

# ClaimVer: Explainable Claim-Level Verification and Evidence Attribution of Text Through Knowledge Graphs

Preetam Prabhu Srikar Dammu<sup>1</sup>, Himanshu Naidu<sup>1</sup>, Mouly Dewan<sup>1</sup>, YoungMin Kim<sup>1</sup>, Tanya Roosta<sup>2,4,\*</sup>, Aman Chadha<sup>3,4,\*</sup> and Chirag Shah<sup>1</sup>

<sup>1</sup>University of Washington

<sup>2</sup>UC Berkeley

<sup>3</sup>Stanford University

<sup>4</sup>Amazon GenAI

## Abstract

In the midst of widespread misinformation and disinformation through social media and the proliferation of AI-generated texts, it has become increasingly difficult for people to validate and trust information they encounter. Many fact-checking approaches and tools have been developed, but they often lack appropriate explainability or granularity to be useful in various contexts. A text validation method that is easy to use, accessible, and can perform fine-grained evidence attribution has become crucial. More importantly, building user trust in such a method requires presenting the rationale behind each prediction, as research shows this significantly influences people’s belief in automated systems. It is also paramount to localize and bring users’ attention to the specific problematic content, instead of providing simple blanket labels. In this paper, we present *ClaimVer*, a *human-centric framework* tailored to meet users’ informational and verification needs by generating rich annotations and thereby reducing cognitive load. Designed to deliver comprehensive evaluations of texts, it highlights each claim, verifies it against a trusted knowledge graph (KG), presents the evidence, and provides succinct, clear explanations for each claim prediction. Finally, our framework introduces an attribution score, enhancing applicability across a wide range of downstream tasks.

## 1 Introduction

Misinformation and disinformation are longstanding issues, but the proliferation of AI tools that can generate information on demand has amplified these issues. Tools for fact-checking are not keeping pace with sophisticated text generation techniques. Even when they are effective, they lack appropriate explainability and granularity to be useful to users. Studies have shown that explanations are crucial for users to build trust in AI systems [Rechkemmer and Yin, 2022; Weitz *et al.*, 2019; Shin, 2021]. There is a need for a novel

\*Work does not relate to position at Amazon.

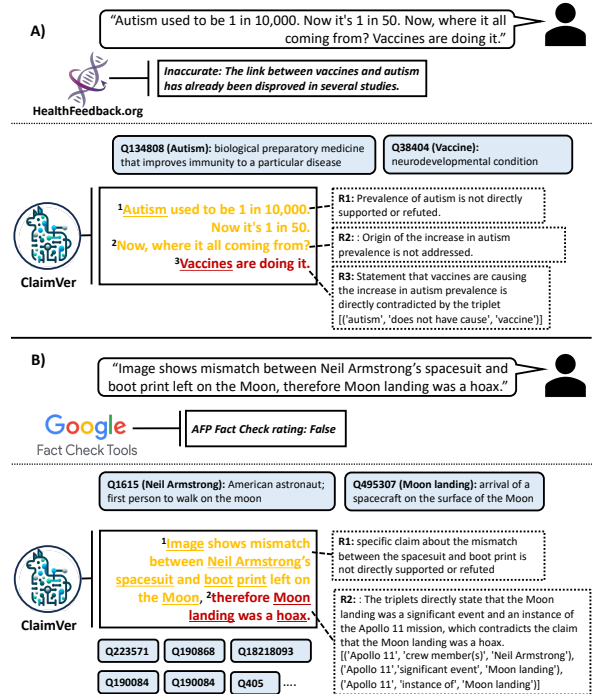


Figure 1: Demonstration of ClaimVer for claim verification and evidence attribution. (A) Text labeled as *Inaccurate* by HealthFeedback and ClaimVer’s predictions, rationale, and evidence. (B) Text labeled as *False* by Google Fact Check Tools and ClaimVer’s outputs. Predictions are color-coded (amber: extrapolatory, red: contradictory);  $R_i$ : rationale; related wiki entities are displayed in boxes.

human-centric approach to text verification that provides usable and appropriately granular explanations that can not only inform but also educate the user.

Most fact-checkers, including widely used ones in deployment, issue blanket predictions that can lead to user misunderstanding. For instance, in Figure 1 (A), we observe that HealthFeedback<sup>1</sup>, a fact-checker for medical text, indicates that a misleading statement about the increase in Autism is inaccurate. However, there are multiple claims made in that text, which are not addressed by this tool. In fact, research

<sup>1</sup><https://healthfeedback.org/>

does show that Autism cases have increased, but this is mostly attributed to increased testing [Russell *et al.*, 2015]. Our method accurately breaks down the text into multiple claims and shows that the specific claim that vaccines are causing autism is indeed incorrect, attributing it to a fact from the Wikidata [Vrandečić and Krötzsch, 2014]. It also provides a clear rationale as to why the first two claims cannot be determined, as there’s no conclusive evidence present in the KG. Such granular predictions, supported by justifications, significantly improve user confidence [Rechkemmer and Yin, 2022; Weitz *et al.*, 2019; Shin, 2021].

Similarly, in Figure 1 (B), we notice that Google Fact Check Tools<sup>2</sup> simply provides a blanket label for an utterance denying the moon landing. In contrast, ClaimVer identifies the exact text span that can be conclusively proven incorrect and proceeds to provide specific information about the Apollo 11 mission and its crew members to refute the claim. All verified entities present in the text, along with their Wiki IDs and descriptions, are displayed for user reference.

Prior research [Rashkin *et al.*, 2023; Yue *et al.*, 2023; Thorne *et al.*, 2019; Aly *et al.*, 2021] typically validates text at the paragraph or sentence level without adequately enhancing user awareness by supplying key details such as rationale, match scores, or evidence. A KG-based approach allows for finer granularity, aiding in pinpointing specific inaccuracies like hallucinations in LLM-generated text or false claims in misleading text. Furthermore, if needed, broader-level metrics can be extracted from this detailed attribution.

The assumption of one-to-one mapping between input and reference texts, prevalent in previous methods [Rashkin *et al.*, 2023; Yue *et al.*, 2023; Thorne *et al.*, 2019; Aly *et al.*, 2021], does not hold if the given text consists of claims that can be mapped to more than one source. In contrast, utilizing a KG, which represents a consolidated body of knowledge, results in a more comprehensive evaluation. While most previous methods may not support scenarios with information spread across various references, querying a KG can yield triplets originally sourced from multiple documents. Additionally, procuring the specific spans of text required to evaluate claims, from large text sources that may span several pages, presents many challenges. On the other hand, a KG captures only the most important relationships as nodes and links, and offers a more efficient way to evaluate the claims.

Prior methods that depend on document indices or vector databases are not easy to maintain or audit. In contrast, existing trusted KGs that are constructed through human curation provide an effective and human-centered approach for evaluating text at scale. Therefore, we leverage KGs to build a framework that realizes our goal of performing fine-grained text verification and evidence attribution. Our framework also generates insights that boost user awareness, thereby fostering increased trust in automated systems.

## 2 Related Work

Research on validating text has been ongoing for the past decade, while the concept of evidence attribution has gained

<sup>2</sup><https://toolbox.google.com/factcheck/explorer>

increased attention in recent years, following the advent of generative models.

Our method integrates fact verification and evidence attribution; therefore, we discuss recent advancements in both domains in this section.

### 2.1 Fact Verification

Fact verification is a task that is closely related to natural language inference (NLI) [Conneau *et al.*, 2017; Schick and Schütze, 2020], in which given a premise, the task is to verify whether a hypothesis is an entailment, contradiction, or neutral. Similarly, in fact verification, the task is to check if a given text can be supported, refuted, or indeterminable, given a reference text. Recent studies in this domain show that LLMs can achieve high performance, and can be considerably reliable for verification tasks, even though they are prone to hallucinations [Guan *et al.*, 2023].

In [Lee *et al.*, 2020], the authors show that the inherent knowledge of LLMs could be used to perform fact verification. Other works [Yao *et al.*, 2022; Jiang *et al.*, 2023b] have shown that using external knowledge is helpful for many reasoning-intensive tasks, and report enhanced performance on HotPotQA [Yang *et al.*, 2018] and FEVER [Thorne *et al.*, 2018]. A wide variety of studies have established LLMs are suitable for fact verification. For example, [Dong and Smith, 2021] enhanced accuracy of table-based fact verification by incorporating column-level cell rank information into pre-training. In FactScore, authors [Min *et al.*, 2023], introduce a new evaluation that breaks a long-form text generated by large language models (LMs) into individual atomic facts and calculates the proportion of these atomic facts that are substantiated by a credible knowledge base.

### 2.2 Evidence Attribution

The distinction between evidence attribution and fact verification lies in the emphasis on identifying a source that can be attributed to the information. This task is becoming increasingly important, as generative models produce useful and impressive outputs, but without a frame of reference to validate them. In [Rashkin *et al.*, 2023], the authors present a framework named AIS (Attributable to Identified Sources) that specifies annotation guidelines and underlines the importance of attributing text to an external, verifiable, and independent source. [Yue *et al.*, 2023] demonstrate that LLMs can be utilized for automatic evaluation of attribution, operationalizing the guidelines presented in [Rashkin *et al.*, 2023]. However, both of these works are primarily designed for the question-answering (QA) task, a primary end-user application for LLMs like ChatGPT [Achiam *et al.*, 2023]. In contrast, our method is not restricted to QA and is designed to work with text in general. Furthermore, while these previous studies focus on sentence or paragraph levels, our approach extends to a more detailed and granular level of analysis.

## 3 Methodology

In this section, we present the methodology for retrieving relevant triplets from the KG, fine-tuning LLM to process text at

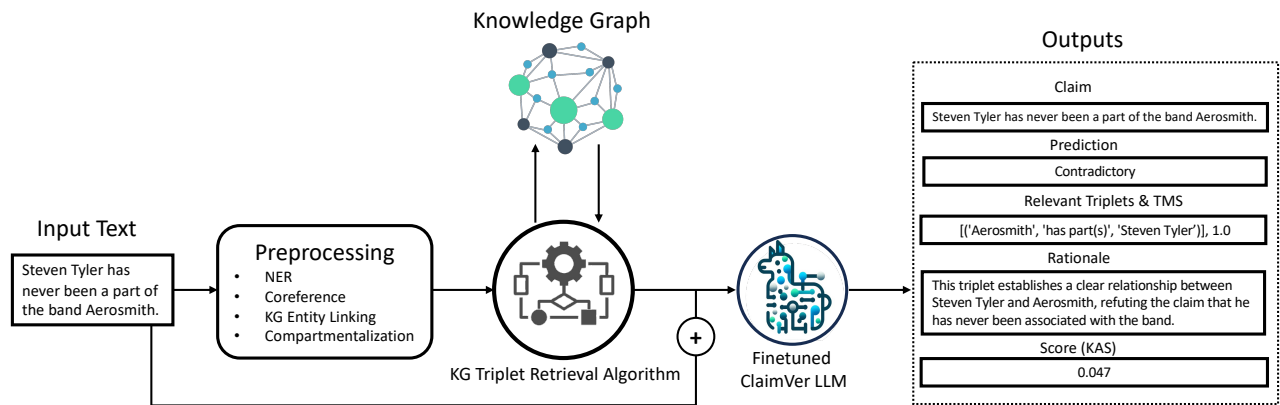


Figure 2: Flow of Operations in the ClaimVer framework. Identified KG entity nodes during preprocessing inform the extraction of relevant triplets by the KG algorithm. Subsequently, these triplets and preprocessed text are then fed to a ClaimVer LLM, fine-tuned to operationalize the objective function. For each claim, the corresponding text span, prediction, relevant triplets, attribution scores, and rationale are generated.

claim-level, verifying claims, tagging evidence for each prediction, and generating a rationale along with an attribution score that reflects the text’s validity.

### 3.1 Preprocessing

Preprocessing involves multiple steps required to make the input text suitable for the subsequent operations. Since the nodes in a KG typically represent entities, performing Named Entity Recognition (NER) is necessary. In our work, we chose Wikidata [Vrandečić and Krötzsch, 2014] as the KG source; thus, we use an NER module suitable for Wiki entities [Gerber, 2023]. However, the framework is sufficiently generic to support any kind of KG that models information in the form of triplets. As our analysis is performed at the claim level, coreference resolution [Lee *et al.*, 2017] becomes a necessary step to form localized claims that are semantically self-contained. If input text exceeds the context length, which depends on design choices, compartmentalization would be required. As a final step in preprocessing, we perform KG entity linking. This step tags all entities in the text that are present in the KG as nodes.

### 3.2 Relevant Triplets Retrieval

Retrieving relevant triplets is a complex problem that has attracted attention from various research communities, and resulted in multiple approaches to address the challenge. While retrieving direct links between two given nodes in a KG is relatively straightforward, identifying complex paths that involve multiple hops is challenging. In our framework, we use Woolnet [Gutiérrez and Patricio, 2023], a multi-node Breadth-First Search (BFS) algorithm, to retrieve the most relevant triplets for a given claim present in the KG. This BFS algorithm initiates from multiple starting points and, at each step, searches for and processes all adjacent neighbors before advancing. It constructs a subgraph of visited nodes, tracking their origins, and distances from each BFS’s start. The algorithm expands each search tree one node at a time until paths intersect or reach a predefined maximum length. Upon intersection, it assesses if the discovered path meets the length criteria. If so, it logs the route, utilizing backtracking to trace

the path to its origins, while ensuring there are no repetitions or cycles, thus maintaining a connection to a starting node. In our experiments, we allow for a maximum of three hops between any two given nodes, and a maximum of four potential paths. Adopting less stringent conditions leads to less relevant triplets.

### 3.3 Objective Function

Previous works on evidence attribution tasks have established definitions for the categorization of input text with reference to a supporting source [Rashkin *et al.*, 2023; Gao *et al.*, 2023; Bohnet *et al.*, 2022; Yue *et al.*, 2023]. Similar to the formulation in [Yue *et al.*, 2023], we use three categories: *Attributable*, *Extrapolatory*, and *Contradictory*. However, there are two main differences that distinguish our approach from previous methods. First, we verify the input text against facts present in a KG, an aggregated information source constructed by integrating numerous data sources into a structure of triplets, instead of relying on a single reference. This approach eliminates the one-to-one dependency between the text and its information source. Second, we perform attribution with finer granularity, specifically at the claim level, involving a subtask of decomposing the input text into individual claims. We define our categories as follows:

- **Attributable:** The triplets fully support the claim.
- **Extrapolatory:** The triplets lack sufficient information to evaluate the claim.
- **Contradictory:** The triplets contradict the claim.

We formulate the objective function of our task as follows:

$$f(input\_text, ret\_triplets) = \{(claim\_span_i, claim\_pred_i, rel\_triplets_i, rationale_i)\}_{i=1}^n \quad (1)$$

where:

- *input\_text*: The input text containing claim(s).
- *ret\_triplets*: A set of *retrieved* triplets for the input text.

- $claim\_span_i$ : The  $i^{th}$  claim extracted as a substring from  $input\_text$ .
- $claim\_pred_i$ : The predicted label for  $claim\_span_i$ .
- $rel\_triplets_i$ : A relevant subset of  $ret\_triplets$  supporting, refuting, or are extrapolatory for  $claim\_span_i$ .
- $rationale_i$ : Justification for  $claim\_pred_i$ .
- $n$ : The total number of claims automatically extracted from  $input\_text$ .

This objective function encompasses two main sub-tasks:

1. Decomposing input text into claims.
2. Generating prediction and corresponding rationale for each claim by identifying relevant supporting triplets.

### 3.4 Fine-tuning LLMs

The objective function shares similarities with the well-studied task of NLI [Conneau *et al.*, 2017; Schick and Schütze, 2020]. LLMs achieve state-of-the-art performance for NLI [Chowdhery *et al.*, 2023], making them a suitable choice to operationalize the objective function. Additionally, [Yue *et al.*, 2023] shows that LLMs can be used to automatically evaluate attribution to a given information source. However, these prior methods do not involve a complex sub-task, which is central to the proposed objective function, i.e., decomposing the input text into text spans that correspond to separate claims in the presence of multiple claims.

It is crucial to perform both claim decomposition and attribution for all claims in a single step, as processing each claim individually can lead to an exponential increase in LLM queries, leading to significantly higher computational costs and latency issues.

In order to perform attribution at the claim level, we need to fine-tune LLMs specifically for the proposed objective function (see §3.3) using a custom dataset. This is necessary because, as of this writing, even the state-of-the-art model, OpenAI’s GPT-4 [Achiam *et al.*, 2023], does not perform satisfactorily right out of the box. Further details on the dataset are provided in §4. On every prediction, a membership check for relevant triplets is performed for additional verification.

We selected five open-source LLMs with diverse sizes, ranging from 2.7B parameters to 13B parameters, to perform the fine-tuning: Phi-2 2.7B [Jawaheripi *et al.*, 2023], Mistral-Instruct 7B [Jiang *et al.*, 2023a], Zephyr-Beta 7B [Tunstall *et al.*, 2023], Solar-Instruct 10.7B [Kim *et al.*, 2023], and Llama2-chat 13B [Touvron *et al.*, 2023]. The models were fine-tuned using LoRA [Hu *et al.*, 2021] with 4-bit quantization and adapters with rank 8 [Dettmers *et al.*, 2024]. The context length was set to 1024 tokens. All models converged after 2 epochs, and high ROUGE-L [Lin, 2004] scores greater than 0.635 were achieved for each model. The instruction prompt used for fine-tuning is presented in Figure. 3.

### 3.5 Computing Attribution Scores

For various downstream tasks, such as ranking and filtering, a continuous score that reflects the validity of a given piece of text with respect to a KG is desirable. We propose the KG Attribution Score (KAS), which accomplishes this task with a high level of granularity, and is detailed in this section.

```
Analyze text against provided triplets, classifying
claims as "Attributable", "Contradictory", or
"Extrapolatory".
Justify your classification using the following
structure:
- "text_span": Text under evaluation.
- "prediction": Category of the text
(Attributable/Contradictory/Extrapolatory).
- "triplets": Relevant triplets (if any, else "NA").
- "rationale": Reason for classification.
For multiple claims, number each component (e.g.,
"text_span1", "prediction1",...). Use "NA" for
inapplicable keys.
Example:
"text_span1": "Specific claim",
"prediction1": "Attributable/Contradictory/Extrapolatory",
"triplets1": "Relevant triplets",
"rationale1": "Prediction justification",
...
Input for analysis:
-Text: {Input Text}
-Triplets: {Retrieved Triplets}
```

Figure 3: Instruction prompt for operationalizing objective function.

#### Claim Scores

$$cs(y_i) = \begin{cases} 2 & \text{if } y_i = \text{Attributable} \\ 1 & \text{if } y_i = \text{Extrapolatory and } |rel\_triplets_i| > 0 \\ 0 & \text{if } y_i = \text{Extrapolatory and } |rel\_triplets_i| = 0 \\ 0 & \text{if } y_i = \text{No attribution} \\ -1 & \text{if } y_i = \text{Contradictory} \end{cases} \quad (2)$$

where,  $y_i$  is  $claim\_pred_i$ .

For each claim, we assign a score that reflects the level of its validity, ranging from -1 (*contradictory*) to 2 (*attributable*). If a claim is predicted to be *extrapolatory*, yet has one or more relevant triplets, we assign that claim a score of 1, as there is still relevant information available even though it may not be sufficient to completely support or refute. However, if there are no triplets at all, along with an *extrapolatory* prediction, we assign 0 as it does not add much information. While decomposing claims, the model might occasionally omit words, typically stop-words, and we assign 0 in those cases as well.

#### Triplets Match Score (TMS)

This score reflects the extent of the match between the relevant triplets and the corresponding claim, and it can also serve as a proxy for prediction confidence. Even though the prediction is made at the claim level, the triplets match score considers word-level matches in the computation. It can be computed as follows:

$$TMS(E(claim\_span_i), E(rel\_triplet_i)) = \alpha \cdot SS(E(claim\_span_i), E(rel\_triplet_i)) + \beta \cdot EPR(E(claim\_span_i), E(rel\_triplet_i)) \quad (3)$$

where,  $E(claim\_span_i)$  and  $E(rel\_triplet_i)$  represent the sets of entities in  $claim\_span_i$  and  $rel\_triplet_i$ , respectively.  $SS$  is the semantic similarity computed using the cosine similarity of text embeddings, and  $EPR$  represents the ratio of entities in  $E(claim\_span_i)$  that are also present in  $E(rel\_triplet_i)$ . The parameters  $\alpha$  and  $\beta$  can be adjusted

Split	Samples	Claims	Claim Labels		
			Att	Ext	Con
Train	3400	5343	2964	1485	894
Test	1000	1675	858	492	325

Table 1: Distribution of fine-tuning dataset. Att: Attributable, Ext: Extrapolatory, Con: Contradictory.

as needed; in our experiments, we use 0.5 for both. In cases where examples of an entity retrieved from the KG are used to support the prediction, instead of the entity itself, we may not have a direct overlap, and thus semantic similarity would be helpful. *EPR* rewards the direct use of the entity, so a balance between both may be ideal in most cases.

### KG Attribution Score (KAS)

For the final KG Attribution Score (KAS), a continuous score between 0 and 1 is desirable, as this facilitates various downstream applications such as ranking, fine-tuning, and filtering. This can be achieved using a Sigmoid function. However, the standard Sigmoid function treats positive and negative scores equally. In most cases, higher penalties should be assigned for erroneous text than rewards for valid text. This requirement can be met using a modified sigmoid function that penalizes mistakes by a factor of  $\gamma$ :

$$\sigma_{\text{mod}}(x, \gamma) = \frac{1}{1 + e^{-\gamma \cdot x}}, \quad (4)$$

where  $\gamma = \begin{cases} \gamma = 3 & \text{if } x < 0, \\ \gamma = 1 & \text{if } x \geq 0, \end{cases}$

In our experiments, we set the value of  $\gamma$  to 3. Finally, the modified Sigmoid function, applied to the summation of triplet match scores and claim scores, is used to generate KAS:

$$\text{KAS} = \sigma_{\text{mod}}\left(\sum_{i=1}^n [TMS_i \cdot cs(y_i)], \gamma\right) \quad (5)$$

## 4 Dataset

Open-domain Question Answering (QA) datasets, such as WikiQA [Yang *et al.*, 2015], HotPotQA [Yang *et al.*, 2018], PopQA [Mallen *et al.*, 2022], and EntityQuestions [Sciavolino *et al.*, 2021], as well as Fact Verification datasets like FEVER [Thorne *et al.*, 2019], FEVEROUS [Aly *et al.*, 2021], TabFacT [Chen *et al.*, 2019], and SEM-TAB-FACTS [Wang *et al.*, 2021], provide texts along with corresponding reference contexts or attributable information sources. However, these datasets significantly differ from the type of data required to train and test our proposed objective function, primarily due to two major factors: First, these datasets predominantly offer samples that are inherently *attributable*. To address this limitation, prior work [Yue *et al.*, 2023] in attribution evaluation introduced new samples by modifying correct answers to generate *contradictory* instances. However, this adjustment alone is not adequate for our use case because, our method requires attribution at the claim level, and necessitates the automatic decomposition of claims. Consequently,

as this task represents a novel challenge, we developed a new dataset that enables effective training and testing of the objective function.

Considering the choice of our KG, which is Wikidata [Vrandečić and Krötzsch, 2014], we opted for WikiQA [Yang *et al.*, 2015] as it is closely associated with the Wiki ecosystem. Given that our method is designed for text validation in general, not limited to question answering, we retain only answers and discard the questions. Subsequently, we processed the answers following the steps detailed in Section 3.1, selecting entries containing two or more Wiki entities. This approach resulted in the exclusion of most single-word answers and other responses that are dependent on their corresponding questions and may lack comprehensibility without them.

We utilize GPT-4 [Achiam *et al.*, 2023] to generate the initial version of the ground truth, as knowledge distillation [Gou *et al.*, 2021] has proven to be an effective strategy. Although GPT-4 can adhere to the instructions (refer to Figure 3) to a reasonable degree and responds in the required format with all necessary keys, it still underperforms in the overall task. The most frequent issue observed is the erroneous assignment of prediction labels. After post-processing, we conducted manual checks to ensure only high-quality samples were retained, as research indicates that high alignment can be achieved with as few as 1,000 samples, provided they are of superior quality [Zhou *et al.*, 2023].

The final dataset is comprised of two splits: the training split, based on the training split of WikiQA [Yang *et al.*, 2015], and a test split, derived from both the test and validation splits. The training split contains 3,400 samples, and since some entries feature multiple claims, there are a total of 5,343 claims within this split. Similarly, the test split includes 1,000 samples and 1,675 claims. The label counts for the claims are tabulated in Table 1.

## 5 Experiments and Results

In this section, we present the evaluation of our claim-level attribution method. The performance metrics of the fine-tuned LLMs, which operationalize the objective function, are presented in Tables 3 and 4. In Table 3, we observe that all models converge and achieve sufficiently high ROUGE-L and ROUGE-1 scores, with *Solar-Instruct* achieving the highest of 0.655 and 0.693 respectively. We also observe that the smaller model, *Phi2* with just 2.7B parameters, is also sufficiently compatible for this task as it attained a decent ROUGE-L score of 0.635.

The average accuracies on the test set, however, vary significantly across the models, with *Solar-Instruct* reporting the highest at 89.31%. The reason behind this variation is that these scores account for both sub-tasks of the objective function (refer §3.3): decomposing input text into claims and generating predictions for each claim. The first task, decomposing text into multiple claims, is somewhat subjective, and there could be multiple valid approaches due to linguistic complexities. We impose a strict strategy while computing accuracy: the text span of the claim, the identified relevant triplets, and the prediction label must all exactly match the ground truth to be considered accurate. Since these LLMs

	Input Text	Relevant Triplets	Prediction (TMS)	Rationale	KAS
1	<sup>1</sup> George O'Malley is a fictional character from the medical drama television series Grey's Anatomy, <sup>2</sup> which airs on the American Broadcasting Company (ABC) in the United States.	1: [{"Grey's Anatomy", 'characters', "George O'Malley"}] 2: [{"Grey's Anatomy", 'original broadcaster', 'American Broadcasting Company'}, {'American Broadcasting Company', 'country', 'United States of America'}]	1: Attributable (0.852) 2: Attributable (0.637)	1: The triplet directly supports the claim that George O'Malley is a character in Grey's Anatomy. 2: The triplets confirm that Grey's Anatomy airs on ABC, which is based in the United States, directly supporting the claim about the show's broadcasting and location.	0.818
2	<sup>1</sup> Bane was portrayed as a tertiary villain by Robert Swenson in Batman & Robin, <sup>2</sup> directed by Joel Schumacher, <sup>3</sup> and Batman's Back	1: [{"Batman & Robin", 'cast member', 'Robert Swenson'}] 2: [{"Batman & Robin", 'director', 'Joel Schumacher'}] 3: NA	1: Attributable (0.788) 2: Attributable (0.882) 3: Extrapolatory (0.0)	1: The triplet directly supports the claim that Robert Swenson was involved in Batman & Robin, which is a requirement for the statement about his portrayal of Bane. 2: The triplet directly supports the claim about the director of Batman & Robin, which is relevant to the context of the film. 3: There are no triplets that directly support or refute the claim about Batman's back	0.752
3	<sup>1</sup> Crater Lake is the main feature of Crater Lake National Park <sup>2</sup> and famous for its deep blue color and water clarity.	1: [{"Crater Lake", 'located in protected area', 'Crater Lake National Park'}] 2: NA	1: Attributable (0.942) 2: Extrapolatory (0.0)	1: The triplet directly supports the claim that Crater Lake is a significant feature within Crater Lake National Park, as it is located within the protected area. 2: There are no triplets provided that directly support or refute the claim about the deep blue color and water clarity of Crater Lake.	0.719
4	<sup>1</sup> Based in Blagnac, France, a suburb of Toulouse, <sup>2</sup> and with significant activity across Europe, <sup>3</sup> airbus produces approximately half of the world's jet airliners.	1: [{"Airbus Operations S.A.S.", 'country', 'France'}, {'Airbus Corporate Jets', 'headquarters location', 'Toulouse'}, {'Blagnac', 'country', 'France'}] 2: NA 3: NA	1: Attributable (0.505) 2: Extrapolatory (0.0) 3: Extrapolatory (0.0)	1: The triplets confirm that Airbus Operations S.A.S. is in France, Airbus Corporate Jets is headquartered in Toulouse, and Blagnac is a suburb of Toulouse in France, supporting the statement about Airbus's location in France and its proximity to Toulouse. 2: The triplets do not provide information about Airbus's activity across Europe 3: The triplets do not provide any information about Airbus's production output or market share	0.583
5	<sup>1</sup> Pope Benedict XVI never appointed anyone significant within the Catholic Church, <sup>2</sup> nor did he ever teach the importance of understanding God's redemptive love.	1: [{"Rutilio del Riego Jáñez", 'appointed by', 'Benedict XVI'}, {'Rutilio del Riego Jáñez', 'religion or worldview', 'Catholic Church'}] 2: [{"God", 'said to be the same as', 'love'}]	1: Contradictory (0.781) 2: Extrapolatory (0.065)	1: The triplets directly contradict the claim by showing that Pope Benedict XVI did indeed appoint someone (Rutilio del Riego Jáñez) who is associated with the Catholic Church, indicating that he did appoint significant individuals within the Church. 2: While the triplets indicate that God is equated with love, it does not directly address whether Pope Benedict XVI taught the importance of understanding God's redemptive love.	0.248
6	<sup>1</sup> Southwest Airlines has never operated any Boeing 737 models.	1: [{"Boeing 737 MAX", 'operator', 'Southwest Airlines'}, {'Boeing 737 #1491', 'operator', 'Southwest Airlines'}]	1: Contradictory (0.933)	1: The triplets directly contradict the claim by indicating that Southwest Airlines has operated both the Boeing 737 MAX and Boeing 737 #1491, which are specific models of the Boeing 737. This refutes the statement that Southwest Airlines has never operated any Boeing 737 models.	0.057

Table 2: Examples of claim-level attribution by the proposed method. The first column shows the numbered claims in the input text. Second column lists relevant triplets for each claim. Predictions and *Triplets Match Score (TMS)* are in the third column, while the rationale behind each prediction is in the fourth column. The *Knowledge Graph Attribution Score (KAS)* is shown in the last column. Model: *Solar-Instruct*.

may decompose the claims slightly differently, as multiple valid options are possible, the accuracy values may appear low even while the objective function is correctly executed. For instance, example 4 in Table 2 has been decomposed into three claims, but the first could arguably be further decomposed to verify whether Blagnac is in France, and whether it is a suburb of Toulouse. Controlling the precise manner of decomposition is challenging, and might necessitate an additional step before the prediction step, involving separate processing for each claim. However, this option could prove to be impractical, as the number of LLM queries could increase exponentially. While the overall accuracy may not fully reflect the models' performance due to the combined assessment of the two sub-tasks, focusing solely on the prediction task could offer better insights into how the models are performing in terms of categorization.

In Table 4, the second column indicates number of claims with text spans exactly matching the ground truth responses. Columns 3 to 6 present the accuracy, precision, recall, and F1 scores for these matching claims. The most performant model is *Solar-Instruct*, with 1052 exact matches out of 1675 claims in the test set. Across all models, the classification

scores in all metrics are above 98%, which clearly demonstrates that the models can reliably differentiate between the classes *attributable*, *extrapolatory*, and *contradictory*.

Table 2 showcases the claim-level attribution performed by our method. Each claim in the input text is numbered and color-coded to reflect its prediction: green for attributable, amber for extrapolatory, and red for contradictory. The examples are sorted in descending order by their KAS scores, which reflect the validity of the text. As expected, we observe more green at the top of the table and more amber and eventually red as we move down. Since the Wiki ecosystem is open-domain, we observe that the examples cover a wide range of topics, demonstrating that the method is adaptable to diverse inputs.

In the first example in 2, the input text is decomposed into two claims, both of which are attributable. The first claim is supported by a single triplet in the KG, while the second claim can be supported by combining two triplets. The second example presents more challenges for evaluation due to its complex sentence structure, but ClaimVer accurately identifies that the third claim regarding Batman's Back is neither supported nor refuted by the triplets, as indicated in the ra-



Model	Size	ROUGE-L	ROUGE-1	Avg. Acc.
Phi-2	2.7B	0.635	0.673	75.16%
Mistral-Instruct	7B	0.645	0.680	83.04%
Zephyr-Beta	7B	0.638	0.676	79.88%
Solar-Instruct	10.7B	<b>0.655</b>	<b>0.693</b>	<b>89.31%</b>
Llama2-Chat	13B	0.6395	0.677	79.41%

Table 3: ROUGE scores and average accuracies on the test set ( $n = 1,000$ ).

Model	#MC	Acc	Prec	Rec	F1
Phi-2	819	98.29	98.33	98.29	98.26
Mistral-Instruct	928	99.35	99.36	99.35	99.35
Zephyr-Beta	757	98.34	98.38	98.34	98.32
Solar-Instruct	<b>1052</b>	<b>99.80</b>	<b>99.81</b>	<b>99.80</b>	<b>99.80</b>
Llama2-Chat	869	99.53	99.54	99.53	99.53

Table 4: Scores on matching claims in the test set ( $n = 1675$ ). #MC: number of matching claims.

tionale. In the third example, we note that the first claim is predicted to be attributable with a high triplet match score of 0.942 since there is a triplet that clearly supports the location description of Crater Lake. However, as there is no information regarding the water characteristics, the second claim is categorized as extrapolatory. In the fourth example, the first claim alone requires three triplets combined as supporting evidence, illustrating the method’s ability to handle complex multi-hop paths within the KG. The second and third claims are predicted to be extrapolatory, since there are no triplets concerning Airbus’s market share, or its activities in Europe, as highlighted in the model’s rationale. It is noteworthy that the context provided in the third claim is crucial for the first claim to be comprehensible, demonstrating why individual claim evaluation may be suboptimal. Interestingly, in the fifth example, the method identifies a specific instance from the KG to refute a general claim, citing the appointment of Rutilio del Riego Jáñez. Similarly, in the sixth example, the method provides specific instances, quoting two distinct Boeing 737 models to demonstrate contradiction with a high triplet match score.

## 6 Discussion

The susceptibility of LLMs to generating factually incorrect statements is an alarming concern as LLM-powered services become increasingly popular for seeking advice and information. The democratization of generative models has also had adverse effects, such as increasing misinformation [Monteith *et al.*, 2024]. To arm end-users with the tools necessary to combat being misinformed, it is crucial to develop text-validation methods that are human-centric, and prioritize user engagement, understanding, and informativeness. We design our method with these principles in mind: we make predictions at the claim level, and identify text spans within the given text, that can be color-coded and presented to the user. The proposed method also generates easily comprehensible explanations along with the prediction and evidence, thus reducing the cognitive burden on the end-user, and making the process user-friendly.

The usability and evaluation of these systems should align with human needs and capabilities. Chatbots, such as ChatGPT [Achiam *et al.*, 2023], serve a wide array of tasks; therefore, the text validation method should be adaptable to various domains. While KGs like Wikidata [Vrandečić and Krötzsch, 2014] are considered open-domain, the implementation of more specialized KGs, along with corresponding routing algorithms may be necessary to support a broader range of topics. For instance, a common-sense KG [Hwang *et al.*, 2020] would be more useful in validating non-factoid answers that involve logic.

Furthermore, the maintenance efficiency of our approach aligns well with the need for sustainable, long-term AI solutions. In a world where information is constantly evolving, the ability to update and maintain AI systems with minimal effort is not just a convenience, but a necessity. This directly ties into the ethical implications of AI, where outdated or incorrect information can lead to harmful decisions. By leveraging existing, well-maintained KGs, we can ensure that AI systems remain accurate and relevant over time.

While there are several advantages associated with using KGs, we also acknowledge the presence of known issues, such as knowledge coverage and the efforts required to keep these sources up-to-date. For our solution, we assume that the KG is up-to-date and possesses adequate coverage. However, this may not always be the case, and thus the most suitable technique should be adopted after considering the specific requirements of a particular use case. Another point to consider is that the proposed method does not provide traditional citations to articles, although it may be possible to retrieve that information from the KG, if information source mapping has been properly maintained.

## 7 Conclusion

In this paper, we present ClaimVer, a framework that facilitates text verification and evidence attribution at the claim level by leveraging information present in KGs. We have prioritized human-centric design principles to make the framework more informative, intuitive, and user-friendly. Additionally, our methodology incorporates design choices that ensure open access, sustainability, and reliability.

ClaimVer presents several advantages, as outlined below:

1. *Human-centric design*: In addition to its primary functions of text verification and evidence attribution, the system generates considerable information conducive to user awareness. This information serves to educate users and enhance their trust in the automated system.
2. *Finer Granularity*: Perform validation at the claim level, enabling localization of hallucinations, or false claims.
3. *Enhanced Coverage*: Eliminate one-to-one mapping between input and reference text, allowing for layered interpretation, and handling of distributed information.
4. *Domain Adaptability*: Flexibility in adapting to new domains by switching to a more suitable KG.
5. *Maintenance Efficiency*: Simplified auditing and updating of the knowledge base, ensuring the data remains current and accurate.

## References

- [Achiam *et al.*, 2023] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [Aly *et al.*, 2021] Rami Aly, Zhijiang Guo, Michael Schlichtkrull, James Thorne, Andreas Vlachos, Christos Christodoulopoulos, Oana Cocarascu, and Arpit Mittal. Feverous: Fact extraction and verification over unstructured and structured information. *arXiv preprint arXiv:2106.05707*, 2021.
- [Bohnet *et al.*, 2022] Bernd Bohnet, Vinh Q Tran, Pat Verga, Roei Aharoni, Daniel Andor, Livio Baldini Soares, Massimiliano Ciaramita, Jacob Eisenstein, Kuzman Ganchev, Jonathan Herzig, et al. Attributed question answering: Evaluation and modeling for attributed large language models. *arXiv preprint arXiv:2212.08037*, 2022.
- [Chen *et al.*, 2019] Wenhui Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyu Zhou, and William Yang Wang. Tabfact: A large-scale dataset for table-based fact verification. *arXiv preprint arXiv:1909.02164*, 2019.
- [Chowdhery *et al.*, 2023] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113, 2023.
- [Conneau *et al.*, 2017] Alexis Conneau, Douwe Kiela, Holger Schwenk, Loic Barrault, and Antoine Bordes. Supervised learning of universal sentence representations from natural language inference data. *arXiv preprint arXiv:1705.02364*, 2017.
- [Dettmers *et al.*, 2024] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient fine-tuning of quantized llms. *Advances in Neural Information Processing Systems*, 36, 2024.
- [Dong and Smith, 2021] Rui Dong and David A Smith. Structural encoding and pre-training matter: Adapting bert for table-based fact verification. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2366–2375, 2021.
- [Gao *et al.*, 2023] Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vincent Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan, et al. Rarr: Researching and revising what language models say, using language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16477–16508, 2023.
- [Gerber, 2023] Emanuel Gerber. spacy module for linking text to wikidata items. <https://github.com/egerber/spaCy-entity-linker>, 2023. Accessed: 2024-02-26.
- [Gou *et al.*, 2021] Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. Knowledge distillation: A survey. *International Journal of Computer Vision*, 129:1789–1819, 2021.
- [Guan *et al.*, 2023] Jian Guan, Jesse Dodge, David Wadden, Minlie Huang, and Hao Peng. Language models hallucinate, but may excel at fact verification. *arXiv preprint arXiv:2310.14564*, 2023.
- [Gutiérrez and Patricio, 2023] Torres Gutiérrez and Cristóbal Patricio. Sistema visual para explorar subgrafos temáticos en wikidata. 2023.
- [Hu *et al.*, 2021] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [Hwang *et al.*, 2020] Jena D. Hwang, Chandra Bhagavatula, Ronan Le Bras, Jeff Da, Keisuke Sakaguchi, Antoine Bosselut, and Yejin Choi. Comet-atomic 2020: On symbolic and neural commonsense knowledge graphs. In *AAAI Conference on Artificial Intelligence*, 2020.
- [Javaheripi *et al.*, 2023] Mojan Javaheripi, Sébastien Bubeck, Marah Abdin, Jyoti Aneja, Sébastien Bubeck, Caio César Teodoro Mendes, Weizhu Chen, Allie Del Giorno, Ronen Eldan, Sivakanth Gopi, et al. Phi-2: The surprising power of small language models. *Microsoft Research Blog*, 2023.
- [Jiang *et al.*, 2023a] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- [Jiang *et al.*, 2023b] Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. Active retrieval augmented generation. *arXiv preprint arXiv:2305.06983*, 2023.
- [Kim *et al.*, 2023] Dahyun Kim, Chanjun Park, Sanghoon Kim, Wonsung Lee, Wonho Song, Yunsu Kim, Hyeonwoo Kim, Yungi Kim, Hyeonju Lee, Jihoo Kim, et al. Solar 10.7 b: Scaling large language models with simple yet effective depth up-scaling. *arXiv preprint arXiv:2312.15166*, 2023.
- [Lee *et al.*, 2017] Kenton Lee, Luheng He, Mike Lewis, and Luke Zettlemoyer. End-to-end neural coreference resolution. *arXiv preprint arXiv:1707.07045*, 2017.
- [Lee *et al.*, 2020] Nayeon Lee, Belinda Z Li, Sinong Wang, Wen-tau Yih, Hao Ma, and Madian Khabsa. Language models as fact checkers? *arXiv preprint arXiv:2006.04102*, 2020.
- [Lin, 2004] Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics.
- [Mallen *et al.*, 2022] Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Hannaneh Hajishirzi, and Daniel



- Khashabi. When not to trust language models: Investigating effectiveness and limitations of parametric and non-parametric memories. *arXiv preprint arXiv:2212.10511*, 2022.
- [Min *et al.*, 2023] Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. Factscore: Fine-grained atomic evaluation of factual precision in long form text generation, 2023.
- [Monteith *et al.*, 2024] Scott Monteith, Tasha Glenn, John R Geddes, Peter C Whybrow, Eric Achtyes, and Michael Bauer. Artificial intelligence and increasing misinformation. *The British Journal of Psychiatry*, 224(2):33–35, 2024.
- [Rashkin *et al.*, 2023] Hannah Rashkin, Vitaly Nikolaev, Matthew Lamm, Lora Aroyo, Michael Collins, Dipanjan Das, Slav Petrov, Gaurav Singh Tomar, Iulia Turc, and David Reitter. Measuring attribution in natural language generation models. *Computational Linguistics*, pages 1–64, 2023.
- [Rechkemmer and Yin, 2022] Amy Rechkemmer and Ming Yin. When confidence meets accuracy: Exploring the effects of multiple performance indicators on trust in machine learning models. In *Proceedings of the 2022 chi conference on human factors in computing systems*, pages 1–14, 2022.
- [Russell *et al.*, 2015] Ginny Russell, Stephan Collishaw, Jean Golding, Susan E Kelly, and Tamsin Ford. Changes in diagnosis rates and behavioural traits of autism spectrum disorder over time. *BJPsych open*, 1(2):110–115, 2015.
- [Schick and Schütze, 2020] Timo Schick and Hinrich Schütze. Exploiting cloze questions for few shot text classification and natural language inference. *arXiv preprint arXiv:2001.07676*, 2020.
- [Sciavolino *et al.*, 2021] Christopher Sciavolino, Zexuan Zhong, Jinhyuk Lee, and Danqi Chen. Simple entity-centric questions challenge dense retrievers. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6138–6148, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.
- [Shin, 2021] Donghee Shin. The effects of explainability and causability on perception, trust, and acceptance: Implications for explainable ai. *International Journal of Human-Computer Studies*, 146:102551, 2021.
- [Thorne *et al.*, 2018] James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. Fever: a large-scale dataset for fact extraction and verification. *arXiv preprint arXiv:1803.05355*, 2018.
- [Thorne *et al.*, 2019] James Thorne, Andreas Vlachos, Oana Cocarascu, Christos Christodoulopoulos, and Arpit Mittal. The FEVER2.0 shared task. In James Thorne, Andreas Vlachos, Oana Cocarascu, Christos Christodoulopoulos, and Arpit Mittal, editors, *Proceedings of the Second Workshop on Fact Extraction and VERification (FEVER)*, pages 1–6, Hong Kong, China, November 2019. Association for Computational Linguistics.
- [Touvron *et al.*, 2023] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [Tunstall *et al.*, 2023] Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, et al. Zephyr: Direct distillation of lm alignment. *arXiv preprint arXiv:2310.16944*, 2023.
- [Vrandečić and Krötzsch, 2014] Denny Vrandečić and Markus Krötzsch. Wikidata: a free collaborative knowledgebase. *Communications of the ACM*, 57(10):78–85, 2014.
- [Wang *et al.*, 2021] Nancy XR Wang, Diwakar Mahajan, Marina Danilevsky, and Sara Rosenthal. Semeval-2021 task 9: Fact verification and evidence finding for tabular data in scientific documents (sem-tab-facts). *arXiv preprint arXiv:2105.13995*, 2021.
- [Weitz *et al.*, 2019] Katharina Weitz, Dominik Schiller, Ruben Schlagowski, Tobias Huber, and Elisabeth André. ”do you trust me?” increasing user-trust by integrating virtual agents in explainable ai interaction design. In *Proceedings of the 19th ACM International Conference on Intelligent Virtual Agents*, pages 7–9, 2019.
- [Yang *et al.*, 2015] Yi Yang, Wen-tau Yih, and Christopher Meek. WikiQA: A challenge dataset for open-domain question answering. In Lluís Màrquez, Chris Callison-Burch, and Jian Su, editors, *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 2013–2018, Lisbon, Portugal, September 2015. Association for Computational Linguistics.
- [Yang *et al.*, 2018] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*, 2018.
- [Yao *et al.*, 2022] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.
- [Yue *et al.*, 2023] Xiang Yue, Boshi Wang, Kai Zhang, Zirui Chen, Yu Su, and Huan Sun. Automatic evaluation of attribution by large language models. *arXiv preprint arXiv:2305.06311*, 2023.
- [Zhou *et al.*, 2023] Chunting Zhou, Pengfei Liu, Puxin Xu, Srinu Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *arXiv preprint arXiv:2305.11206*, 2023.