

# ECom-Bench: Can LLM Agent Resolve Real-World E-commerce Customer Support Issues?

**Haoxin Wang<sup>1,2</sup>, Xianhan Peng<sup>1,2</sup>, Xucheng Huang<sup>1</sup>,**  
Yizhe Huang<sup>1</sup>, Ming Gong<sup>1</sup>, Chenghan Yang<sup>1,2</sup>, Yang Liu<sup>1</sup>, Ling Jiang<sup>1</sup>  
<sup>1</sup>Xiaoduo AI Lab, <sup>2</sup>Shanghai Jiao Tong University  
{wanghaoxin, williampen, scottyang}@sjtu.edu.cn  
{huangxucheng, gongming, liuyangfoam, huangyizhe}@xiaoduotech.com

## Abstract

In this paper, we introduce ECom-Bench, the first benchmark framework for evaluating LLM agent with multimodal capabilities in the e-commerce customer support domain. ECom-Bench features dynamic user simulation based on persona information collected from real e-commerce customer interactions and a realistic task dataset derived from authentic e-commerce dialogues. These tasks, covering a wide range of business scenarios, are designed to reflect real-world complexities, making ECom-Bench highly challenging. For instance, even advanced models like GPT-4o achieve only a 10–20% pass@3 metric in our benchmark, highlighting the substantial difficulties posed by complex e-commerce scenarios. The code and data have been made publicly available at <https://github.com/XiaoduoAILab/ECom-Bench> to facilitate further research and development in this domain.

## 1 Introduction

E-commerce has become a cornerstone of the global economy, with online transactions accounting for a significant share of retail sales. However, as the scale of e-commerce continues to expand, the demand for efficient customer service has surged. Traditional rule-based intelligent customer service systems, which have alleviated some of the pressure on human customer service, are now encountering performance bottlenecks. These systems typically handle only single-turn interactions, lack

the ability to understand context, cannot call external tools, and are unable to process multimodal information such as images or videos, which significantly limits their application in business (Zhang et al., 2021).

Recently, intelligent customer service solutions based on large language models (LLMs) have emerged as a research hotspot (Pandya and Holia, 2023). The strong language generation and context understanding capabilities of LLMs provide significant advantages in handling complex interactive tasks. However, the form and architecture of this new paradigm remain unclear, and deploying LLM-based solutions directly poses considerable risks. To date, current evaluation frameworks fall short in comprehensively assessing LLM-based customer service agent in e-commerce scenarios. For example, some frameworks lack e-commerce-specific tasks, such as ToolBench (Xu et al., 2023); others support only single-turn dialogues, like WebShop (Yao et al., 2022); and some omit multimodal task evaluations, as seen in  $\tau$ -Bench (Yao et al., 2024). These limitations hinder a full assessment of LLM-based agent in complex multimodal scenarios. Table 1 provides a detailed comparison of ECom-Bench with other related frameworks.

To address these issues, we propose ECom-Bench, the first benchmark framework for evaluating multimodal LLM agent in e-commerce customer support scenarios. ECom-Bench aims to fill the gap in existing research by combining persona-

Feature	ECom-Bench	$\tau$ -Bench	WebShop	ToolBench	EcomScriptBench
User Simulation	✓	✓	×	×	×
Persona Driven	✓	×	×	×	×
Multi Turn	✓	✓	×	×	×
Multimodal	✓	×	×	×	×
Tools	✓	✓	×	✓	×

Table 1: Comparison of ECom-Bench with other simulation environments.

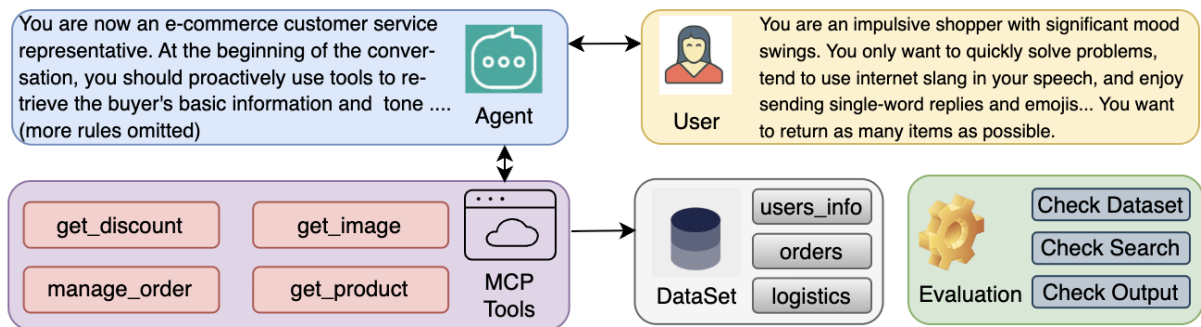


Figure 1: Framework of ECom-Bench

driven user simulation based on real e-commerce customer interactions and a realistic task dataset derived from authentic e-commerce dialogues to provide a rigorous and comprehensive evaluation platform. This framework not only enhances the accuracy and authenticity of user simulation but also evaluates agent capabilities across a wide range of business scenarios.

The main contributions of this work are:

1. **ECom-Bench Framework:** We introduce ECom-Bench, the first open-source benchmark framework for evaluating multimodal LLM agent with Model Context Protocol(MCP) tools in e-commerce customer support scenarios.
2. **Persona-driven User Simulation:** We present a dataset of persona information meticulously collected from real e-commerce customer interactions. This dataset enables the simulation of user behaviors with LLMs, capturing diverse customer personalities and interaction patterns to enhance the accuracy and authenticity of user simulation in e-commerce contexts.
3. **Realistic Task Dataset:** We develop a dataset of task instances derived from authentic e-commerce dialogues. This dataset encompasses a wide range of business scenarios, including product inquiries, order management, and multimodal content recognition. The dataset consists of 53 manually verified tasks, with 18 involving multimodal interactions, ensuring representativeness and alignment with real-world business requirements.

## 2 Related Work

Recently, research focus has progressively shifted from dataset performance metrics to evaluating agent capabilities in specific environments. However, relevant research exhibits several limitations: some focus on internal operations within simulation environments while neglecting tool usage capabilities (Shridhar et al., 2021)(Zhou et al., 2023)(Yang et al., 2023)(Wang et al., 2025); others simulate real interaction scenarios but only address single-turn interactions(Yao et al., 2022)(Budzianowski et al., 2018)(Soltau et al., 2023); and some provide extensive tools but lack domain-specific data(Qin et al., 2023).

The significant discrepancy between LLMs and real user behaviors severely constrains the development of user simulation. Although some studies have attempted to solve this problem by providing LLMs with memory, emotions, and personality traits, they are either limited to multi-agent interactions in sandbox environments (Wang et al., 2023) or focus exclusively on user preference issues (Xiang et al., 2024), rendering them inadequate for intelligent customer service application. Chan et al. (Chan et al., 2024) validated the potential of persona-driven synthetic data for training and testing, but their personas differ significantly from e-commerce customer characteristics, which is challenging for LLMs to accurately simulate e-commerce customer behavior.

Evaluation in vertical domains is crucial for LLM deployment. Consequently, numerous domain-specific datasets have emerged, spanning medical (Pal et al., 2022)(Zhou et al., 2023), legal (Colombo et al., 2024)(Guha et al., 2023), and telecommunications (Lee et al., 2024) sectors. All these customized datasets emphasize the importance of domain-specific evaluation. While  $\tau$ -Bench offers a multi-domain dynamic interaction

benchmark that includes retail, its test scenarios lack depth, and its user simulations remain relatively simplistic.

### 3 Benchmark Construction

ECom-Bench comprises user simulation classes, database JSON files, tools, task instances, and domain-specific documents. All data are based on historical online dialogues, primarily sourced from conversations related to home and appliance products. We first use LLMs to generate preliminary content, which is then reviewed by humans to ensure consistency and authenticity.

**User personas.** Grounded in established principles from consumer psychology and behavioral studies, and leveraging user modeling methodologies from recommender systems, we meticulously designed our user personas. These personas encompass consumer types, personality traits (such as emotion, attention to detail, patience, trust, and rights awareness), and behavioral traits (including questioning style, communication style, and interaction patterns). This comprehensive design enables LLMs to interact with customer service agent in a manner closely resembling that of real human customers. By inputting online dialogues into LLMs and employing chain-of-thought reasoning (Wei et al., 2022), we systematically analyzed the data across these dimensions, generating hundreds of high-quality user personas.

**Database.** To eliminate ambiguity in usage, we define and annotate essential data classes, including products, logistics, orders, invoices, and more. We then prompt LLMs with the company’s internal data and data classes to generate comprehensive and accurate datasets. These datasets are subsequently reviewed and processed by human experts to ensure reliability and validity.

**Tools.** Drawing on existing business tools, insights from experienced customer service staff, and strategic business planning, we have meticulously defined 21 categories of tools. These tools encompass a wide range of functionalities, including product inquiry, order modification, returns and exchanges, and more. Some tools are composite, supporting multiple operations, while others require a strict operational sequence to ensure proper execution. To better align with real-world scenarios, we have introduced a multimodal tool that enables agent to "see" images. All tools are implemented in strict compliance with the MCP.

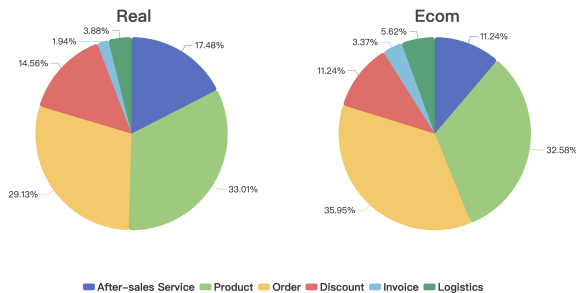


Figure 2: Distribution of task instances

**Task instances.** To ensure that simulated user trajectories closely mirror those of real customers, we enlisted trained customer service personnel to analyze authentic online conversations and extract the core problems that customers aimed to resolve. These extracted intents are then transformed into standardized task instances for LLMs to learn from. Our task instances are designed to align with real business scenarios in terms of both category and distribution, ensuring representativeness. To further evaluate agent capabilities, we also manually introduced certain variations into the tasks (see Figure 2).

**Domain documentation.** We provide the agent with operation documents tailored for e-commerce scenarios and prompt the simulated user with a detailed customer guideline as a partial description of the world model.

**Evaluation.** We evaluate each task across three dimensions: (1) comparing the database state against the ground truth expected state; (2) verifying whether the model invoked the required information retrieval tools; (3) examining whether the output contains the necessary keywords. This rule-based reward can objectively reflect the agent’s decision-making capabilities and hallucination tendencies, while also allowing for stochastic variation in dialogues: user may express themselves in different way, yet the database state and key search steps remain consistent after the customer service agent completes the same task. Given that customer service agent must handle identical issues reliably and consistently, we introduce the  $\text{pass}^k$  (Yao et al., 2024) to measure agent robustness across  $k$  independent and identically distributed trials for the same problem. Let  $n$  denote the total number of trials for a task,  $c$  denote the number of successful trials, and  $k$  denote the number of trials under

evaluation. Then we have:

$$\text{pass}^k = \mathbb{E}_{\text{task}} \left[ \frac{\binom{c}{k}}{\binom{n}{k}} \right]$$

## 4 Experiments

This section presents the performance evaluation of various agents using the ECom-Bench framework. For each trial  $j$  (where  $j = 1, 2, \dots, k$ ), we conduct the  $j$ -th execution of the evaluation. The score for trial  $j$  means the chance that all  $j$  i.i.d. task trials are successful, averaged across tasks.

### 4.1 Setup

We evaluate various proprietary models, LLMs like GPT-4o, DeepSeek-V3, Doubao-1.5-Pro-32k, Qwen-Max, and Moonshot-V1-32k, Multimodal Large Language Models (MLLMs) like Qwen-VL-Max, Doubao-Pro-Vision. Agents are built using LangGraph, which incorporates the ReAct, limited to 20 turns and 600 seconds per task (temperature=0.3). Qwen-Max simulates users, while Moonshot-V1-128k-Vision-Preview processes images as a tool.

### 4.2 Main Results

#### 4.2.1 Overall Model Performance

Table 2 shows significant differences in the performance of various models on e-commerce customer service tasks. Such tasks require models to strictly follow domain-specific rules and accurately invoke the appropriate tools. Some models, such as Qwen-Max, demonstrate more natural language expression when simulating user roles, but are more prone to hallucinations during operational execution, which negatively impacts their overall scores. In contrast, models with stronger instruction-following capabilities and more controllable behaviors (such as GPT-4o and DeepSeek-

Model	Pass <sup>1</sup>	Pass <sup>2</sup>	Pass <sup>3</sup>
<b>GPT-4o</b>	44.03	26.42	16.98
<b>Doubao-1.5-Pro-32k</b>	38.99	23.90	16.98
<b>DeepSeek-V3</b>	36.48	21.38	15.09
<b>Qwen-Max</b>	15.09	8.81	5.66
<b>Moonshot-V1-32k</b>	14.47	7.55	5.66
<b>Doubao-Pro-Vision*</b>	20.75	13.21	11.32
<b>Qwen-VL-Max*</b>	8.18	4.4	3.77

Table 2: Model performance comparison with pass<sup>3</sup> (\* denotes MLLMs)

V3) exhibit more stable task completion rates. GPT-4o achieves the highest pass<sup>1</sup> and pass<sup>3</sup> scores, demonstrating that it has higher success rates across all tasks and exhibits more consistent performance in solving the same task repeatedly. This superior performance indicates that GPT-4o not only handles diverse tasks with greater ease, but also maintains more stable and reliable problem-solving capabilities when faced with identical challenges multiple times.

#### 4.2.2 Architecture Comparison

The decoupled MLLM-as-Tool approach (GPT-4o: pass<sup>3</sup>=16.98) outperforms end-to-end MLLM-as-Planner (Qwen-VL-Max: pass<sup>3</sup>=3.77, Doubao-Pro-Vision: pass<sup>3</sup>=11.32), revealing challenges in integrated multimodal planning.

#### 4.2.3 Consistency and Robustness

In business scenarios, it is vital for agent to reliably solve the same problems consistently. We run all tasks eight times. As shown in Figure 4, GPT-4o achieves the best performance with an average success rate of 43.87% across all tasks. However, as illustrated in Figure 3, the pass<sup>8</sup> drops to 7.55% (even worse than Doubao-1.5-Pro-32K), indicating that as the number of trials  $K$  increases, the probability of reliably and consistently solving the same task multiple times decreases significantly. This demonstrates that current models still have substantial room for improvement in terms of consistency and robustness.

#### 4.2.4 Error Analysis

We analyzed 159 trajectories (53 tasks executed 3 times) generated by the Doubao-1.5-Pro-32k agent and categorized the errors into three main types:

- Wrong Argument:** The agent uses the correct tool but fills in some arguments incor-

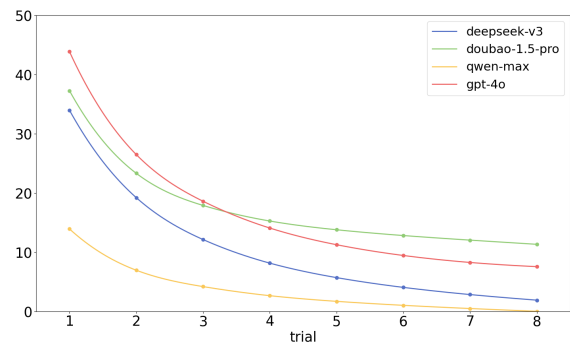


Figure 3: pass<sup>8</sup> across models

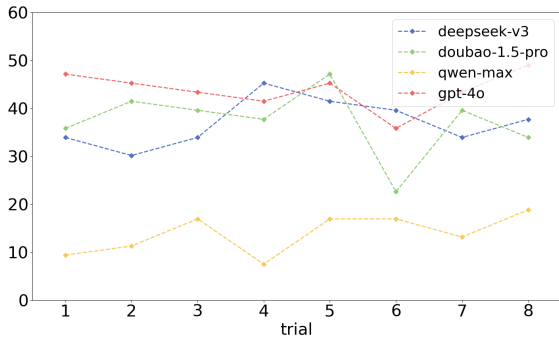


Figure 4: pass<sup>1</sup> across models

rectly (e.g., calls the order query tool with the correct order ID but fills in other irrelevant parameters).

2. **Wrong Decision:** The agent fails to choose the appropriate tool to solve the problem or processes the wrong object when invoking tools (e.g., calls the order query tool but queries an incorrect order ID; or calls wrong tools).
3. **Partially Resolved:** The agent only solves part of the problem (e.g., when handling return and refund requests, it only processes the return but fails to execute the refund; or proposes actions to be taken but does not actually invoke the tools to perform the operations).

As shown in Figure 5, Doubao-1.5-Pro is more prone to wrong decision errors, indicating that the current agent still faces difficulties in understanding user intents and tool usage. Additionally, when

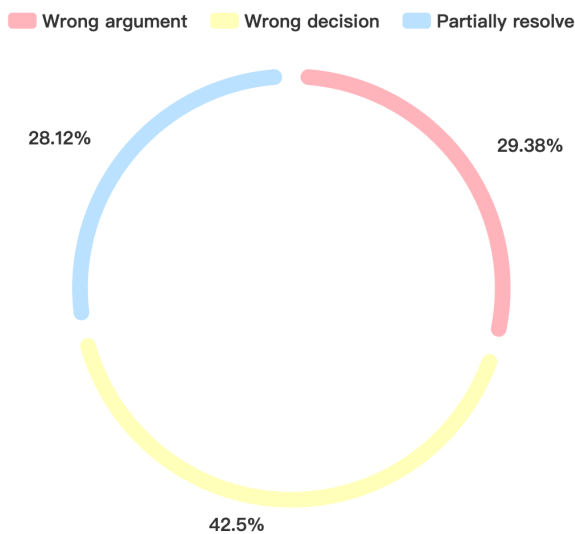


Figure 5: Breakdown trajectories of Doubao-1.5-Pro

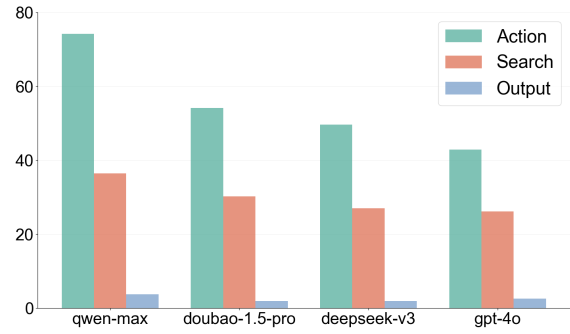


Figure 6: Failure rates in Action, Search, and Output dimensions across evaluated models

multiple entities are present in the context window, the agent tends to overlook differences between entities and becomes confused during processing.

#### 4.2.5 Dimensional Evaluation

We primarily evaluate agent from three aspects: action, search, and output. The search focuses on the execution status of search-related tools, mainly reflecting whether agent exhibits hallucination phenomena. The action examines database operation results, mainly indicating whether agent can correctly and completely execute all required steps. The output concentrates on the agent's output keywords, mainly reflecting whether agent can properly invoke multimodal tools to recognize image content and generate key words. As shown in Figure 6, all agents demonstrate higher failure rates in the Action dimension, which can be attributed to the inherent difficulty of database modification tasks. Meanwhile, models with superior overall performance (e.g., GPT-4o) exhibit lower error rates across all three dimensions, showing a greater tendency to obtain information through invoking external tools rather than directly generating answers, as well as stronger tool usage awareness and factual consistency. At the same time, compared to weaker models (e.g., Qwen-Max), the performance gap between action and search is also smaller. This reflects the higher stability of advanced models in handling different tasks and stronger robustness when dealing with complex operations. This result indicates that model performance improvements are often more pronounced in challenging domains such as database operations, where the gap between strong and weak models becomes larger.



Group	Emotion	Patience	Style	Pass^1	Pass^2	Pass^3
1	Calm	High	Clear	58.3	40.0	30.0
2	Angry	High	Clear	55.0	38.3	30.0
3	Calm	Low	Clear	48.1	31.0	20.0
4	Calm	High	Vague	45.0	29.9	20.0
5	Angry	Low	Clear	46.6	29.9	25.0
6	Angry	High	Vague	45.0	29.9	30.0
7	Calm	Low	Vague	45.0	26.7	20.0
8	Angry	Low	Vague	40.0	20.0	10.0

Table 3: Ablation study results on persona-driven user simulation.

## 5 Ablation Study on Persona-driven User Simulation

To investigate the impact of different persona traits on user model performance, we conduct an ablation study by manipulating three key dimensions of user personas: **emotional state** (calm vs. angry), **patience level** (high vs. low), and **communication style** (clear and professional vs. vague and colloquial).

### 5.1 Experiment Setup

Eight different combinations of the above dimensions are tested. Qwen-Max simulates users and GPT-4o serves as agent. Temperature of user model is set to 0.3 (default value). We pick 20 tasks significantly influenced by persona settings to conduct experiments. The performance is measured by pass^3.

### 5.2 Result Analysis

The results, summarized in Table 3, indicate that the extreme characteristics of personas markedly impact agent performance:

- **Emotional State:** Angry emotional characteristics can cause agent to generate hallucinated responses rather than invoke tool-based responses when under pressure.
- **Patience Level:** Low patience levels prevent users from completing multi-turn interactions necessary to provide agent with required information, thereby impacting task completion.
- **Communication Style:** Users unfamiliar with technical terminology find it difficult to provide relevant keywords directly, which affects the model’s tool calling and parameter specification.

## 6 Conclusion

This paper introduces ECom-Bench, the first (to our best knowledge) benchmark including user simulation, personas, multimodal tasks, and MCP tools in e-commerce customer service, and evaluates the reliability and consistency of Large Language Models (LLMs) in solving real-world scenario problems using our designed tasks. We also compare the effectiveness of Multimodal Large Language Models serving as agent versus tools, and examine the impact of user personas on agent performance.

The ablation study on persona-driven user simulation yields valuable insights for benchmark development. The results demonstrate that persona characteristics (e.g., emotional states and patience levels) significantly influence agent reactions. This underscores the importance of accurate user behavior simulation, and the potential impact of real-world customers’ complex behavioral patterns on agent performance.

Our results demonstrate that current large models remain unstable when addressing identical problems. These findings highlight the importance of the benchmark in advancing the development of intelligent customer service agent. Future research will focus on further expansion of sub-scenarios and more realistic scenario simulation, such as incorporating memory mechanisms, optimizing user personas, or better task design.

### Limitations

This work has several limitations: (1) Scope limitation : Our research is confined to the home appliances and furniture vertical domain and has not been extended to other e-commerce categories. While different e-commerce platforms exhibit considerable similarity in interaction patterns, the cross-category generalization capability of our con-

structured benchmark requires further evaluation. (2) Data constraints : Given the sensitivity of commercial data, we employed a large language model-based data synthesis approach, resulting in data that is not entirely authentic. Furthermore, since our benchmark focuses on evaluating agents' capability to address real-world challenges rather than database complexity, we maintain a relatively small-scale dataset with simplified operational methods. (3) Limited task instances : We include 53 task instances, which are representative in terms of typicality and scenario coverage. However, the scale remains limited due to the time-consuming and labor-intensive nature of manual creation. (4) User simulation construction : The construction of user persona feature dimensions primarily relies on domain expertise, lacking sufficient theoretical validation and empirical analysis. Additionally, user simulation depends on LLM-generated behavioral data, which may not fully capture authentic conversational dynamics.

## References

- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2018. [MultiWOZ - a large-scale multi-domain wizard-of-oz dataset for task-oriented dialogue modelling](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5016–5026, Brussels, Belgium. Association for Computational Linguistics.
- Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and Dong Yu. 2024. [Scaling synthetic data creation with 1,000,000,000 personas](#). *arXiv preprint arXiv:2406.20094*.
- Pierre Colombo, Telmo Pessoa Pires, Malik Boudiaf, Dominic Culver, Rui Melo, Caio Corro, André F. T. Martins, Fabrizio Esposito, Vera Lúcia Raposo, Sofia Morgado, and Michael Desa. 2024. [Saullm-7b: A pioneering large language model for law](#). *arXiv preprint arXiv:2403.03883*.
- Neel Guha, Julian Nyarko, Daniel E. Ho, Christopher Ré, Adam Chilton, Aditya Narayana, Alex Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel N. Rockmore, Diego Zambrano, Dmitry Talisman, Enam Hoque, Faiz Surani, Frank Fagan, Galit Sarfaty, Gregory M. Dickinson, Haggai Porat, Jason Hegland, and 21 others. 2023. [Legalbench: A collaboratively built benchmark for measuring legal reasoning in large language models](#). *Preprint*, arXiv:2308.11462.
- Sunwoo Lee, Dhammiko Arya, Seung-Mo Cho, Gyoung-eun Han, Seokyoung Hong, Wonbeom Jang, Seojin Lee, Sohee Park, Sereimony Sek, Injee Song, Sungbin Yoon, and Eric Davis. 2024. [TelBench: A benchmark for evaluating telco-specific large language models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 609–626, Miami, Florida, US. Association for Computational Linguistics.
- Ankit Pal, Logesh Kumar Umaphathi, and Malaikandan Sankarasubbu. 2022. [Medmcqa : A large-scale multi-subject multi-choice dataset for medical domain question answering](#). *arXiv preprint arXiv:2203.14371*.
- Keivalya Pandya and Mehfuza Holia. 2023. [Automating customer service using langchain: Building custom open-source gpt chatbot for organizations](#). *Preprint*, arXiv:2310.05421.
- Yujia Qin, Shengnan Zeng, Michail Burtsev, Mark Steedman, Yixin Zhang, Chi Zhang, and William Yang Wang. 2023. [Tool learning with foundation models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16443–16469.
- Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. 2021. [ALFWorld: Aligning Text and Embodied Environments for Interactive Learning](#). In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Hagen Soltau, Izhak Shafran, Mingqiu Wang, Abhinav Rastogi, Wei Han, and Yuan Cao. 2023. [DSTC-11: Speech aware task-oriented dialog modeling track](#). In *Proceedings of the Eleventh Dialog System Technology Challenge*, pages 226–234, Prague, Czech Republic. Association for Computational Linguistics.
- Lei Wang, Jingsen Zhang, Xu Chen, Yankai Lin, Ruihua Song, Wayne Xin Zhao, and Ji-Rong Wen. 2023. [Recagent: A novel simulation paradigm for recommender systems](#). *arXiv preprint arXiv:2306.02552*.
- Weiqi Wang, Limeng Cui, Xin Liu, Sreyashi Nag, Wenju Xu, Chen Luo, Sheikh Muhammad Sarwar, Yang Li, Hansu Gu, Hui Liu, Changlong Yu, Jiaxin Bai, Yifan Gao, Haiyang Zhang, Qi He, Shuiwang Ji, and Yangqiu Song. 2025. [Ecomscriptbench: A multi-task benchmark for e-commerce script planning via step-wise intention-driven product association](#). *Preprint*, arXiv:2505.15196.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2022. [Chain-of-thought prompting elicits reasoning in large language models](#). In *Advances in Neural Information Processing Systems*, volume 35.
- Wei Xiang, Hanfei Zhu, Suqi Lou, Xinli Chen, Zhenghua Pan, Yuping Jin, Shi Chen, and Lingyun Sun. 2024. [Simuser: Generating usability feedback by simulating various users interacting with mobile](#)

applications. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, CHI '24, New York, NY, USA. Association for Computing Machinery.

Qiantong Xu, Fenglu Hong, Bo Li, Changran Hu, Zhengyu Chen, and Jian Zhang. 2023. *On the tool manipulation capability of open-source large language models*. *Preprint*, arXiv:2305.16504.

John Yang, Akshara Prabhakar, Karthik Narasimhan, and Shunyu Yao. 2023. *Intercode: Standardizing and benchmarking interactive coding with execution feedback*. In *NeurIPS 2023 Datasets and Benchmarks Track*.

Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. 2024.  *$\tau$ -bench: A benchmark for tool-agent-user interaction in real-world domains*. *Preprint*, arXiv:2406.12045.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022. *Webshop: Towards scalable real-world web interaction with grounded language agents*. In *Advances in Neural Information Processing Systems*, volume 35, pages 11602–11615.

Min Zhang, Chuanxiong Ren, G. Alan Wang, and Zhen He. 2021. *Challenges and opportunities in e-commerce operations: A systematic review*. *Electronic Commerce Research and Applications*, 48:101–115.

Shuyan Zhou, Frank F. Xu, Hao Zhang, Graham Neubig, Daniel Fried, and Nazneen Fatema Rajani. 2023. *Webarena: A realistic web environment for building autonomous agents*. *arXiv preprint arXiv:2307.13854*.

## A Appendix: Complete Multimodal Dialogue Trajectories

### A.1 User Settings

#### User Profile:

*Consumer Type:* Value-oriented

#### Personality Traits:

*Emotion:* Calm and rational in problem-solving.

*Attentiveness:* Highly detail-oriented regarding product specifics.

*Patience:* Demonstrates patience without frustration during waiting periods.

*Trust:* Trusts customer service guidance and follows instructions diligently.

*Rights Awareness:* Proactively seeks information to safeguard legitimate interests.

#### Behavioral Traits:

*Inquiry Style:* Direct and unambiguous questions.

*Communication Tone:* Concise and efficient, avoiding excessive formality or familiarity.

*Response Pace:* Maintains appropriate silence until receiving clear replies.

#### User Objectives:

*Intent 1:* Price protection request with image proof. <https://dd-static.jd.com/ddimgp/jfs/t20260623/297694/2/8257/158939/682a8a2aF3a87a48a/4c7eb1c5a42e8829.jpg>

*Intent 2:* Request price difference refund to JD E-card

*Intent 3:* Schedule installation service for Thursday

*Intent 4:* Request personal invoice issuance

### A.2 Failed Multimodal Dialogue Example

**User:** Hello, I previously purchased a product from your store, but I noticed the price has dropped. Here's an image showing the current lower price for your reference.

<https://dd-static.jd.com/ddimgp/jfs/t20260623/297694/2/8257/158939/682a8a2aF3a87a48a/4c7eb1c5a42e8829.jpg>

**Assistant:** Dear customer, your order has successfully qualified for price protection. The difference of 301.35 yuan has been refunded to your JD E-Card.

**User:** Okay, thank you. Next, I'd like to schedule an installation service for this product, preferably on Thursday.

**Assistant:** Dear customer, I couldn't find the installation service information for the product you purchased.

**User:** Could you check the installation service information? Here's the product link: <https://item.jd.com/123456789.html>

**Assistant:** Dear customer, I'm sorry but I couldn't find the product information corresponding to the link you provided.

**User:** Could you check my order information? My order ID is 1234567890.

**Assistant:** Dear customer, I couldn't find any information for order ID 1234567890.

**User:** One last question - I'd like you to issue a personal invoice for me.

**Assistant:** Dear customer, I'll need your order ID, invoice title, and contact number.

**User:** Order ID: 1234567890. Invoice title: Zhang San. Contact: 13800138000.

**Assistant:** Dear customer, I couldn't find any information for order ID 1234567890.

**User:** ###STOP###

### A.3 Successful Multimodal Dialogue Example

**Assistant:** Dear customer, how may I assist you?

**User:** Hello, I noticed the price has dropped. Here's a screenshot showing the current price. <https://dd-static.jd.com/ddimgp/jfs/t20260623/297694/2/8257/158939/682a8a2aF3a87a48a/4c7eb1c5a42e8829.jpg>

**Assistant:** Thank you for the screenshot. The current price is 1707.65 yuan (original: 2009 yuan). Would you like to apply for price protection?

**User:** Yes, please refund the difference to my JD E-Card.

**Assistant:** The difference of 301.35 yuan has been refunded to your JD E-Card.

**User:** Next, I'd like to schedule installation for Thursday.

**Assistant:** Your installation has been scheduled for Thursday between 9:00-12:00.

**User:** Finally, I'd like a personal invoice.

**Assistant:** Please provide the invoice title and contact num



ber.

**User:** The title should be my name. For contact number, please check the order.

**Assistant:** Your personal invoice has been issued (title: 'Lin Yunpei').

**User:** ###STOP###