":
 "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]
```

2268 ReAct Step 2

Model	Human
 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 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] 	 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]

ReAct Step 3

Aodel	• Thought: N/A
Thought : I need to use the python inter-	Action: N/A
years using the population growth rate and the population of Japan in 2023.	Action Input:
(Label: ✓)	population =
Action: python_interpreter (Label: ✓)	11.7 total_bi
Action Input:	* (population total_deaths_
<pre>import numpy as np population = 124352000</pre>	= death_rate_; * (population
<pre>population_growth_rate = 0.0012</pre>	population_gr (total_births
<pre>population_in_10_years = population * (1 +</pre>	<pre>- total_death / population</pre>
<pre>population_growth_rate) **10 print (population_in_10_years)</pre>	population_in = population
(Label: 🗡)	population_gr
Observation: 125852307.84945872	print (populat
	Observation: 11

```
24620000
thousand =
ths_extrapolated
er_thousand
/ 1000)
xtrapolated
er_thousand
/ 1000)
wth_rate =
extrapolated
_extrapolated)
10_years
(1 +
wth_rate)**10
on_in_10_years)
08275.35965854
```

ReAct Step 4

Model
• Thought : I have enough information to answer the question (Label : ✓)
• Action: finish (Label: ✓)
<pre>• Action Input: {"answer": ["118,408,275"]} (Label: X)</pre>
• Observation: ["118, 408, 275"]



2376 G 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.

G.1 DATE

Description: Returns the current date (e.g., January 1, 2024).

2384 2385 Input Example:

2386 1 2387

2389

2390

2393 2394

2391 ²

2391 3 **2392** 4

2381

2382

2388 Output Example:

{ }

{

}

[]

{

{

}

```
"result": "Today is Friday, February 16, 2024.",
"error": ""
```

Description: Calculates expressions including basic arithmetic and brackets.

"operation": "2*32-4+456+(1+2)+3+(1/2*3+3+(1+2))"

Parameters:

G.2 CALCULATOR

Input Example:

Output Example:

"error": "",

"result": "529.5"

2395 2396 2397

2398

2399

2400 2401

2402

2403 1

2404 2 **2405** ³

3 }

2406

2407

2408

2409 ₂ **2410** ₃

2411 4 2412

2413

Parameters:

```
2414
    1
       [
2415
    2
          {
2416<sub>3</sub>
            "Input Name": "operation",
            "Type": "String",
2417 4
            "Description": "Computes numerical expressions involving float
2418 5
                numbers and operators like +, -, *, /, ^.""
2419
         }
    6
2420
       ]
```

2424 2425

2426

2427

2428

2430 G.3 CURRENT WEATHER

2432 Description: Retrieves current daily averages for temperature, rainfall, and hours of precipitation
 2433 for a specified city and country. Does not return historical data.

2435 Input Example:

{

{

```
2436
```

2440

2434

```
2437 2 "city_name": "London",
2438 3 "country_code": "GB"
2439 4 }
```

Output Example:

"error": "",

```
2442
2443 1
2444 <sup>2</sup>
2445 <sub>4</sub>
```

```
"result": [
     3
2445
     4
              {
2446 5
                 "date": "2024-03-25 00:00:00",
2447 6
                 "temperature (F)": "47.78615"
                 "total rain (mm)": "1.4000001",
2448<sup>7</sup>
                 "total snowfall (mm)": "0.0",
2449<sup>8</sup>
2450 <sup>9</sup><sub>10</sub>
                 "precipitation hours (hours)": "4.0"
              },
2451<sub>11</sub>
               {
                 "date": "2024-03-26 00:00:00",
2452 12
                 "temperature (F)": "48.374897",
2453<sup>13</sup>
                 "total rain (mm)": "8.2",
2454<sup>14</sup>
2455<sup>15</sup><sub>16</sub>
                 "total snowfall (mm)": "0.0",
                 "precipitation hours (hours)": "11.0"
2456<sub>17</sub>
              },
2457 18
               {
                 "date": "2024-03-27 00:00:00",
2458<sup>19</sup>
                 "temperature (F)": "47.217274",
2459<sup>20</sup>
                 "total rain (mm)": "2.3999999",
2460<sup>21</sup><sub>22</sub>
                 "total snowfall (mm)": "0.0"
2461 23
                 "precipitation hours (hours)": "4.0"
2462 24
               }
            ]
2463<sup>25</sup>
```

2465 2466 Parameters:

}

2464²⁶

```
2467 1
        ſ
2468 <sub>2</sub>
             "Input Name": "city_name",
2469 3
             "Type": "String",
2470<sup>4</sup>
             "Description": "The name of the city."
2471 <sup>5</sup>
     6
           },
2472 7
           {
2473 8
             "Input Name": "country_code",
             "Type": "Two Alphabet-Number",
2474 9
             "Description": "The country code (ISO 3166-2). The list can be found
2475<sup>10</sup>
                 here: https://en.wikipedia.org/wiki/ISO_3166-2"
2476
    11
           }
2477<sub>12</sub>
        1
2478
2479
2480
2481
2482
2483
```

2484 G.4 HISTORICAL WEATHER 2485

2486 **Description:** Retrieves daily averages for temperature and precipitation starting from the 1940s for 2487 a given city. Note: 5-day data delay, meaning you cannot get current weather data for the last 5 days. 2488

Input Example:

2493

2494

2495 2496

2497

2515 2516

2517

{

```
{
          "city_name": "London",
          "country_code": "GB",
2492<sup>3</sup>
          "start_date": "2023-03-09",
    4
          "end_date": "2023-03-21"
       }
```

Output Example:

```
2498
```

```
2499 <sub>2</sub>
            "error": "",
             "result": [
2500 3
2501 4
               {
                  "date": "2024-03-09 00:00:00",
2502 <sup>5</sup>
     6
                  "temperature (F)": "48.102356",
2503
                  "total rain (mm)": "0.4",
      7
2504 8
                  "total snowfall (mm)": "0.0",
2505 9
                  "precipitation hours (hours)": "2.0"
               },
2506<sup>10</sup>
2507<sup>11</sup>
              . . .
2508<sup>12</sup><sub>13</sub>
               {
                  "date": "2024-03-23 00:00:00",
2509<sub>14</sub>
                  "temperature (F)": "43.373596",
                  "total rain (mm)": "1.09999999",
2510<sub>15</sub>
                  "total snowfall (mm)": "0.42000002",
2511<sup>16</sup>
2512<sup>17</sup>
                  "precipitation hours (hours)": "3.0"
2513<sup>18</sup><sub>19</sub>
               }
            ]
2514<sub>20</sub>
         }
```

Parameters:

```
2518 1
         [
2519 2
           {
              "Input Name": "city_name",
2520<sup>3</sup>
              "Type": "String",
2521<sup>4</sup>
              "Description": "The name of the city."
     5
2522
     6
           },
2523 7
           {
2524 8
              "Input Name": "country_code",
              "Type": "Two Alphabet-Number",
2525 <sup>9</sup>
              "Description": "The country code (ISO 3166-2). The list can be found
2526^{10}
                  here https://en.wikipedia.org/wiki/ISO_3166-2"
2527<sub>11</sub>
           },
2528<sub>12</sub>
           {
2529<sub>13</sub>
              "Input Name": "start_date",
              "Type": "Date Format",
2530 14
              "Description": "The start date in YYYY-MM-DD format"
2531<sup>15</sup>
2532<sup>16</sup><sub>17</sub>
           },
           {
2533<sub>18</sub>
              "Input Name": "end_date",
2534 19
              "Type": "Date Format",
              "Description": "The start date in YYYY-MM-DD format"
2535 20
2536<sup>21</sup>
2537<sup>22</sup>
        1
```

2538 G.5 WIKI SEARCH 2539

Input Example:

2540 Description: Searches Wikipedia and returns a summary of the top pages matching the query.

2541 2542

2545

2546

2543 1 { "query": "covid-19", **2544**² "num_results": "1" 3 4 }

2547 **Output Example:**

```
2548
2549 1 {
255
255
```

2550 2	"error": "",
2551 ³	"result": [
2552 4	{
5	"title": "COVID-19",
2553 ₆	"summary": "Coronavirus disease 2019 (COVID-19) is a contagious
2554	disease caused by the coronavirus SARS-CoV-2. The first known
2555	case was identified in Wuhan, China, in December 2019. Most
2556	scientists believe the SARS-CoV-2 virus entered into human
2557	populations through natural zoonosis, similar to the SARS-CoV- $\!\!\!$
2558	and MERS-CoV outbreaks, and consistent with other pandemics in
2000	human history. Social and environmental factors including
2559	climate change, natural ecosystem destruction and wildlife
2560	trade increased the likelihood of such zoonotic spillover. The
2561	disease quickly spread worldwide, resulting in the COVID-19
2562	pandemic. The symptoms of COVID-19 are variable but often
2563	include fever, fatigue, cough, breathing difficulties, loss of
2303	smell, and loss of taste. Symptoms may begin one to fourteen
2564	days after exposure to the virus. At least a third of people
2565	who are infected do not develop noticeable symptoms. Of those
2566	who develop symptoms noticeable enough to be classified as
2567	patients, most (81%) develop mild to moderate symptoms (up to
2568 -	inita pheumonia), truncatea"
2569 。	
2569 8	

2572 **Parameters:**

}

[1

```
2573
2574
```

2575 $_{3}^{-}$

2576 4

2577 5 **2578**⁶

2579

2581₁₀

2582 11

2

7

8 2580 9

2570 9

2571

```
{
    "Input Name": "query",
    "Type": "String",
    "Description": "The search query."
  },
  {
    "Input Name": "num_results (Optional)",
    "Type": "Integer",
    "Description": "Number of search results to return."
  }
]
```

2583¹² 2584

2585 2586

2587

2588

```
2589
```

2590

2592 G.6 INTRADAY STOCK INFO 2593

2594 Description: Provides intraday time series data for specified equities.

"timestamp": "2024-07-16 19:00:00",

"open_market_value": "234.6520", "high_market_value": "234.7200",

"low_market_value": "234.2200",

"volume": "38722"

"volume": "24098"

"volume": "14364524"

"close_market_value": "234.3200",

"timestamp": "2024-07-16 18:00:00",

"open_market_value": "234.6220",

"high_market_value": "234.7500",

"close_market_value": "234.7000",

"timestamp": "2024-07-08 16:00:00",

"open_market_value": "227.8100",

"high_market_value": "227.8800", "low_market_value": "226.0630",

"close_market_value": "227.6400",

"low_market_value": "234.5050",

```
2595
        Input Example:
2596
```

```
2597 1
         {
            "symbol": "AAPL",
2598<sup>2</sup>
            "interval": "60min"
     3
2599
     4
         }
2600
```

Output Example:

{

{

. . .

{

}

Parameters:

]

}

"error": "",

"result": [

{

```
2601
2602
2603 1
2604<sup>2</sup>
2605
2606
2607 6
2608 7
2609<sup>8</sup>
2610 <sup>9</sup>
2611<sup>1</sup><sub>11</sub>
2612<sub>12</sub>
2613 13
2614<sup>14</sup>
2615<sup>15</sup>
2616<sub>17</sub>
2617<sub>18</sub>
2618 19
2619<sup>20</sup>
2620<sup>21</sup>
2621 <sub>23</sub>
2622<sub>24</sub>
2623 25
```

```
5
                              },
          16
                              },
2624<sup>26</sup>
2625<sup>27</sup>
```

2628 2629

28

2626⁻₂₉

2627₃₀

```
2630
         ſ
2631<sup>1</sup>
     2
           {
2632
              "Input Name": "symbol",
     3
2633 4
              "Type": "String",
              "Description": "The ticker symbol of the equity."
2634 5
2635 <sup>6</sup>
           },
           {
2636
              "Input Name": "interval",
2637 9
              "Type": "String",
2638<sub>10</sub>
              "Description": "Data point interval (1min, 5min, etc.)."
263911
           },
2640<sup>12</sup>
           {
              "Input Name": "month (optional)",
2641<sup>13</sup>
              "Type": "String",
    14
2642<sup>14</sup><sub>15</sub>
              "Description": "You can use the month parameter (in YYYY-MM format)
2643
                  to query a specific month in history."
2644 16
           }
        ]
2645<sup>17</sup>
```

"timestamp": "2024-07-16",

"timestamp": "2024-07-15", "open_market_value": "236.4800",

"timestamp": "2024-07-10",

"volume": "43234278"

"volume": "62631252"

"volume": "62627687"

"open_market_value": "235.0000",

"high_market_value": "236.2700",
"low_market_value": "232.3300",

"close_market_value": "234.8200",

"high_market_value": "237.2300",

"close_market_value": "234.4000",

"open_market_value": "229.3000",

"high_market_value": "233.0800",

"low_market_value": "229.2500",
"close_market_value": "232.9800",

"low_market_value": "233.0900",

2646 G.7 DAILY STOCK INFO

Input Example:

2648 **Description:** Returns daily time series data for specified equities.

2649 2650

2651 1 {
2652 2 "symbol": "AAPL",
2653 3 "number_of_days": 5
2654 }

Output Example:

{

},

{

},

{

}

]

}

. . .

"error": "",

"result": [

```
2655
2656
2657 <sub>1</sub>
2658 2
2659<sup>3</sup>
           4
2660
2661
           6
2662 7
2663 8
2664 <sup>9</sup>
2665<sup>10</sup>
          11
2666<sup>1</sup><sub>12</sub>
2667<sub>13</sub>
2668 14
2669<sup>15</sup>
2670<sup>16</sup><sub>17</sub>
2671<sub>18</sub>
2672<sub>19</sub>
2673 20
2674<sup>21</sup>
2675<sup>22</sup>
2676<sub>24</sub>
2677<sub>25</sub>
2678 26
2679<sup>27</sup>
2680<sup>28</sup>
          29
2681<sup>2</sup>/<sub>30</sub>
2682
2683
2684
```

```
Parameters:
         [
2685 1
2686<sup>2</sup>
            {
              "Input Name": "symbol",
2687
              "Type": "String",
     4
2688 5
              "Description": "The ticker symbol of the equity."
2689 6
           },
2690 7
           {
              "Input Name": "number_of_days",
2691<sup>8</sup>
2692<sup>9</sup><sub>10</sub>
              "Type": "Integer",
              "Description": "The number of days before today to return data for."
2693<sub>11</sub>
           }
2694_{12}
        ]
2695
2696
2697
2698
2699
```

2700 G.8 STOCK SYMBOL SEARCH

2702 **Description:** Searches for stock tickers based on provided keywords.

```
2703
2704 Input Example:
```

```
2705 1 {
2706 2 "keywords": "tesla"
2707 3 }
```

"error": "",

"symbol": "TSLA",

"name": "Tesla Inc",
"type": "Equity",

"region": "United States",

"market_open": "09:30",
"market_close": "16:00"

"match_score": "0.8889"

"market_open": "08:00",

"timezone": "UTC+02",

"symbol": "TL01.FRK",

"market_open": "08:00",

"timezone": "UTC+02",

"Input Name": "keywords",

the ticker symbol for"

"Type": "String",

"currency": "EUR",

"market_close": "20:00",

"match_score": "0.3846"

"type": "Equity",
"region": "Frankfurt",

"name": "TESLA INC. CDR DL-001",

"currency": "EUR",

"market_close": "20:00",

"match_score": "0.7143"

"timezone": "UTC-04",

"currency": "USD",

"symbol": "TLO.DEX", "name": "Tesla Inc",

"type": "Equity", "region": "XETRA",

"result": [

Output Example:

{

},

{

},

{

. . .

{

```
2709
2710
2711 1
2712<sup>2</sup>
2713<sup>3</sup>
             4
2714
2715<sub>6</sub>
2716 7
2717<sup>8</sup>
2718 <sup>9</sup>
2719<sup>10</sup><sub>11</sub>
2720<sub>12</sub>
2721 13
2722<sup>14</sup>
2723<sup>15</sup>
2724<sup>16</sup><sub>17</sub>
2725_{18}
2726 19
2727<sup>20</sup>
2728<sup>21</sup>
2729<sup>--</sup><sub>23</sub>
2730<sub>24</sub>
2731 25
2732<sup>26</sup>
2733<sup>27</sup>
2734<sup>28</sup><sub>29</sub>
2735<sub>30</sub>
2736 31
2737<sup>32</sup>
2738<sup>33</sup>
2739<sup>34</sup><sub>35</sub>
2740<sub>36</sub>
2741 37
2742<sup>38</sup>
2743<sup>39</sup>
2744
2745
```

2708

Parameters:

{

}

]

}

[

}

```
2746
2747 <sup>1</sup>
2748 <sup>2</sup>
3
2749 <sup>4</sup>
2750 <sup>5</sup>
2751 <sup>6</sup>
2753 <sup>7</sup>
```

"Description": "Keywords to search, , e.g., company name, to retrieve

2754 G.9 PYTHON 2755

{

2761 ³ }

}

Description: Runs a python interpreter on a code snippet.

2757 2758

2759 1 **2760** ²

"code": "print(4 + 5)"

2763 Output Example:

"result": "9", "error": ""

Input Example:

```
2764 1 {
2765 2
2766 3
```

2767 ⁴ 2768

2769

Parameters:

```
2770 1
        Γ
2771 <sub>2</sub>
           {
2772 3
             "Input Name": "code",
             "Type": "String",
2773 4
             "Description": "The code snippet that we want to run on a python
2774 <sup>5</sup>
                  interpreter."
2775
     6
           }
2776
    7
        ]
```

2777 2778 2779

2780

G.10 WOLFRAM ALPHA

2781 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.

2784 Input Example:

Output Example:

2785 2786

2

3 }

2793³₄}

{

2787 2788

2789

```
2790
```

2791¹ **2792**²

2794 2795

2796

```
"error": "",
"result": "47 years 5 months 13 days"
```

"query": "what is Ronaldo's age?"

Parameters:

```
Γ
2797<sup>1</sup>
2798<sup>2</sup>
             "Input Name": "query",
     3
2799
             "Type": "String",
     4
2800 5
             "Description": "The query to perform computations/searches on. When
2801
                 unsure of your query search, try searching yourself on the
                 website!"
2802
2803<sup>6</sup>
          }
        ]
2804
2805
2806
```

2808 G.11 GOOGLE SEARCH 2809

2810 **Description:** Performs a Google search and returns snippet results, without linked page details 2811 Google is often used for popular culture, location-awareness and crowdsourcing.

2812 **Input Example:** 2813

2817 2818

2819 2820 1

{

```
2814<sup>1</sup>
            "query": "What is the capital of France?",
2815<sup>2</sup>
            "location": "Paris"
     3
2816
     4
```

Output Example:

```
2822<sup>3</sup>
                        {
2823
         5
2824 6
2825 7
2826 8
2827
2828 9
2829<sub>10</sub>
2830 11
                        {
2831 12
2832<sup>13</sup>
2833<sup>17</sup><sub>15</sub>
        14
2834
2835
2836 16
2837<sup>17</sup>
2838<sup>18</sup><sub>19</sub>
                        {
2839<sub>20</sub>
2840 21
2841<sup>22</sup>
2842
2843<sup>23</sup>
2844
2845 24
2846<sup>25</sup>
                   1
2847^{26}
2848<sup>27</sup>
               }
2849
2850
2851
               Γ
2852 2
                    {
2853<sup>3</sup>
2854 <sup>4</sup>
         5
2855 6
                   },
2856 7
                   {
2857 8
2858 9
```

```
"error": "",
2821<sup>2</sup>
         "result": [
             "source": "en.wikipedia.org",
             "date": "None",
             "title": "Paris",
             "snippet": "Paris is the capital and largest city of France. With
                an official estimated population of 2,102,650 residents as of 1
                 January 2023 in an area of more than ...",
             "highlight": "Paris"
           },
             "source": "home.adelphi.edu",
            "date": "None",
             "title": "Paris facts: the capital of France in history",
             "snippet": "Paris facts: Paris, the capital of France. Paris is the
                 capital of France, the largest country of Europe with 550 000
                km2 (65 millions inhabitants).",
             "highlight": "Paris"
          },
           . . .
            "source": "britannica.com",
             "date": "None",
             "title": "France | History, Maps, Flag, Population, Cities, Capital
                , & ...",
             "snippet": "Get a special academic rate on Britannica Premium. The
                capital and by far the most important city of France is Paris,
                one of the world's preeminent cultural ... ",
             "highlight": "Paris"
          },
```

Parameters:

```
"Input Name": "query",
             "Type": "String",
             "Description": "The search query."
             "Input Name": "location (Optional)",
             "Type": "String",
             "Description": "The geographical location for the search (optional)."
2859<sup>10</sup>
          }
2860<sup>11</sup><sub>12</sub>
    11
        ]
```