AgentInstruct: Agent Instructs Large Language Models to be General Zero-Shot Reasoners

Nicholas Crispino, Kyle Montgomery, Fankun Zeng, Dawn Song, Chenguang Wang



Large Language Models







Large language models exhibit advanced performance on language understanding tasks





Zero-Shot Reasoning

[1st prompt] Reasoning Extraction

Q: On average Joe throws 25 punches per minute. A fight lasts 5 rounds of 3 minutes. How many punches did he throw? A: Let's think step by step.



In one minute, Joe throws 25 punches. In three minutes, Joe throws 3 * 25 = 75 punches. In five rounds, Joe throws 5 * 75 = 375 punches.

Zero-shot reasoning methodologies have proven useful for specific task domains



[2nd prompt] **Answer Extraction**





Language Models as Agents



Agents utilize planning to accomplish more complex tasks







Zero-Shot AgentInstruct: Overview **Zero-Shot Zero-Shot AgentInstruct**

Passage: Alan Rickman & Emma Thompson give good performances with southern/New Orleans accents in this detective flick. It's worth seeing for their scenes- and Rickman's scene with Hal Holbrook.

These three actors mannage to entertain us no matter what the movie, it seems. The plot for the movie shows potential, but one gets the impression in watching the film that it was not pulled off as well as it could have been.

The fact that it is cluttered by a rather uninteresting subplot and mostly uninteresting kidnappers really muddles things.

The movie is worth a view- if for nothing more than entertaining performances by Rickman, Thompson, and Holbrook.

Sentiment: **Positive** X

Zero-Shot CoT

Passage: Alan Rickman & Emma Thompson... Let's think step by step.

...What is the author's conclusion?

The author concludes that the movie is worth a view, mainly for the entertaining performances of the three main actors, despite some flaws in the plot.

Therefore, the correct label among {'Negative', 'Positive'} (just the label) to the original task is below.

Sentiment: Positive X

Our agent produces task-specific instructions to instruct the reasoning process of LLMs







Zero-Shot AgentInstruct: Agent Instructions

Name: IMDB; Labels: Positive, Negative; Input Examples: "Passage: ..."

<u>Thought 1</u>: I need to create instructions for the IMDB dataset. The IMDB dataset is a large database of information related to films, television programs, video games, and streaming content online...
 <u>Action 1</u>: Ask_about_dataset [Can you provide more details about the IMDB Movie Reviews dataset? Specifically, how is the sentiment of the reviews determined and how are the reviews structured?]
 <u>Observation 1</u>: The IMDB Movie Reviews dataset is a binary sentiment analysis dataset that consists of 50,000 reviews from the Internet Movie Database (IMDB). These reviews are labeled as either positive or negative...

<u>Thought 2</u>: The IMDB Movie Reviews dataset is used for binary sentiment analysis...The sentiment of the reviews is determined based on the score given by the reviewer. Only highly polarizing reviews are considered...This information will be useful in creating instructions for the dataset.
<u>Action 2</u>: Finish [Instructions:
1. Understand the Dataset: The IMDB Movie Reviews dataset contains 50,000 reviews... These reviews are labeled as either positive or negative...
2. Analyze the Passage: Pay attention to the language, the tone of the review, and any explicit or implicit expressions of sentiment towards the movie...]

An example of our agent creating task-specific instructions for the IMDB dataset

Agent





Datasets



We evaluated our method on 29 datasets spanning generation, classification, and reasoning





Results



Zero-shot AgentInstruct shows a significant improvement over all 29 datasets on average



Results



Zero-shot AgentInstruct shows an improvement on each category of tasks, especially reasoning



Results



Zero-shot AgentInstruct achieves SoTA results on 20 of the 29 datasets

| Shot CoT 🛛 Zero-Shot AgentInstruct | | | |
|------------------------------------|--|--|--|
| 69.0 (20) | | | |
| (a) Overall | | | |
| 69.2 (9) | | | |
| eration Datasets | | | |
| 68.7 (11) | | | |
| sification Datasets | | | |
| 83.4 (10) | | | |
| | | | |

(d) Reasoning Datasets

Ablation

Zero-Shot AgentInstruct w/o Agent Instructions w/o Input Examples w/o Labels w/o GPT-4

All components of zero-shot AgentInstruct are important

| AddSub | IMDB | NarrativeQA |
|-------------|-------------|-------------|
| 79.5 | 94.0 | 65.0 |
| 73.2 | 89.0 | 62.3 |
| 72.4 | 88.0 | 60.1 |
| 74.9 | 93.8 | 63.9 |
| 75.2 | 92.6 | 63.5 |

Comparison on GPT-4



Zero-shot AgentInstruct is a cost-effective alternative to using agents directly

Context Length



Zero-shot AgentInstruct is sensitive to the context window

Scaling



Llama-2-70b-chat with Zero-shot AgentInstruct outperforms zero-shot ChatGPT by 10.2%



Comparison Methods: Few-Shot



Zero-shot AgentInstruct performs near the level of few-shot prompting

Comparison Methods: Self-Consistency



Zero-shot AgentInstruct exceeds the performance of self-consistency

Conclusion

- Zero-Shot AgentInstruct: Combine an autonomous agent generating instructions with CoT reasoning
- AgentInstruct outperforms zero-shot and zero-shot CoT on generation, classification, and reasoning tasks
- State-of-the-art on 20 of 29 datasets
- Code: <u>https://github.com/wang-research-lab/agentinstruct</u>

