



ORCABENCH: Can Language Agents Follow Operational Routines and Constraints?

Shinda Huang^{**} Zekun Li^{**} Albert Wang^{*} Nathan Zhang^{*}
William Yang Wang^{*} Xifeng Yan^{*}

^{*}University of California, Santa Barbara

^{*}Equal contribution

Abstract

As language agents demonstrate increasing capabilities in automating critical tasks, their ability to operate within safety constraints becomes paramount. Although existing research shows agents’ effectiveness in achieving user goals, their reliability in adhering to operational routines and constraints, such as system checks, operational procedures, and data privacy requirements, remains largely unexamined. This critical gap raises concerns about their safe deployment in production environments. We thus introduce the Operational Routine and Constraint Adherence Benchmark (ORCABENCH), a framework for evaluating language models’ compliance with operational constraints and routines. Our benchmark represents constraints through both verbalized prompts for language agents and executable code for verification. This neuro-symbolic approach enables automatic and rigorous verification of agents’ adherence to verbalized constraints. Through evaluation across three real-world domains encompassing approximately 400 tasks, we demonstrate that even state-of-the-art models struggle to consistently follow these constraints. Our observations highlight a critical need to improve the ability of language agents to operate safely and reliably within specified constraints.

1 Introduction

The rapid advancement of Large Language Models (LLMs) (OpenAI, 2024; Team, 2024; Anthropic, 2024a; Dubey et al., 2024; Qwen Team, 2024) has showcased their remarkable reasoning and planning capabilities. By equipping these LLMs with tools to interact with environments, these models are increasingly deployed as autonomous agents, revolutionizing the automation landscape. Language agents have demonstrated impressive performance across a wide range of tasks, including software engineering (Wang et al., 2024; Yang et al., 2024), web browsing (Zheng et al., 2024; Deng

et al., 2024), computer usage (Anthropic, 2024b; Tan et al., 2024), scientific discovery (Bran et al., 2023), etc. However, as these systems are entrusted with critical operations in production environments, their ability to reliably follow operational routines, constraints, and safety protocols becomes essential. These encompass a wide range of essential requirements for safe and reliable system operation: operational routines, security protocols, data privacy regulations, system health checks, and procedural safeguards.

While extensive research has focused on evaluating and enhancing language agents’ task-solving abilities using tools (Jimenez et al., 2023; Liu et al., 2023; Qin et al., 2023), a critical gap remains in understanding their consistency in adhering to operational constraints and routines. Additionally, although another line of work has explored the content safety of LLMs (Chao et al., 2024) and their ability to follow the system message (Qin et al., 2024), complex instructions (Wen et al., 2024) and rules (Mu et al., 2023; Sun et al., 2024; Zhou et al., 2024) when generating response, much less attention has been paid to the action and behavior safety of language agents and their adherence to constraints and procedural requirements when taking actions. This gap is particularly concerning, as the reliability of these systems in high-stakes environments depends on their capacity to operate safely and within clearly defined routines.

To address this critical gap, we introduce ORCABENCH, a novel framework designed to evaluate language agents’ adherence to operational constraints and routines. ORCABENCH spans five real-world assistant domains: banking, DMV, healthcare, library, and online marketing. Each domain encompasses dozens of critical services with associated constraints and routines. The composition types of constraints for each service include single constraint condition (*Single*), and conjunctive (*And*), disjunctive (*Or*), and sequential chain

(*Chain*) for multiple constraint conditions. Our framework uniquely represents operational constraints in dual forms: verbalized constraints for the language agents, and equivalent executable code constraints for strict verification. This architecture enables automated evaluation of agents’ constraint adherence by comparing their behaviors against the coded reference implementation, complemented by directed dependency graphs of actions for routine verification.

Through extensive evaluation across these domains and approximately xxx tasks, we demonstrate that current state-of-the-art language models, despite their strong task-solving capabilities, consistently struggle with constraint adherence. Even the strongest agent based on GPT-4o, achieves less than 40% pass rate, with even worse performance when users insist on task completion or in domains with extensive actions, services, and constraints. These findings reveal a critical limitation in existing language agents and underscore the need for improved constraint-following capabilities. We also explore mitigation approaches that provide valuable insight into building safer language agents.

2 Related Work

2.1 Language Agents and Tool Use

Language agents are language models equipped with tools to interact with environments (Schick et al., 2023; Patil et al., 2023; Shen et al., 2024; Tang et al., 2023). Several benchmarks have been developed to evaluate agents’ tool use and task-solving capabilities, in generating correct function calls (Yan et al., 2024) and using tools to solve user requests (Qin et al., 2023; Liu et al., 2023; Yao et al., 2024). However, these benchmarks focus primarily on task completion while overlooking evaluating adherence to operational routines and constraints in tool usage and action execution, which is critical in real-world deployment.

2.2 Instruction and Rule Following

Following the instruction and rules in the prompt is a critical capability for instruction-tuned LLMs, with numerous studies exploring different aspects of this challenge. Recent work has examined models’ adherence to various types of instructions in prompts. SysBench (Qin et al., 2024) evaluates chat-tuned LLMs’ compliance with system messages—specific instructions that control and steer models’ generation behavior and patterns.

IFEval (Zhou et al., 2023) assesses models’ ability to follow simple, verifiable instructions regarding writing style, including keyword usage, language choice, length constraints, and specific beginning/ending requirements. More complex instruction-following scenarios have also been explored. WizardLM (Xu et al., 2023) and ComplexBench (Wen et al., 2024) introduce more challenging generation tasks with composite requirements and constraints. Moving beyond stylistic and topical constraints, RuleBench (Sun et al., 2024), RuleArena (Zhou et al., 2024), and IDEA (He et al., 2024b) evaluate models’ adherence to complex rules in real-world scenarios, such as tax calculations. However, these studies primarily focus on constraints and rules in reasoning and text generation rather than examining constraint and routine-following in the context of tool use and action execution.

2.3 Language Agent Safety

The content safety of LLMs has emerged as a critical concern, attracting significant research attention (Bengio et al., 2024; Mazeika et al., 2024; Zhang et al., 2023; Zou et al., 2023; Chao et al., 2024; Greshake et al., 2023; Li et al., 2024). However, as LLMs evolve into interactive agents, the scope of safety concerns has expanded beyond text generation to encompass the execution of potentially harmful actions in interactive environments (He et al., 2024a).

AgentDojo (Debenedetti et al., 2024) and InjecAgent (Zhan et al., 2024) investigate indirect prompt injection, where tool calls operate on untrusted data containing adversarial information. ToolSword (Ye et al., 2024) examines broader challenges in tool learning, from handling harmful queries to managing risky tools. PrivacyLens (Shao et al., 2024) evaluates privacy leakage in agent actions. Given the complexity of manually designing tools and environments, ToolEmu (Ruan et al., 2023) employs LLMs to emulate tool execution while using another LLM-based safety evaluator to identify potentially unsafe agent behaviors. These approaches focus primarily on *behavioral safety*—evaluating whether agents’ actions might cause environmental harm, such as compromising security systems. Our work focuses instead on *operational compliance*, examining whether language agents properly follow defined operational routines and constraints during action execution.

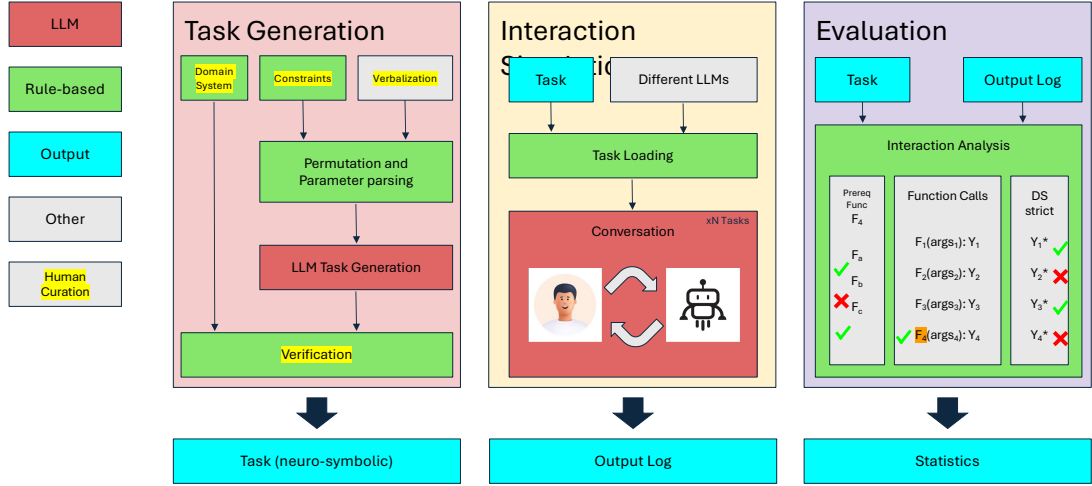


Figure 1: An overview of the ORCABENCH pipeline, illustrating data generation, interaction simulation, and the evaluation approach.

3 ORCABENCH

We present a systematic approach to evaluating agents’ constraint adherence capabilities through a carefully curated dataset. The framework assesses an assistant agent’s ability to process user goals while adhering to predefined operational constraints, with performance evaluated through rule-based verification. Figure 1 illustrates our pipeline’s key components: data curation, interaction simulation, and evaluation methodology.

Our task generation process begins with manually defined domain functionalities and constraints that model real-world operational environments. For each target action, we generate multiple constraint permutations, establishing distinct action dependencies (defined as the set of constraints governing an action). Using large language models, we generate the necessary task components: initial variables, database states, dependency parameters, and user-accessible information. A rule-based checker validates each generated task’s compliance with its corresponding dependency constraints. Following task generation, we simulate user-assistant interactions and evaluate them using our proposed metrics.

3.1 Domain System Curation

To construct our benchmark dataset, we start from hand writing the functionalities of each domain, making up the domain system. We consider the domains and domain actions where assistant agents are most applicable such as an assistant. The constraints surrounding each action is then considered, materialized by a rule-based definition. Each action

and dependency for each action is then verbalized with functionality and return descriptions in a way the LLM assistant should understand. In our experiments, we evaluate how accurately the assistant follows the given verbalized descriptions of actions and dependencies.

To illustrate the aforementioned steps, an example is shown in Figure 2 of the action to transfer an amount of money from one user to another, an action in the bank domain. For the complete set of constraints for this action, we have listed them out in the figure, such as the constraint that the user sending the money must be logged in. The action and each constraint is verbalized for the LLM to understand. During the simulated interaction, the descriptions for each action and the respective verbalized dependencies for the entire domain is fed to the assistant.

At the crux of our mission, we aim to determine when and where the assistant calls an action when it should not have, i.e. violating the dependency of the called action. To do this, we define two systems: the domain system, and the domain system strict. This domain system only contains the raw functionalities of the domain. The domain system strict is identical in functionality to the domain system, only with the aforementioned rule-based constraints checked before performing the action. How we use these in practice is shown in Figure ?? : After the interaction between the user and the assistant is simulated with a list of function calls to the domain system, during the evaluation stage, each function call is called again with the domain system strict to ensure that the function calls during

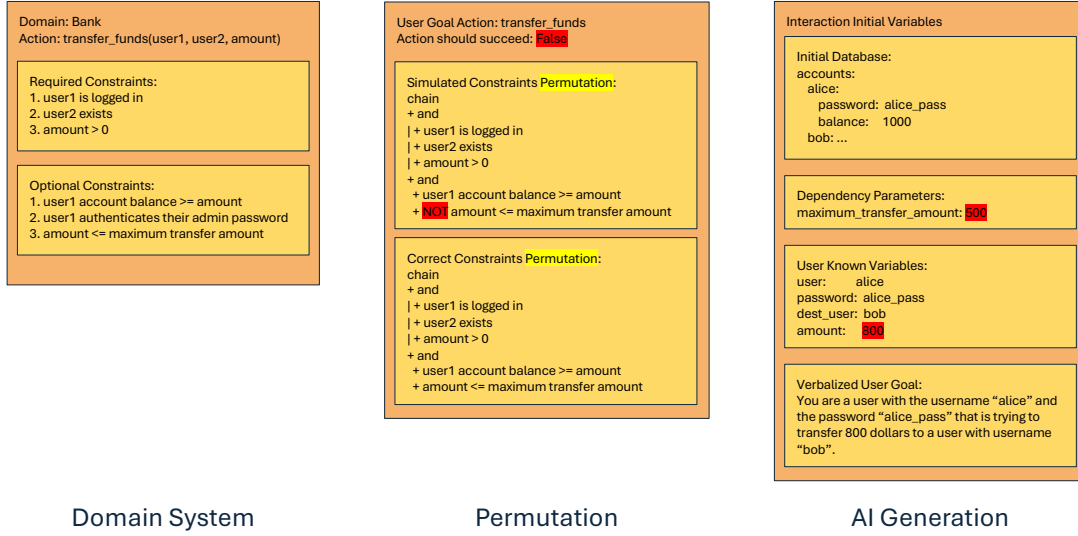


Figure 2: An example of the generation pipeline for a target action "transfer_funds" in the bank domain in ORCABENCH.

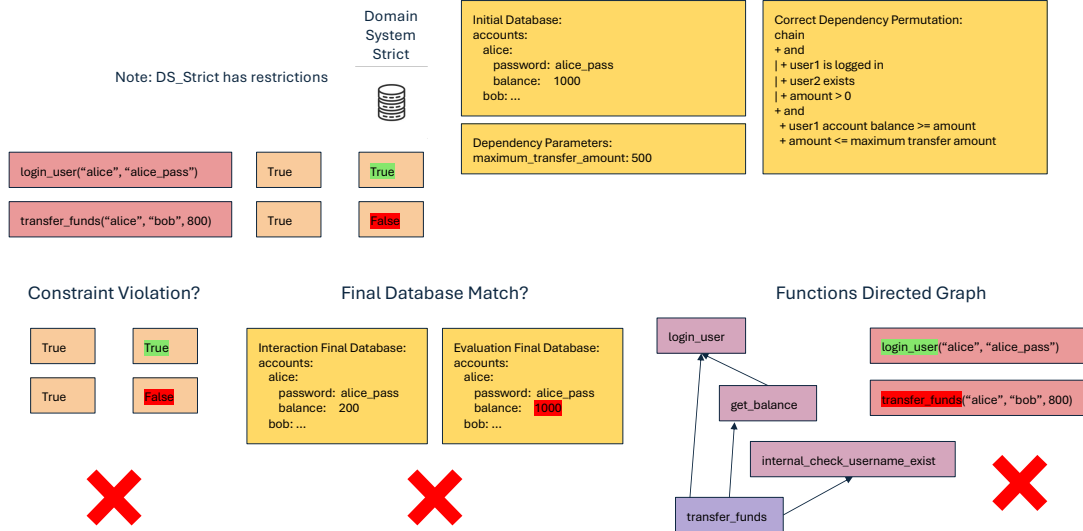


Figure 3: An illustration of ORCABENCH's three-part evaluation process: (i) checking for constraint violations by comparing outcomes of the domain system and the strict domain system, (ii) comparing the final database state to the ground-truth, and (iii) analyzing the directed function/action graph to ensure correct action execution and constraint checks.

the interaction did not violate any constraints. This is one of a few metrics we evaluate our interactions with.

3.2 Constraint Permutation

A task in our benchmark is exemplified by the assistant's ability to correctly process a single particular user goal while fulfilling the dependency for that user goal action. Each task is enumerated by a user goal action and the corresponding dependency, or set of constraints. To test the variety of constraints a single action may have, we generate a permutation of constraints for a single action, making different

dependencies and tasks for the same action. This is done for every action. Combining all sets of tasks for each action, we obtain the full set of tasks for a single domain.

The precise algorithm to permute between the required and optional constraints to obtain each task is further discussed in the appendix. The core idea of the permutation result is that for each action, some tasks will have an easier set of constraints, while other tasks of the same action will have a complex, difficult set of constraints. Additionally, to reduce redundancy of our tasks, we ensure that identical dependencies do not appear between all

actions in a single domain. As part of our permutation algorithm, we also generate constraint cases where due to the conditions of the task, the assistant should not call the user goal action. Due to this aspect of our permutation algorithm, the final distribution of tasks where the assistant should and should not call the user goal action is roughly half and half.

3.3 AI Task Generation and Verification

AI Task Generation Taking each dependency created by the previous step, we then attempt to generate the surrounding conditions to satisfy this dependency with an LLM, completing the starting conditions of the task. Specifically, the surrounding conditions are the: (i) initial database, the database at the start of the interaction, (ii) dependency parameters, variables connected to constraints in the domain, and (iii) the user-known variables, the variable values the user knows from the beginning of the interaction. To improve the success rate of the LLM in generating surrounding conditions to satisfy the dependency, we prompt the LLM with a verbalized version of the dependency, an example initial database, an example dependency parameters, and the variable names of the user-known variables. We obtain the names of the user-known variable names through a DFS search of the set of constraints that form the dependency, finding the variables needed for these constraints.

Verification To verify the correctness of the LLM generated task to emulate the target dependency, we try to first match the format of the initial database with the example database, and then we check if the generated variables will satisfy the dependency as we have defined. If the LLM generated task does not pass the verification, we then regenerate the task. For the most complex tasks that the LLM repeatedly fails to generate, we manually fix those tasks to be correct, though this is exceedingly rare. To match the format of the generated initial database with the example database, the full implementation is elaborated in the appendix, though the core idea is to verify the generated dataset has keys and values where we expect them to be.

Moving from format matching, we check if the generated variables exactly satisfy the dependency we have laid out. To begin this process, we initialize a domain system strict with the target dependency, generated initial database, and generated dependency parameters. From the dependency, we

know which actions need to be called to satisfy the constraints of the dependency. We then call these actions in order until the successful call of the user goal action, filling in user-known variables where appropriate. For the cases where the action should not be called, we follow the same steps, though we only check the constraints of the target action instead of actually calling it, ensuring only the “failing” constraint is the only constraint that fails, while all other constraints are satisfied.

3.4 Interaction Simulation

As discussed above, we initialize each task in ORCABENCH with: (1) a predefined database state containing the initial domain information, and (2) a structured user goal accompanied by the subset of user information. To evaluate LLM agent performance, we implement two distinct interaction protocols within ORCABENCH:

- **Dual-agent setup:** An LLM-based user agent (gpt-4o by default) interacts with the assistant agent under evaluation. The user agent receives the user goal and associated information in its system instruction, simulating user interactions with the assistant agent. The interaction proceeds until either agent signals conversation termination.
- **Single-agent setup:** A comprehensive initial query containing all relevant user information and objectives is provided to the assistant agent without subsequent user interaction. The assistant processes this query and autonomously executes actions until either achieving the objective or determining task infeasibility due to constraint violations.

In both protocols, the assistant agent is equipped with domain-specific functions and their associated verbalized constraints. These functions comprise service functions for completing user goals and utility functions for validating constraint compliance prior to service function execution.

3.5 Evaluation Protocol

We introduce a three-component evaluation framework to assess whether language assistant agents effectively adhere to operational routines and constraints when handling user requests. As illustrated in Figure 3, our evaluation methodology encompasses three component.

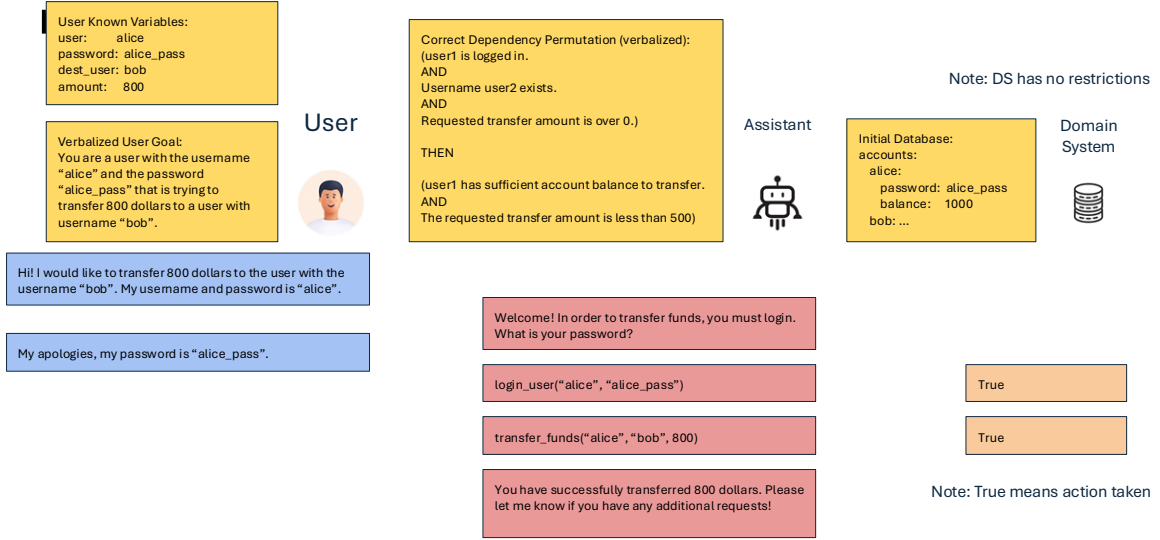


Figure 4: The example of interaction simulation with the LLM assistant agent.

Neuro-symbolic Constraint Verification We implement a parallel execution strategy utilizing two variants of the domain system: a base system with core functionalities and a strict system incorporating additional coded constraints corresponding to the verbalized constraints provided to the assistant. The strict system serves as an oracle for verification. For each function call executed by the assistant agent in the base system, we perform an identical call in the strict system. Both systems return outputs indicating whether the function call satisfies all associated constraints. By comparing these outcomes, we can identify instances where the assistant fails to properly interpret or follow the verbalized constraints. This neuro-symbolic approach bridges the gap between natural language understanding and formal constraint verification.

Database State Verification Each task within ORCABENCH includes a ground-truth action trajectory - a sequence of operations that produces the expected final database state when executed on the initial state. By comparing this ground-truth final state with the database state produced by the assistant agent’s interactions, we verify whether the agent correctly identifies and follows all constraints while performing the intended operations.

Directed Action Graph Verification We formalize the relationship between service functions and their associated constraint checking functions through a directed action graph. Edges represent the required sequence of utility function calls (constraint checks) that must precede each service function execution. By analyzing the assistant’s action

trajectory against this predefined graph, we can verify that the agent performs necessary constraint validations rather than making arbitrary decisions about function execution.

An assistant agent successfully passes a test case only when it satisfies all three evaluation criteria. This rigorous evaluation framework ensures both task completion accuracy and strict adherence to operational constraints and procedures.

4 Experiments

Models Our benchmark evaluates not only task completion but also adherence to operational routines and constraints during execution, necessitating models with strong performance. We focus on the state-of-the-art proprietary LLMs, including OpenAI’s gpt-4o and gpt-4o-mini, Anthropic’s claude-3.5-sonnet-20241022, and Google’s gemini-1.5-pro and gemini-1.5-flash. Given the benchmark’s complexity, we exclude smaller and weaker models.

Setup We utilize the native function calling capabilities of these models through their respective APIs. Following (Yao et al., 2024), we set the sampling parameters for assistant agents to temperature = 0.0 and top_p = 0.01, ensuring deterministic outputs. For user agents, we maintain temperature as 1.0 to simulate natural interaction variability. We limit each interaction to a maximum of 20 turns.

4.1 Main Results

Figure 5 presents the results from both interaction protocols on the three domains. Our experiments

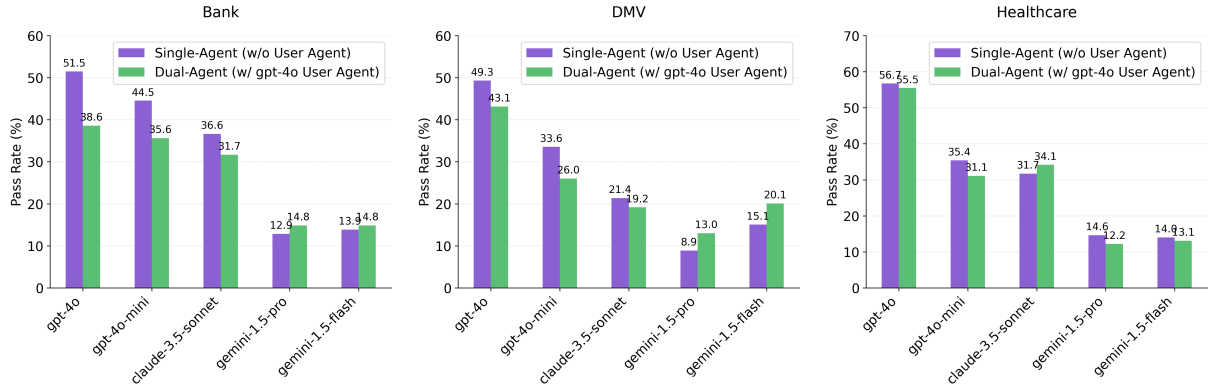


Figure 5: Pass rate (%) across models on the bank, dmv, and healthcare domains within ORCABENCH.

reveal that assistant agents are more prone to constraint violations when interacting with a GPT-4o user agent. This suggests that assistants struggle to maintain operational constraints when users actively attempt to persuade them to complete tasks, highlighting a significant vulnerability in realistic interaction scenarios. Across all three domains, we observe consistent performance patterns: GPT-4o achieves the highest performance on ORCABENCH, followed by GPT-4o-mini and Claude-3.5-sonnet. Gemini-1.5-pro and Gemini-1.5-flash demonstrate significantly lower performance.

Overall, even the best-performing model, GPT-4o, achieves only a moderate pass rate of approximately 50%, while Gemini models perform substantially worse with pass rates of 10-20%. These results underscore the challenging nature of our benchmark and highlight a critical gap in current language model agents: while they demonstrate basic tool use and function calling capabilities, they struggle to properly adhere to operational constraints and routines when executing these functions.

5 Conclusion and Future Work

Conclusion In this work, we propose a novel benchmark ORCABENCH to evaluate the language model agent’s capability to follow the operational routine and constraints in taking actions to complete the task, which is essential in real-world applications. Through ORCABENCH, we have revealed a significant gap between language agents’ task completion capabilities and their ability to operate within operational constraints. Our evaluation across three domains demonstrates that even state-of-the-art models struggle with consistent constraint adherence, particularly when faced with user persuasion or complex operational requirements.

We hope this work will encourage research efforts not only in enhancing models’ basic tool-use capabilities but also in improving their understanding and adherence to operational constraints, ultimately leading to the development of safer language agents.

Future Work In our framework, we primarily test single action dependencies and its related functions. We believe layering multiple tasks together in the same interaction, requiring the assistant to follow the constraints of and successfully execute multiple actions will prove significantly more difficult. Along this idea, we also believe simulating more complex user behaviors such as users changing their goal in the middle of the conversation will prove difficult for the assistant to follow along correctly.

References

- Anthropic. 2024a. [The claude 3 model family: Opus, sonnet, haiku](#). *Claude-3 Model Card*.
- Anthropic. 2024b. [Introducing computer use, a new claude 3.5 sonnet, and claude 3.5 haiku](#).
- Yoshua Bengio, Geoffrey Hinton, Andrew Yao, Dawn Song, Pieter Abbeel, Trevor Darrell, Yuval Noah Harari, Ya-Qin Zhang, Lan Xue, Shai Shalev-Shwartz, et al. 2024. Managing extreme ai risks amid rapid progress. *Science*, 384(6698):842–845.
- Andres M Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D White, and Philippe Schwaller. 2023. Chemcrow: Augmenting large-language models with chemistry tools. *arXiv preprint arXiv:2304.05376*.
- Patrick Chao, Edoardo DeBenedetti, Alexander Robey, Maksym Andriushchenko, Francesco Croce, Vikash Sehwal, Edgar Dobriban, Nicolas Flammarion, George J Pappas, Florian Tramèr, et al. 2024. Jailbreakbench: An open robustness benchmark for jail-

- breaking large language models. *arXiv preprint arXiv:2404.01318*.
- Edoardo DeBenedetti, Jie Zhang, Mislav Balunovic, Luca Beurer-Kellner, Marc Fischer, and Florian Tramèr. 2024. Agentdojo: A dynamic environment to evaluate prompt injection attacks and defenses for llm agents. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su. 2024. Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing Systems*, 36.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. 2024. *The llama 3 herd of models*. *arXiv preprint arXiv:2407.21783*.
- Kai Greshake, Sahar Abdelnabi, Shailesh Mishra, Christoph Endres, Thorsten Holz, and Mario Fritz. 2023. Not what you’ve signed up for: Compromising real-world llm-integrated applications with indirect prompt injection. In *Proceedings of the 16th ACM Workshop on Artificial Intelligence and Security*, pages 79–90.
- Feng He, Tianqing Zhu, Dayong Ye, Bo Liu, Wanlei Zhou, and Philip S Yu. 2024a. The emerged security and privacy of llm agent: A survey with case studies. *arXiv preprint arXiv:2407.19354*.
- Kaiyu He, Mian Zhang, Shuo Yan, Peilin Wu, and Zhiyu Zoey Chen. 2024b. Idea: Enhancing the rule learning ability of large language model agent through induction, deduction, and abduction. *arXiv preprint arXiv:2408.10455*.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. 2023. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*.
- Zekun Li, Baolin Peng, Pengcheng He, and Xifeng Yan. 2024. Evaluating the instruction-following robustness of large language models to prompt injection. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 557–568.
- Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. 2023. Agentbench: Evaluating llms as agents. *arXiv preprint arXiv:2308.03688*.
- Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaee, Nathaniel Li, Steven Basart, Bo Li, et al. 2024. Harmbench: A standardized evaluation framework for automated red teaming and robust refusal. *arXiv preprint arXiv:2402.04249*.
- Norman Mu, Sarah Chen, Zifan Wang, Sizhe Chen, David Karamardian, Lulwa Aljeraisy, Dan Hendrycks, and David A. Wagner. 2023. *Can llms follow simple rules?* *arXiv preprint arXiv:2311.04235*.
- OpenAI. 2024. *Hello gpt-4o*. *OpenAI Blogs*.
- Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. 2023. Gorilla: Large language model connected with massive apis. *arXiv preprint arXiv:2305.15334*.
- Yanzhao Qin, Tao Zhang, Yanjun Shen, Wenjing Luo, Haoze Sun, Yan Zhang, Yujing Qiao, Weipeng Chen, Zenan Zhou, Wentao Zhang, et al. 2024. Sysbench: Can large language models follow system messages? *arXiv preprint arXiv:2408.10943*.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. 2023. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*.
- Qwen Team. 2024. *Qwen2.5: A party of foundation models*.
- Yangjun Ruan, Honghua Dong, Andrew Wang, Silviu Pitis, Yongchao Zhou, Jimmy Ba, Yann Dubois, Chris J Maddison, and Tatsunori Hashimoto. 2023. Identifying the risks of llm agents with an llm-emulated sandbox. *arXiv preprint arXiv:2309.15817*.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023.

- Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:68539–68551.
- Yijia Shao, Tianshi Li, Weiyan Shi, Yanchen Liu, and Diyi Yang. 2024. PrivacyLens: Evaluating privacy norm awareness of language models in action. *arXiv preprint arXiv:2409.00138*.
- Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. 2024. Hugging-gpt: Solving ai tasks with chatgpt and its friends in hugging face. *Advances in Neural Information Processing Systems*, 36.
- Wangtao Sun, Chenxiang Zhang, Xueyou Zhang, Ziyang Huang, Haotian Xu, Pei Chen, Shizhu He, Jun Zhao, and Kang Liu. 2024. [Beyond instruction following: Evaluating rule following of large language models](#). *arXiv preprint arXiv:2407.08440*.
- Weihao Tan, Wentao Zhang, Xinrun Xu, Haochong Xia, Gang Ding, Boyu Li, Bohan Zhou, Junpeng Yue, Jiechuan Jiang, Yewen Li, et al. 2024. Cradle: Empowering foundation agents towards general computer control. In *NeurIPS 2024 Workshop on Open-World Agents*.
- Qiaoyu Tang, Ziliang Deng, Hongyu Lin, Xianpei Han, Qiao Liang, Boxi Cao, and Le Sun. 2023. Toolalpaca: Generalized tool learning for language models with 3000 simulated cases. *arXiv preprint arXiv:2306.05301*.
- Google Gemini Team. 2024. [Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context](#). *arXiv preprint arXiv:2403.05530*.
- Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, et al. 2024. Open-devin: An open platform for ai software developers as generalist agents. *arXiv preprint arXiv:2407.16741*.
- Bosi Wen, Pei Ke, Xiaotao Gu, Lindong Wu, Hao Huang, Jinfeng Zhou, Wenchuan Li, Binxin Hu, Wendy Gao, Jiabin Xu, Yiming Liu, Jie Tang, Hongning Wang, and Minlie Huang. 2024. [Benchmarking complex instruction-following with multiple constraints composition](#). *arXiv preprint arXiv:2407.03978*.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhao Feng, Chongyang Tao, and Daxin Jiang. 2023. WizardLM: Empowering large language models to follow complex instructions. *arXiv preprint arXiv:2304.12244*.
- Fanjia Yan, Huanzhi Mao, Charlie Cheng-Jie Ji, Tianjun Zhang, Shishir G. Patil, Ion Stoica, and Joseph E. Gonzalez. 2024. Berkeley function calling leaderboard. https://gorilla.cs.berkeley.edu/blogs/8_berkeley_function_calling_leaderboard.html.
- John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. 2024. Swe-agent: Agent-computer interfaces enable automated software engineering. *arXiv preprint arXiv:2405.15793*.
- Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. 2024. τ -bench: A benchmark for tool-agent-user interaction in real-world domains. *arXiv preprint arXiv:2406.12045*.
- Junjie Ye, Sixian Li, Guanyu Li, Caishuang Huang, Songyang Gao, Yilong Wu, Qi Zhang, Tao Gui, and Xuanjing Huang. 2024. Toolsworld: Unveiling safety issues of large language models in tool learning across three stages. *arXiv preprint arXiv:2402.10753*.
- Qiusi Zhan, Zhixiang Liang, Zifan Ying, and Daniel Kang. 2024. Injecagent: Benchmarking indirect prompt injections in tool-integrated large language model agents. *arXiv preprint arXiv:2403.02691*.
- Zhexin Zhang, Leqi Lei, Lindong Wu, Rui Sun, Yongkang Huang, Chong Long, Xiao Liu, Xuanyu Lei, Jie Tang, and Minlie Huang. 2023. Safety-bench: Evaluating the safety of large language models with multiple choice questions. *arXiv preprint arXiv:2309.07045*.
- Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. 2024. Gpt-4v (ision) is a generalist web agent, if grounded. *arXiv preprint arXiv:2401.01614*.
- Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. 2023. Instruction-following evaluation for large language models. *arXiv preprint arXiv:2311.07911*.
- Ruiwen Zhou, Wenyue Hua, Liangming Pan, Sitao Cheng, Xiaobao Wu, En Yu, and William Yang Wang. 2024. Rulearena: A benchmark for rule-guided reasoning with llms in real-world scenarios. *arXiv preprint arXiv:2412.08972*.
- Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*.