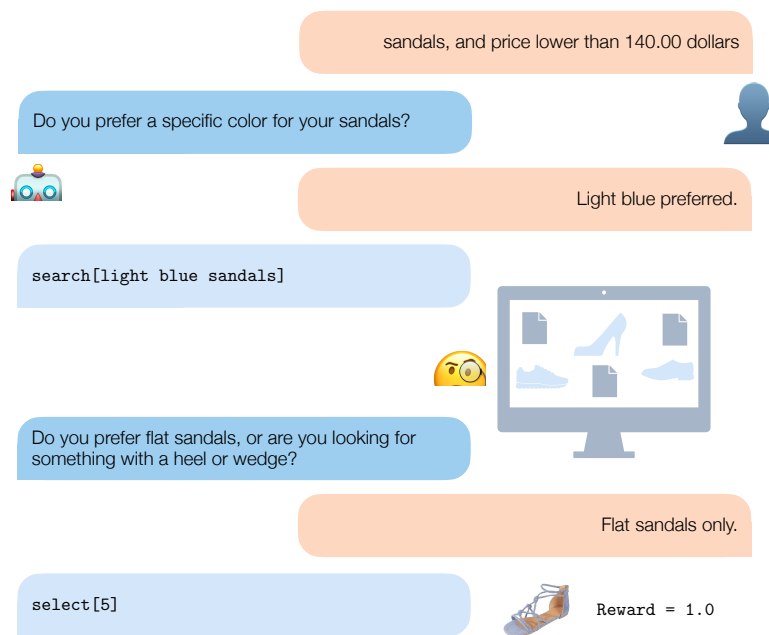


Motivation

- LLM-as-agent has become a popular idea as they can follow instructions and complete many nontrivial tasks.
- Does LLM agents deliberately seek information from the environment, and how strategic is their plan?
- Can we evaluate this automatically in a realistic environment?

Key idea

- Many tasks assume the initial instruction contains sufficient information for an agent to complete the task
- Create a information need by simplifying the instruction and build a communication channel to release information through interactions
- e.g. “a nonslip sandals for my wife that is blue in color, 5.5 size”



Dataset

| | WebShop | ChatShop |
|-------------|---------|----------|
| # Vocab | 2871 | 1166 |
| Avg. Length | 15.1 | 2.3 |

- We repurpose **WebShop**, which is a large-scale web-shopping task with millions of products crawled from amazon.com
- We process the 1500 goal instructions in the dev and test sets of Webshop and obtain the simplified instructions using GPT-3.5 and few-shot prompts.

Experiments

WebShop can be solved non-interactively

- Retriever:** use a BERT-based relevance model for (instruction, product) pairs trained with contrastive loss
- Procedure:** rerank top 50 products obtained from a BM25 search engine using the instruction as query.
- Results:** we achieve a 78.3% success rate and 87.2 average rewards on the dev set, which is superior to the reported 59.6% success rate and 82.1 average rewards of human expert annotator.
- Zero-shot LLMs:** prompting GPTs also leads to near 80 average rewards

| | CodeLlama | GPT-3.5 | GPT-4 |
|------------|-----------|---------|-------|
| None | 34.3 | 43.4 | 48.8 |
| Open-ended | - | 40.6 | 49.7 |
| Instance | - | 40.4 | 51.3 |
| Full Info | 64.5 | 76.0 | 80.1 |

Table 1: Avg. rewards of (*auto q*) agents under different settings of information disclosure. CODELLAMA cannot perform under the interactive settings without advanced prompting strategies.

| Strategy | CodeLlama | | GPT-3.5 | | GPT-4 | |
|-------------------|-----------|------|---------|-------------|-------|------|
| | w/o | w/ | w/o | w/ | w/o | w/ |
| CoT | | | | | | |
| <i>no q</i> | 34.3 | 30.1 | 43.4 | 45.6 | 48.8 | 47.5 |
| <i>auto q</i> | - | - | 40.6 | 62.7 | 49.7 | 59.2 |
| <i>all q</i> | 25.6 | 29.4 | 63.7 | 61.3 | 63.0 | 66.3 |
| <i>interleave</i> | 18.8 | 28.9 | 64.3 | 68.2 | 60.5 | 68.1 |

Table 2: Avg. rewards of agents with different strategies and the open-ended communication channel. *no q* is the non-interactive baselines.

An interactive setup

- Info seeking task:** the agent starts with a simplified instruction
- Agent and Shopper:** a shopper with the intent to purchase an item and an agent that assists the shopper in finding the correct product
- Action Space:**
 - 1) `search[query]`: search a BM25 search engine to get a ranked list of products;
 - 2) `select[index]`: finalize recommendation when a single product is determined;
 - 3) `question[content]`: when more information is needed for a precise decision, the agent can interact with the shopper for further clarification.
- Communication Channel:** 1) open-ended text-based interaction and, 2) instance-based comparison
- Advanced Prompting Strategy:** using heuristics to enforce search and question actions
- Results**
 - basic prompting strategy is inadequate to incentivize the agents to interact with the environment, LLMs are satisfy with partial information
 - CoT or ScratchPad prompting generally help with interaction
 - GPT-3.5 surprisingly outperforms GPT-4
 - the gap between the best agent and the no-interactive full info baseline remains significant

LLM versus Human Shopper

| | GPT-3.5 | GPT-4 |
|-----------|---------|-------|
| Simulated | 59.0 | 62.8 |
| Human | 58.2 | 63.4 |

Table 3: Avg. rewards of LLM agents with simulated and human shoppers over 50 sessions.

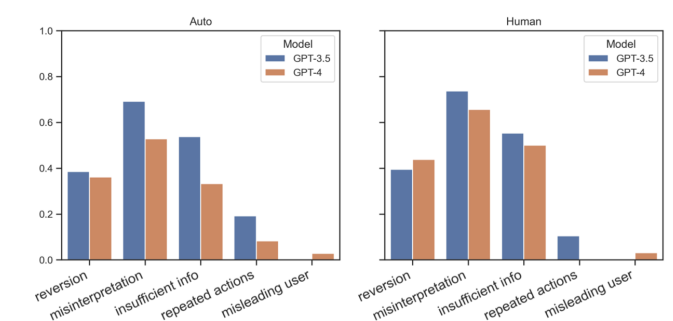


Figure 2: Relative frequency of error types in the LLM agents' failed trajectories with simulated and human shoppers.

Quality of simulation

- LLM agents performance with the simulated shopper and the human shopper are consistent.
- The distributions of automatic categorized failure patterns between the two environments are also similar

Conclusions

- Some agent tasks can be framed as non-interactive retrieval tasks and better solved by smaller models
- Scaling up model size doesn't naturally improve information seeking ability

References

- Bachman et al., “Towards information-seeking agents.” Arxiv’16
- Yao et al., “Webshop: Towards scalable real-world web interaction with grounded language agents.” NeurIPS’22
- Andreas, “Language Models as Agent Models.” EMNLP’22 Findings