Citation Evaluation Task Guidelines

Project Motivation

Verifying the reliability of responses from language models, like ChatGPT, can be difficult for users. We are exploring the extent to which citations to sources can help users verify outputs.

HIT Explanation

Navigate to the HIT website linked in the MTurk HIT and enter your username specified on the MTurk HIT page. Then, the HIT website will show 5 tasks for evaluation. Each task is about a query and a response. The evaluation for each task has several parts, as described below. At the end of the five tasks, you will receive a one-time-use completion code. Enter this code into the MTurk HIT and submit for the reward.

Please note that the responses are timed, which is an important measurement for the project. You do not need to rush or feel pressured, but please plan to avoid multitasking during the HIT.

Task-level Questions (not timed)

For each task, you will answer the fluency and utility questions once. Consider the full text of the response, even if parts are highlighted.

Fluency Question

To what extent is the response fluent and coherent?

Multiple choice options:

- 1: The response has noticeable misprints or disfluent transitions and sentences
- 2: The response has no misprints and mostly smooth transitions and sentences
- 3: The response has no misprints and all of the sentences flow nicely together

Utility Question

To what extent does the response seem to be a useful answer to the query? Multiple choice options:

- 1: The response includes too many irrelevant details or the query is not addressed
- 2: The response is only a partially satisfying answer to the query
- 3: The response is concise and seems to be a satisfying answer to the query

Skills

Before we review the 2 questions that follow the Fluency and Utility Questions, we will review 4 necessary skills.

Skill 1: Understanding the information conveyed by a sentence

Often, the information communicated by a sentence in the response is exactly what's written in the sentence. However, there are some cases where information is implied.

In the example below, the highlighted sentence starts off with "This recognition", which refers to "The inaugural Nobel Prize in Physics" from the previous sentence. This isn't explicitly written in the highlighted sentence, but it matters to its meaning.

User Query: who got the first nobel prize in physics

Cited System Response:

The inaugural Nobel Prize in Physics was bestowed upon Wilhelm Röntgen, a German physicist, for his groundbreaking discovery of X-rays [1]. This recognition in 1901 highlighted the significant impact of his work on both the field of physics and medical diagnostics, setting a precedent for future awards [1].

In the example below, the highlighted response is not a full sentence. You may interpret the highlighted response as Old Trafford is the name of the Manchester United stadium, even though this isn't explicitly written in the response.

User Query: what is the name of manchester united stadium

Cited System Response:

Old Trafford [1].

Skill 2: Finding citations

For these questions, we highlight one cited sentence in the response at a time. We also provide the source of each citation in that sentence. Find citations by matching the citation number at the end of a highlighted sentence in the response to the citation number at the start of a highlighted sentence in a source.

See the example below, where we can find citation [2]. Note that you can scroll through the sources on the right hand side.

User Query: where does the autobahn have no speed limit

Cited System Response:

"Much of the system has no speed limit for some classes of vehicles" [1]. "More than half of the total length of the German autobahn network has no speed limit" [2].

Source: en.wikipedia.org/wiki/Autobahn

[2] More than half of the total length of the German autobahn network has no speed limit, about one third has a permanent limit, and the remaining parts have a temporary or conditional limit. Some cars with very powerful engines

Skill 3: Evaluating citations

Both questions require determining whether a citation supports some information. When doing this, you should consider not only the highlighted text of the citation, but also the entire source it comes from. When determining whether a citation "supports" information in the sentence (X), you can judge whether it is correct to state: "According to the source, X." You do not need to access external websites.

In the example below, the source of citation [3] *supports* information in the highlighted sentence in the response. This is the case even though the relevant sections of the source happen to lie just after the highlighted section.

User Query: where did the idea of fortnite come from

Cited System Response:

Fortnite's origins trace back to Epic Games, evolving significantly from its initial concept [1]. Originally envisioned as a different genre focused on surviving a zombie apocalypse, Fortnite: Save The World was designed as a cooperative third-person shooter with survival elements [3].

Source: thegamer.com/fortnite-history-explained/

Fortnite has become a global phenomenon and is growing quickly within the realm of esports. Its revenue stream for 2019 is expected to reach \$3.5 billion, and it shows no signs of losing its grip over the Battle Royale genre. However, few fans know how Fortnite became the phenomenon it is today. [1] The history of Fortnite begins with Epic Games, but it was originally far from the game we know and love today. [2] Instead, Fortnite began as an entirely different game about surviving a zombie apocalypse... sort of. [3] Fortnite: Save The World The game was first revealed at the 2011 Spike Video Game Awards, where Epic's former design director Cliff Bleszinski presented a rough trailer for the game. Fortnite: Save The World was initially conceived as a co-op third-person shooter survival game.

In the example below, the source of citation [1] *does not* support the highlighted sentence in the response.

User Query: how many countries are a part of opec

Cited System Response:

As of January 2020, OPEC has 13 member countries [1].

Source: en.wikipedia.org/wiki/OPEC

countries[edit] [1] As of January 2020, OPEC has 12

member countries: five in the Middle East (West Asia), six in

Africa, and one in South America. According to the U.S.

Skill 4: Evaluating quotations

Sometimes, the response will contain quotations. You can trust that every quotation is word-for-word from the cited source. However, there is still the possibility that the quote is taken out of context. Even though the quotations are exact, still use Skills 1, 2, and 3 to evaluate quoted citations.

In the example below, the source of citation [2] *does not* support the highlighted sentence in the response, even though it's an exact quote. This is because the quote in the response implies that electrolytically isolating lithium consumes more than three-quarters of its production; that's not what the source originally meant by the quotation.

User Query: what is the main mineral in lithium batteries

Cited System Response:

"Lithium metal is isolated electrolytically from a mixture of lithium chloride and potassium chloride" [1]. "These uses consume more than three-quarters of lithium production" [2].

Source: en.wikipedia.org/wiki/Lithium

staged thermonuclear weapons.[8] Lithium and its compounds have several industrial applications, including heat-resistant glass and ceramics, lithium grease lubricants, flux additives for iron, steel and aluminium production, lithium metal batteries, and lithium-ion batteries. [2] These uses consume more than three-quarters of lithium production.[citation]

Sentence-level Questions (timed)

For each sentence in a task's response, we ask two questions about whether cited sources support the sentence.

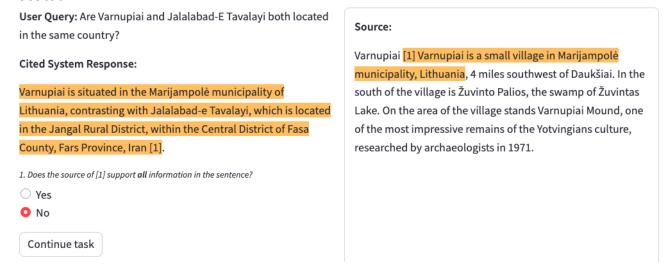
Question 1: Citation Coverage (timed)

1. Do the sources of the citations together support **all** information in the sentence?

If the sentence contains information not supported by the sources of its citation(s), then select "No". Otherwise, select "Yes".

We are interested in this question because we want to know whether everything in the highlighted sentence in the response is supported by the sources of its citations.

In the example below, the source of citation [1] does not support all information in the highlighted sentence. The source does not specify where Jalalabad-e Tavalayi is located.



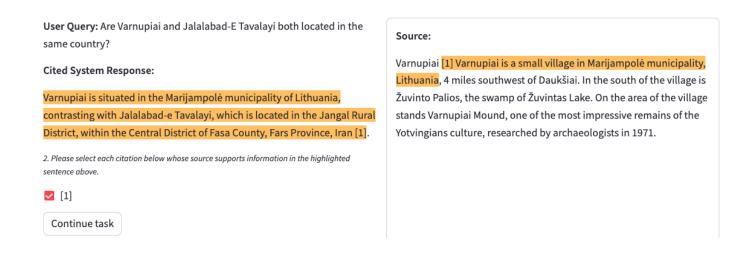
Question 2: Citation Correctness (timed)

2. Please select each citation below whose source supports information in the highlighted sentence above.

A checklist of citations (that match the ones in the highlighted sentence of the response) will be made available. Check the box for a citation if its source supports some information in the sentence, even if another citation supports the same information.

We are interested in this question because we want to identify any citations that do not help to support the highlighted sentence in the response.

The example below shows how the source of citation [1] supports information in the highlighted sentence in the response, even though it doesn't cover the full sentence.



Tips for using the website

- 1. Once annotations are inputted, they are final. Please do not use the back button.
- Once you login, please finish all tasks and submit the completion code to the MTurk HIT website before closing the window. Submitting the completion code is critical for receiving the reward.

Looking forward to working with you!

Feel free to contact us through the MTurk site with any questions.