

Factcheck-GPT: End-to-End Fine-Grained Document-Level Fact-Checking and Correction of LLM Output

Anonymous ACL submission

Abstract

The increased use of large language models (LLMs) across a variety of real-world applications calls for mechanisms to verify the factual accuracy of their outputs. In this work, we present a holistic end-to-end solution for annotating the factuality of LLM-generated responses, which encompasses a multi-stage annotation scheme designed to yield detailed labels concerning the verifiability and factual inconsistencies found in LLM outputs. We build an annotation tool to speed up the labelling procedure and ease the workload of raters. It allows flexible incorporation of automatic results in any stage, e.g. automatically-retrieved evidence. We further construct an open-domain document-level factuality benchmark in three-level granularity: claim, sentence and document. Preliminary experiments show that FacTool, FactScore and Perplexity.ai are struggling to identify false claims, with the best F1=0.63 by GPT-4. Annotation tool, benchmark and code are available at URL withheld.

1 Introduction

Large language models (LLMs) have demonstrated impressive capabilities in terms of generating naturally-sounding answers over a broad range of human inquiries (OpenAI, 2023). However, ChatGPT and other text generation models still frequently produce content that deviates from real-world facts (Menick et al., 2022; Bang et al., 2023; Borji, 2023; Guiven, 2023; Augenstein et al., 2023). This degrades the system performance and undermines its reliability, representing a significant bottleneck in their deployment especially for high-stakes applications, e.g., clinical, legal and financial settings (Chuang et al., 2023).

Before LLMs, most prior work investigate human-written sentence-level fact-checking (e.g. FEVER-based series) (Guo et al., 2022; Thorne et al., 2018), and detect hallucinations of conditional text generation for specific tasks, such as

abstract summarisation, dialogue generation and machine translation (Ji et al., 2023). They are either highly task-specific with references or focusing on short statements, making it hard to directly apply to open-domain generations of arbitrary length in LLM settings. Though some work identify fake news and assess the reliability of a website, they take the entire document and website features as a whole, instead of verifying individual component claims (Zellers et al., 2019; Baly et al., 2020).

While for outputs of LLMs, not only the overall factuality, we also expect every single claim of the generation to be checked and attributed by one or more pieces of evidence from reliable sources, such as encyclopedia or web articles (Gao et al., 2022). This requires understanding context, extracting context-independent claims from a document, identifying what to check (check-worthiness), verifying and correcting claims (and the document).

There have been an increasing number of approaches for LLM generated text fact-checking. However, they were not end-to-end, and only looked into some (subset) of these problems. For example, most studies only detect whether the response contain hallucinations, such as SelfCheck-GPT (Manakul et al., 2023), FACTOR (Muhlgay et al., 2023) and FactScore (Min et al., 2023), they do not correct the errors. FacTool (Chern et al., 2023) makes corrections, but requires human-decomposed claims as the input of the system. RARR (Gao et al., 2022) and CoVe (Dhuliawala et al., 2023) conduct both detection and correction, but in coarse granularity — editing over the whole document. Compared to claim-level verification, it cannot spot false spans precisely and tends to result in poor preservation of the original input.

In contrast, we propose a comprehensive end-to-end fine-grained solution, covering eight steps that may occur in the automatic detection and correction of factual errors, shown in Figure 1. Our work is applicable for both human-written text and

the outputs of LLMs, with an emphasis on long documents, while it is motivated by LLMs.

In sum, the main contribution of this work can be summarised into four folds:

- We propose a holistic end-to-end solution for automatically detecting and correcting factual errors for the outputs of generative LLMs.
- We build an annotation tool for efficient construction of factuality benchmark. It allows flexible customisation of annotations for any component subtasks, and supports semi-auto annotation by incorporating results from automatic methods, such as automatically-decomposed claims and automatically-retrieved evidence.
- Using the tool, we construct the first document level claim-based fact-checking benchmark of LLMs, in terms of both detection and revision, facilitating the evaluation and analysis of existing and future fact checkers.
- We evaluate the popular checkers FacTool and FactScore against the annotated examples, and find large headroom for improvement in LLM fact-checking. We open sourced the annotation tool, data and code at URL withheld.

2 Framework and Subtasks

We frame the automated detection and correction of factual errors for outputs of LLMs into eight subtasks: (1) decomposition; (2) decontextualisation; (3) checkworthiness identification; (4) evidence retrieval and collection; (5) stance detection; (6) correction determination; (7) claim correction and (8) final response revision. Figure 1 presents the overview of the whole procedure, coupled with an example flowing through each subtask.

(1) Decompose Given a response R generated by a LLM, it is infeasible to fact-check the whole document at once, especially when it is long. The first step is to decompose R into a set of context-independent atomic statements, with no information lost or distorted in this process. Decomposed statements should be checkable independently without preceding and following context.¹

(2) Decontextualise Sentences in a response might be context-dependent, with discourse and coreference relations existing between statements.

¹Statements are assumed to be checkable if relevant documents exist in publicly-available data sources.

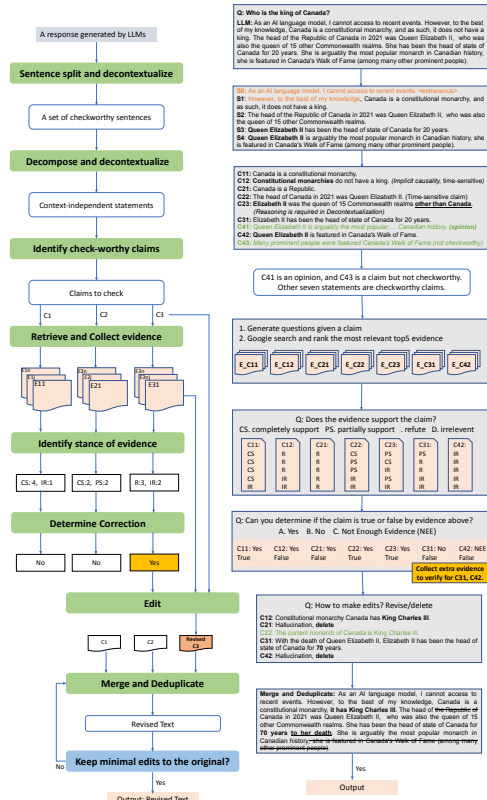


Figure 1: **Left:** Fact-checking pipeline for a response generated by LLMs. **Right:** An example workflow.

For example, it is invalid to check statement *It does not have a king* before replacing “It” with “Canada” or “Constitutional monarchy” (see Figure 1). In addition to coreference relation, for the sentence S2, it is not reasonable to check the claim *Queen Elizabeth II is also the queen of 15 other Commonwealth realms*. Instead, the claim should be reframed to *Queen Elizabeth II is the queen of 16 Commonwealth realms (including Canada)* or *Queen Elizabeth II was the queen of 15 Commonwealth realms other than Canada*.

The concept of “context-independent” is straightforward, while the notion of “atomic” is subjective and ambiguous. This poses challenges: how to determine the granularity of *atomic claim*? when and where to break down a response? For example, S1: *Canada is a constitutional monarchy, and as such, it does not have a king*, can be fact-checked as one statement, or be decomposed into two claims: *Canada is a constitutional monarchy* and *Canada does not have a king*. In our work, we first split a document into sentences, and then from sentence to claims, with each claim containing only one property or fact to verdict.

(3) Identify Checkworthy Claims Not all statements in a response require fact-checking, such as subjective opinions and actual commonsense as obvious as *sun rises from the east*. Each statement in this framework will be identified whether it is checkworthy or not. However, checkworthiness is subjective to determine. Hassan et al. (2015) defined checkworthy claims as those for which the general public would be interested in knowing the truth. In the context of fact-checking LLMs outputs, we assume users who ask LLMs questions are interested in knowing the truth of all factual claims in the corresponding answer.

We specifically classify a statement into four categories: factual claim, opinion, not a claim (e.g. questions, exclamations, imperatives), and others (e.g. *As a language model, I cannot access personal information*). Afterwards, a set of check-worthy factual claims needs to be verified by retrieving and collecting evidence.

It is worthwhile to note that regarding checkworthiness, we not only take account of objective fact against subjective judgement, other aspects such as the role (importance) of the claim to the response may also be a crucial criteria for its checkworthiness. For example, the sentence S1 in Figure 1 needs more attention than the last sentence S4. We label the importance level of both decomposed sentences and claims by labels: *most important*, *intermediate* and *less important*.

(4) Retrieve and Collect Evidence Evidence can be retrieved by a search engine like Google, or deep retrieval from a closed document collection such as Wikipedia, or using the parametric knowledge of a LLM. Search queries can be questions covering different aspects of the claim, entities in the claim or even claim itself.

(5) Identify Stance of Evidence With retrieved evidence for a claim, how to identify the stance of the evidence against the claim. RARR (Gao et al., 2022) achieved this by assessing whether answers depending on the evidence and the claim are same or not, in terms of a given query. If they are same, then the evidence supports the claim, otherwise refutes. Previous work also employs natural language inference (NLI) model to classify whether the claim can be entailed by evidence, or is controversial against evidence, or is irrelevant.

However, some evidence may neither refute nor fully support a claim. This mainly results from

the fact that it is always possible that the evidence supports part of the claim. For example, for the claim *Elon Musk is the founder, CEO and chief engineer of SpaceX*, evidence *Elon Musk is the CEO of SpaceX, Tesla and Twitter* falls into this category. The evidence supports the factual statement of *Elon Musk is the CEO of SpaceX*, but it does not provide information regarding whether *Elon Musk is the founder and chief engineer of SpaceX*.

Therefore, we incorporate *partially support* in addition to *support*, *refute* and *irrelevant*. Concretely, *support* means that the evidence entails the claim. *Partial support* refers to the scenario where part of the information presented in a claim appears in the evidence. *Refute* means that the evidence mention the same event as the claim, but a clear opposite fact contrasting to a part or the whole facts presented in a claim. *Irrelevant* refers to the situation that the evidence does not mention anything about the fact described in the claim, such that it neither supports nor refutes the claim.

Sometimes, it is ambiguous to distinguish between *refute* and *irrelevant*. We highlight that the evidence shows a clear opposite fact under *refute* stance, while the evidence does not include relevant facts mentioned in the claim under *irrelevant*. Overall, we aim to identify whether a piece of evidence supports, partially supports, refutes the claim, or is irrelevant to the claim in this subtask.

(6) Determine Correction Given a claim, there will be more than one pieces of related evidence. Most of the time, they hold consistent stances except for irrelevance, but sometimes, some support, some partially support while some refute (see Figure 8). How to aggregate conflicting stances and further decide how to make correction for the claim is an open question. In practice, when evidence paragraphs conflict with each other, we will take the reliability of the evidence source into consideration, meanwhile, retrieve extra information to judge which one is more dependable.

A label often used is *not-enough-evidence* if there is insufficient information to make the veracity prediction, e.g., all retrieved evidence is irrelevant or intricate contradictory evidence (Atanasova et al., 2022). So we set three labels in terms of factuality: true, false and not-enough-evidence.

(7) Edit Claims With the principle that revised claims should preserve the text’s original intent and style. without adding or changing unnecessary

additional information, we include edit operations: delete the whole claim, replace X to Y, delete X, where X and Y are meta information in a claim.

(8) Revise Response After revision, we merge statements in original order, including non-check-worthy statements, true claims and correctly-revised claims. Finally, we delete reduplicative content if applicable, and output a factually-correct and fluent response.

Are solutions under this framework effective to identify the factuality of LLM-generated responses?

To answer this question, a human-annotated factuality benchmark — LLM responses with true/false labels for component claims and the revision, is needed to evaluate whether the automatic fact-checkers can recognise and rectify false statements.

FELM (Chen et al., 2023) collected 184 world-knowledge responses, but only annotated sentence-level *true or false* labels without correction. Their findings show that detection performance tends to be improved when utilising claim-based segmentation methods compared with sentences. To this end, we annotate a claim-based document-level fact-checking dataset in Section 3.

3 Dataset Construction

We annotate a dataset serving for few-shot demonstration examples and a benchmark evaluating the effectiveness of approaches for LLM fact-checking subtasks or the whole pipeline.

3.1 Data Collection

What kind of LLM generations are we most concerned about? In the context of detecting and correcting factual errors, we focus on generations in which the majority of statements are objective facts rather than subjective opinions whose veracity is not checkable. Additionally, we are more interested in questions where LLMs are prone to hallucinate or produce factual errors in responses.

The whole annotation process is extremely time-consuming, about 15-30 minutes for an instance even if with intermediate results from automatic methods to ease the procedure. This requests us to cherry-pick examples that highly satisfy two criteria — **fact-intensive** and **factually-false**.

Sources To this end, we start from hallucinations posted by ChatGPT users in Twitter and further collect data by in-house brainstorming with preliminary verification, resulting in 45 examples.

Inspired by collecting data from brainstormed questions, we decide to select more questions from dolly-15k, which is brainstormed by thousands of Databricks employees with eight categories. We select 563 examples from closed QA and 528 from open QA by ChatGPT response length and the semantic similarity with gold answers, thus 1,136 (question, response) pairs in total with the 45 from the first source (see more in Appendix B.1).

Data Selection We select factually-false responses by estimating the percentage of incorrect claims in a response with four steps.

Sentence and claim split: given the whole response as the context and the first sentence (initialised by NLTK tokenizer), we instruct ChatGPT by three demonstration examples to guide it first break the input sentence into independent atomic claims, and then continue the decomposition of the next sentence until the end of the response (see the prompt in Appendix B.5).

Evidence collection: given an atomic claim, we first prompt ChatGPT to generate search queries for the claim, and then Google search engine is used to get relevant web pages. Retrieved documents are split into passages by sliding windows, and a re-ranker combining lexical and semantic similarity is used to identify the most relevant passages for the given query, in which Sentence-BERT (Reimers and Gurevych, 2019) serves for semantic embeddings. We aggregate evidence for all queries and select the top-5 evidences per atomic claim.

FactScore calculation: FactScore is an automatic factuality metric, measuring the percentage of atomic claims supported by knowledge sources in a generation (Min et al., 2023). We use the gathered evidences as input, along with the claim, and an instruction-tuned LLM as the verifier to verdict.

Example selection: we keep all 45 pairs from the first source and dolly examples whose FactScore is less than 0.2, resulting in 33 closed QA pairs and 37 open questions, in total of 115 examples. We remove a similar question, and four questions where the LLM did not provide helpful answers due to its inherent disability to access real-time data, eventually annotating 110 examples in our first annotation stage. More cases of multiple languages would be annotated in the next stage.

3.2 Annotation

As many studies illustrated, annotating a LLM factuality benchmark is a highly challenging

task (Chen et al., 2023; Li et al., 2023). Our preliminary trials in which authors manually annotate labels of all subtasks empirically confirm the pain.

Preliminary Trial Take-away Manually annotating the whole process and typing results into a *json* file expose three major difficulties: (1) retrieving supportive or contradictory evidence takes time and demands annotator’s strong skills in searching relevant and filtering out unreliable information, especially for non-common knowledge (e.g. *most popular bottled water brand in Israel*); (2) lengthy responses require good reading comprehension ability and patience; (3) certain domains such as gene and astronomy require domain knowledge, otherwise it is hard to search valid evidence and determine whether it is true or false.

Taking the factors mentioned above into consideration, we design and build an annotation tool to support the efficient construction of LLM factuality benchmark. Annotators can edit and assign labels based on intermediate outputs of automatic methods, click buttons instead of typing to copy-paste text, select, and download annotated results.

Annotation Tool includes all subtasks and supports semi-auto annotation by incorporating the results of automatic methods, such as automatically-decomposed claims and automatically-retrieved evidence, to ease the annotation process and reduce the workload (see interfaces in Appendix G).

We perform the whole annotation in three steps: (1) decomposition, decontextualization and check-worthiness detection; (2) evidence stance identification and claim correction; (3) claim merge and deduplication.

Between the step (1) and (2), we incorporate an automatic evidence retrieval system to provide annotators a set of most relevant snippets of text as evidence for each checkworthy atomic claim, generally five pieces. They are selected and ranked by semantic relevance degree against the claim throughout a large number of documents, similar to the evidence collection in data selection above.

Then, annotators determine the stance of each piece of evidence against the claim. With the evidence from automatic system, if annotators still cannot determine the factuality of a claim, they are requested to collect relevant evidence manually. This to some extent alleviates the system bias.

Quality Control To guarantee the annotation quality, instead of employing crowd-sourcing an-

	example	sent	cw_sent	claim	cw_claim	evid
size	94	311	277	678	661	3,305

Table 1: Statistics of the dataset. *cw_sent*=checkworthy sentences, *cw_claim*=checkworthy claims, *evid*=the total pieces of evidence, five for each *cw_claim*.

notators, we perform an in-house annotation by ten annotators who are Master’s and PhD students, postdocs and professors and familiar with the task of automatic fact-checking.

Two annotators as a group are responsible for 22 examples. For each step, annotators first independently finish individual annotations, and then consolidate their results with the group partner. In consolidation, partners discuss their disagreements until reaching consensus. For cases where it is hard to reach agreement even with the participation of the third rater, we discard it. Three steps are rigorously conducted serially. Annotators start the second step only after they finish the consolidation of the first step. Collecting evidence and judging stances is the most time- and patience-consuming step. To ensure the quality, we incorporate the third rater when consolidating the second-step annotations in case of unintentional mistakes.

3.3 Data Analysis

During annotation, we remove another 16 examples (see details in Appendix B.4), resulting in 94 instances. Statistics shows in Table 1.

Statistics 277 sentences contain factual statements among 311. There are 678 atomic claims, where 661 claims are checkworthy, 16 are opinions and one is *not-a-claim*. For each checkworthy claim, five pieces of evidence are collected, resulting in 3,305 (claim, evidence, stance) triplets.

How many examples are factually incorrect? 61 examples contain factual errors and 31 examples are factually correct. Specifically, 53 examples contain false claims, and 19 examples contain claims in which annotators cannot verify the statement due to insufficient evidence despite the manual search. Generally, one example contain 0-5 false claims, and six answers have >5 incorrect claims.

16 sentences among 331 are deleted. 12 are total hallucinations, e.g., *Trump was the second black president*. 4 sentences are removed due to lacking enough evidence to support its factual correctness.

Table 2 shows that more incorrect responses appear in in-house collected questions, followed by

Source	In-house	Closed-QA	Open-QA	All
Collected	45	33	35	115
Annotated	39	30	25	94
False	38	16	8	61

Table 2: False responses over three question sources.

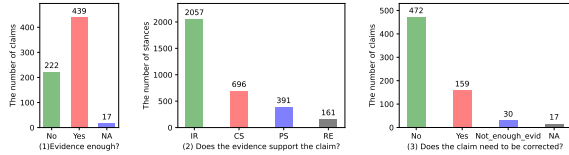


Figure 2: **Claim analysis:** (1) whether raters can determine the factuality of a claim depending on the automatically-collected evidence (*Yes/No*); (2) does the evidence support the claim (*CP*:completely support, *PS*: partially support, *RE*: refute, *IR*: irrelevant); (3) does the claim need to be corrected. NA (17) refers to 16 opinion-claims + 1 *not-a-claim*.

dolly closed questions that require knowledge to obtain a unique correct answer. Less errors occur in dolly open questions, in which correct answers are not unique, e.g., *How do you play an E major chord on a guitar?* It has diverse correct answers requiring more general knowledge.

Claims in 678 claims, 419 and 227 are labelled as the most and intermediate important claims, only 32 fall into *not-important*, indicating that users concern almost the whole response given their importance. We analyse annotations of 661 checkworthy claims from two perspectives.

Can raters determine the factuality of a claim depending on the automatically-collected evidence? For 439 claims, annotators can determine *true or false* with automatic evidence, while 222 claims (one-third) need further manual retrieval to make judgements.

We wonder what makes most collected evidence invalid for these 222 claims (see Figure 6). Almost half of them (97) are either factually-incorrect claims (76) or undetermined claims without sufficient evidence in spite of manual retrieval (21). The remaining 125 true claims basically fall into domain knowledge and information that is less known by the external people given a country, region, company or an individual.

This may suggest the ineffectiveness of the automatic evidence retrieval methods, particularly collecting evidence conditioned on false premises (claims). However, it may also reply that not all

facts have been presented by textual descriptions directly. Some facts are unknown by the public, and some require connecting and reasoning knowledge from multiple sources, e.g., *did Aristotle use a laptop?* (Geva et al., 2021).

How many claims need to be corrected? In Figure 2, about a quarter (159/661) claims are factually incorrect and need to be corrected. 30 claims are undetermined due to inadequate related information and knowledge even with manual retrieval. It is hard to obtain reliable related information about these cases by searching publicly-available sources. They involve expert-level knowledge (e.g., gene, water memory, black hole) and minor details of an individual, organisation or country (personal awards and preferences, revenue of a company), which are only known by a small group of people, such as domain experts or internal individuals who are familiar with the event.

Original vs. revised responses We quantify the difference between the original responses and the human-revised responses over the 61 false responses, showing that the normalised edit distance is 0.354, word overlap is 0.715, while semantically, BERTScore-F1 is 0.955 and cosine similarity based on SimCSE (Roberta-large) is 0.912.

Summary The dataset consists of 94 ChatGPT (prompt, response) pairs. Each sample has detailed labels concerning the verification: elements of decontextualised sentences, atomic claims, the importance degree of sentence.claim to the response, five pieces of evidence for a claim, the relationship between a claim and evidence, factual label (*true or false*) and revised version of claims, sentences and the response.

4 Baseline

The whole pipeline involves several steps, we first compare the automatic and human-annotated decomposition of atomic claim, and then evaluate five subtasks: (1) identify whether the sentence contains factual statement; (2) detect the checkworthiness of a claim by categories of *factual*, *opinion*, *not a claim* and *other*; (3) judge the stance of a given evidence against a claim, whether it *supports*, *partially supports*, *refutes* or is *irrelevant* to the claim; (4) determine whether a claim is factually true or false, give a claim without “gold evidence”, if false, revise it into a correct one; (5) edit a list of originally-true or revised claims into a new re-

sponse, given the original response, to correct the factual errors while preserving the linguistic features and style of the original.

Other steps are excluded because they are either relatively easy for current techniques (e.g., splitting a document into sentences), or results of automatic approaches have been compared against human annotations in data analysis, such as the relevance or quality of the automatically-retrieved evidence.

4.1 Automatic vs. Manual Decomposition

For 66/277 check-worthy sentences, the number of decomposed atomic claims are different between automatic breaking-down by ChatGPT and manual annotations. Amongst, more claims decomposed by the automatic method than human for 48 sentences, and less claims for 18 sentences. This exhibits that human annotators add extra claims on only a small number of sentences. For most cases, automatic approach decomposes sentences into equal number of claims or even more fine-grained than humans.

For the rest 211 sentences, human and ChatGPT decompose the sentence into the same number of claims, 521 claims are involved. This enables pairwise claim comparison between the human annotation and automatic method. We calculate the lexical similarity and distance: normalised edit distance=0.11, n-gram distance=0.11 and word overlap=0.88, demonstrating high agreement between human annotation and ChatGPT decomposition.

4.2 Checkworthiness

We apply ChatGPT to identify if decomposed sentences and claims are verifiable objective facts or statements containing personal opinions.

Subtask 1 and 2 We identify whether a sentence contains factual statement by a binary label (*yes* or *no*), and whether a claim is check-worthy by four labels (*factual claim*, *opinion*, *not-a-claim* and *other*). The accuracy for subtask 1 by majority guess (always checkworthy) will be $277/311=0.891$ and the baseline for subtask 2: claim classification is $661/678 = 0.975$. They are superior to using the prompt based on ChatGPT: the accuracy is 0.814 and 0.932 respectively. However, this is mainly attributed to the extremely-unbalanced data. Practically, our aim is to make distinctions. It’s critical to consider recall: ChatGPT is much better than the majority guess (see Table 3).

Confusion matrix in Figure 9 shows that 46

Task	Method	Acc	Prec	Recall	F1-macro
1	Always-checkworthy	0.891	0.445	0.500	0.471
1	ChatGPT	0.814	0.637	0.740	0.660
2	Always-checkworthy	0.975	0.325	0.333	0.329
2	ChatGPT	0.932	0.314	0.534	0.319

Table 3: **Checkworthiness** detection by majority guess: Always-checkworthy vs. ChatGPT zero-shot prompt. *average*=“macro” is used in precision (Pred), recall and F1 calculation.

Method	Acc	Prec	Recall	F1-macro
Four-label space				
Random guess	0.255	0.258	0.264	0.215
ChatGPT-zeroshot	0.365	0.402	0.439	0.332
Three-label space				
ChatGPT-zeroshot	0.567	0.506	0.588	0.483
LLaMA2-zeroshot	0.401	0.407	0.384	0.299
RoBERTa-large-mnli	0.607	0.536	0.609	0.512

Table 4: **Stance** detection by ChatGPT and LLaMA2 zero-shot prompt. Three-label space merges complete and partial support into one.

checkworthy sentences are identified as non-checkworthy, accounting for 15%. Factual claims could be recognised into any of the four labels, and real opinions tend to be identified as factual claims, even more than the opinion.

4.3 Verification

Subtask 3 classifies whether the evidence fully supports, partly supports, refutes or is irrelevant to the claim, given a (*claim*, *evidence*) pair. We use zero-shot prompting based on ChatGPT and LLaMA2 (7B), and find that LLaMA2 barely predicts partial support, so we merge *complete support* and *partial support* into a single label *support*. As results shown in Table 4, three labels are easier for models to predict with higher accuracy, but its absolute F1-score is still less than 0.5, revealing the challenges to distinguish the relationship between claim and evidence by LLM in-context learning, especially on the label of *refute*. Both LLaMA2 and ChatGPT show around-0.1 F1 (see Table 10). We further use a fine-tuned NLI model (*RoBERTa-large-mnli*) to predict the stance, where entailment, contradiction and neutral correspond to labels of support, refute and irrelevant respectively. It performs better than zero-shot ChatGPT, mainly being superior to predicting the label of *support*.

Subtask 4 determines whether the claim is true or false leveraging evidences retrieved from external knowledge sources. We evaluate the verifica-

Verifier	Source	Label = True			Label = False		
		Prec	Recall	F1	Prec	Recall	F1
Random	NA	0.79	0.43	0.56	0.18	0.52	0.27
Always True	NA	0.81	1.00	0.88	0.00	0.00	0.00
Always False	NA	0.00	0.00	0.00	0.19	1.00	0.33
Inst-LLAMA	Wiki	0.87	0.74	0.80	0.34	0.56	0.42
Inst-LLAMA	Web	0.88	0.80	0.84	0.40	0.56	0.47
GPT-3.5-Turbo	Wiki	0.87	0.67	0.76	0.31	0.60	0.41
GPT-3.5-Turbo	Web	0.89	0.74	0.81	0.37	0.62	0.46
Perplexity.ai	Web	0.93	0.73	0.83	0.40	0.76	0.53
AutoGPT-4	Web	0.90	0.71	0.79	0.52	0.80	0.63

Table 5: **Verification results** on our benchmark: judge whether a claim is factually true or false with external knowledge (Wikipedia or Web articles) as evidence.

tion methods used in FactScore (Min et al., 2023) and FacTool (Chern et al., 2023), with varying evidence sources: Wikipedia (September 2023 dump) and web articles searched by Google.

Table 5 shows that false claims tend to be identified less accurately than true claims across all approaches, implying that it is more difficult to detect factual errors than the correct statements. GPT-4 performs the best on false claims with F1=0.63, and then Perplexity.ai verifier by 0.53, followed by Instruction-LLaMA with web articles as evidence (F1=0.47/0.84), and verifying using GPT-3.5-Turbo exhibits slight declines. This reveals that current mainstreaming SOTA fact-checkers still have large room to improve on verification, particularly on false claims. Performance using Wikipedia as source is inferior to using web articles, this is largely limited by the knowledge coverage of Wikipedia, esp. on open-domain benchmarks.

4.4 Revision

Subtask 5 Given the original false response, a list of revised true claims, ChatGPT and GPT-4 are prompted to revise the responses with/without the question, resulting in four revised responses.

Which revised response is better? We evaluate by human and the intrinsic metrics. BERTScore measures semantic preservation between gold reference answers and the edit-distance measures style preservation between original responses.

In human evaluation, we use the criteria: whether the revised response (1) contain factual errors? (2) keep the style feature of the original response as much as possible? (3) is it natural, coherent and smooth as an answer? Criteria (1) is the most important, followed by (2) and (3). For instance, only *A* and *B* are factually correct, while *A* preserves more of the original response, thus *A* is better. If some responses are totally the same, raters can choose more than one. We collect 66

Prompt	model	Edit-dis↓	WO↑	BS-F1↑	STS↑	Human
no-ques	ChatGPT	0.207	0.864	0.953	0.937	10
no-ques	GPT-4	0.275	0.789	0.954	0.931	28
with-ques	ChatGPT	0.222	0.853	0.956	0.941	13
with-ques	GPT-4	0.286	0.776	0.953	0.935	15

Table 6: **Revision evaluation** by intrinsic metrics and human (how many responses are preferred). Edit distance (**Edit-dis**) and word overlap (**WO**) between revised and the original responses. BERTScore (**BS-F1**) and semantic textual similarity (**STS**) based on SimCSE between the revised responses and human annotations.

preference labels for 61 examples.

In case of personal preference bias from one or two raters, six raters are invited to choose their preferred response and provide a brief reason. We also shuffled four revisions and show by “revision_x” (x=0,1,2,3), masking the real setting name to avoid possible inherent biases.

In Table 6, intrinsic metric results show that responses revised by ChatGPT (GPT-3.5-turbo) is better than GPT-4, which is against our experience and observation (see examples in Appendix F). Human assessment exhibits that 43 GPT-4 responses are preferred by raters and 23 from ChatGPT. Human is more satisfied with revisions without questions 38 vs. 28. This somewhat reflects the ineffectiveness of intrinsic evaluation metrics.

Take-Away ChatGPT shows promising results in atomic-claim decomposition, but low F1-score in checkworthiness detection. Also, verification remains challenging, especially identifying false claims even if harnessing external knowledge. GPT-4 can generate sounding revised responses based on true statements. It’s still an open-question in terms of how to evaluate the quality of revised responses by effective intrinsic metrics.

5 Conclusion

We propose a holistic fact-checking framework for open-domain generations of LLMs, and annotate a factuality benchmark to facilitate the evaluation of automatic fact-checkers. Experiments show that current verifiers are struggling to identify open-domain false claims with the best F1=0.63 even if using external knowledge. Additionally, intrinsic metrics based on edit distance and semantic similarity are ineffective to evaluate the edited responses against true evidence and the original response, misaligning with human preferences. It is worth more exploration in the future work.

679 Limitations

680 Three major limitations are identified in this work:

681 **Small-scale dataset** It consists of only 94 (ques-
682 tion, response) pairs, we plan to scale up the dataset
683 in English, Chinese and Arabic in the future work.

684 **Inter-claim dependencies** This reflects at three
685 challenges. First, current approaches including our
686 solution are unable to checking the overall logic
687 correctness of a procedure, such as how to cook,
688 whether some steps are out of order. Second, if
689 the first claim is invalidated, maybe the entire text
690 needs to be deleted. Third, it is hard to decontextu-
691 alize implicit claims, e.g., “other 15 realms”, which
692 means there are 16 realms.

693 **Quality of evidence** More than half automati-
694 cally retrieved evidences are irrelevant. Improving
695 the relevance of retrieved evidence is critical to the
696 accuracy of fact-checking.

697 Ethics and Broader Impact

698 We identify two major risks of the framework and
699 benchmark:

700 **Biases:** The automatic atomic-claim decomposi-
701 tion and evidence retrieval systems incorporated in
702 the fact-checking annotations may introduce biases,
703 which can affect the annotation results.

704 Besides, the dataset does not cover all types of
705 claims. Limited scope and coverage may be more
706 effective in certain domains, possibly leading to
707 inaccurate or unfair assessments in certain domains
708 for automatic fact-checkers. The responses gener-
709 ated by LLMs might also inherit some biases
710 present in the involved LLMs.

711 **The cost of making an error:** The goal of fact-
712 checking is to improve the reliability of the LLM
713 outputs, If post-hoc fact-checking methods under
714 this framework always make errors, practitioners
715 may lose faith in the accuracy of the fact-checking
716 results, which can affect efforts to maintain public
717 trust in fact-checking systems.

718 **Broader impact:** The proposed framework is not
719 limited to checking the output of LLMs; it is appli-
720 cable to checking any kind of document, including
721 human-written.

References

- Pepa Atanasova, Jakob Grue Simonsen, Christina Li-
oma, and Isabelle Augenstein. 2020. [Generating fact
checking explanations](#). In *Proceedings of the 58th
Annual Meeting of the Association for Computational
Linguistics*, pages 7352–7364, Online. Association
for Computational Linguistics. 723
724
725
726
727
728
- Pepa Atanasova, Jakob Grue Simonsen, Christina Li-
oma, and Isabelle Augenstein. 2022. [Fact Checking
with Insufficient Evidence](#). *Transactions of the Asso-
ciation for Computational Linguistics*, 10:746–763. 729
730
731
732
- Isabelle Augenstein. 2021. [Towards Explainable Fact
Checking](#). Dr. Scient. thesis, University of Copen-
hagen, Faculty of Science. 733
734
735
- Isabelle Augenstein, Timothy Baldwin, Meeyoung Cha,
Tanmoy Chakraborty, Giovanni Luca Ciampaglia,
David Corney, Renee DiResta, Emilio Ferrara, Scott
Hale, Alon Halevy, Eduard Hovy, Heng Ji, Filippo
Menczer, Ruben Miguez, Preslav Nakov, Dietram
Scheufele, Shivam Sharma, and Giovanni Zagni.
2023. [Factuality Challenges in the Era of Large Lan-
guage Models](#). 736
737
738
739
740
741
742
743
- Isabelle Augenstein, Christina Lioma, Dongsheng
Wang, Lucas Chaves Lima, Casper Hansen, Chris-
tian Hansen, and Jakob Grue Simonsen. 2019. [Multi-
FC: A real-world multi-domain dataset for evidence-
based fact checking of claims](#). In *Proceedings of
the 2019 Conference on Empirical Methods in Natu-
ral Language Processing and the 9th International
Joint Conference on Natural Language Processing
(EMNLP-IJCNLP)*, pages 4685–4697, Hong Kong,
China. Association for Computational Linguistics. 744
745
746
747
748
749
750
751
752
753
- Ramy Baly, Georgi Karadzhov, Jisun An, Haewoon
Kwak, Yoan Dinkov, Ahmed Ali, James Glass, and
Preslav Nakov. 2020. [What was written vs. who
read it: News media profiling using text analysis
and social media context](#). In *Proceedings of the 58th
Annual Meeting of the Association for Computational
Linguistics*, pages 3364–3374, Online. Association
for Computational Linguistics. 754
755
756
757
758
759
760
761
- Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wen-
liang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei
Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu,
and Pascale Fung. 2023. [A multitask, multilingual,
multimodal evaluation of chatgpt on reasoning, hal-
lucination, and interactivity](#). *CoRR*, abs/2302.04023. 762
763
764
765
766
767
- Oscar Barrera, Sergei Guriev, Emeric Henry, and Ekate-
rina Zhuravskaya. 2020. [Facts, alternative facts, and
fact checking in times of post-truth politics](#). *Journal
of public economics*, 182:104123. 768
769
770
771
- Ali Borji. 2023. [A categorical archive of chatgpt fail-
ures](#). *CoRR*, abs/2302.03494. 772
773
- Shiqi Chen, Yiran Zhao, Jinghan Zhang, I-Chun Chern,
Siyang Gao, Pengfei Liu, and Junxian He. 2023.
[Felm: Benchmarking factuality evaluation of large
language models](#). *arXiv preprint arXiv:2310.00741*. 774
775
776
777

778	I-Chun Chern, Steffi Chern, Shiqi Chen, Weizhe Yuan,	<i>Conference of the North American Chapter of the</i>	834
779	Kehua Feng, Chunting Zhou, Junxian He, Graham	<i>Association for Computational Linguistics: Human</i>	835
780	Neubig, and Pengfei Liu. 2023. Factool: Factual-	<i>Language Technologies</i> , pages 3670–3686, Seattle,	836
781	ity detection in generative AI - A tool augmented	United States. Association for Computational Lin-	837
782	framework for multi-task and multi-domain scenar-	guistics.	838
783	ios . <i>CoRR</i> , abs/2307.13528.		
784	Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu,	839
785	Kim, James R. Glass, and Pengcheng He. 2023. Dola:	Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea	840
786	Decoding by contrasting layers improves factuality	Madotto, and Pascale Fung. 2023. Survey of halluci-	841
787	in large language models . <i>CoRR</i> , abs/2309.03883.	nation in natural language generation . <i>ACM Comput.</i>	842
788		<i>Surv.</i> , 55(12):248:1–248:38.	843
789	Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu,		
790	Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Ja-	Junyi Li, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun	844
791	son Weston. 2023. Chain-of-verification reduces hal-	Nie, and Ji-Rong Wen. 2023. Halueval: A large-	845
792	lucination in large language models . <i>arXiv preprint</i>	scale hallucination evaluation benchmark for large	846
	<i>arXiv:2309.11495</i> .	language models . <i>CoRR</i> , abs/2305.11747.	847
793	Thomas Diggelmann, Jordan L. Boyd-Graber, Jan-		
794	nis Bulian, Massimiliano Ciaramita, and Markus	Tianyu Liu, Xin Zheng, Baobao Chang, and Zhifang Sui.	848
795	Leippold. 2020. CLIMATE-FEVER: A dataset for	2021. Towards faithfulness in open domain table-	849
796	verification of real-world climate claims . <i>CoRR</i> ,	to-text generation from an entity-centric view . In	850
797	abs/2012.00614.	<i>Thirty-Fifth AAAI Conference on Artificial Intelli-</i>	851
798	Angela Fan, Aleksandra Piktus, Fabio Petroni, Guil-	<i>gence, AAAI 2021, Thirty-Third Conference on In-</i>	852
799	laume Wenzek, Marzieh Saeidi, Andreas Vlachos,	<i>novative Applications of Artificial Intelligence, IAAI</i>	853
800	Antoine Bordes, and Sebastian Riedel. 2020. Gen-	<i>2021, The Eleventh Symposium on Educational Ad-</i>	854
801	erating fact checking briefs . In <i>Proceedings of the</i>	<i>vances in Artificial Intelligence, EAAI 2021, Vir-</i>	855
802	<i>2020 Conference on Empirical Methods in Natural</i>	<i>tual Event, February 2-9, 2021</i> , pages 13415–13423.	856
803	<i>Language Processing (EMNLP)</i> , pages 7147–7161,	AAAI Press.	857
804	Online. Association for Computational Linguistics.		
805	Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony	Yixin Liu, Alex Fabbri, Pengfei Liu, Yilun Zhao, Liny-	858
806	Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vin-	ong Nan, Ruilin Han, Simeng Han, Shafiq Joty,	859
807	cent Y Zhao, Ni Lao, Hongrae Lee, Da-Cheng	Chien-Sheng Wu, Caiming Xiong, and Dragomir	860
808	Juan, et al. 2022. Attributed text generation via	Radev. 2023. Revisiting the gold standard: Ground-	861
809	post-hoc research and revision . <i>arXiv preprint</i>	ing summarization evaluation with robust human	862
810	<i>arXiv:2210.08726</i> .	evaluation . In <i>Proceedings of the 61st Annual Meet-</i>	863
811	Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot,	<i>ing of the Association for Computational Linguistics</i>	864
812	Dan Roth, and Jonathan Berant. 2021. Did aristotle	<i>(Volume 1: Long Papers)</i> , pages 4140–4170, Toronto,	865
813	use a laptop? a question answering benchmark with	Canada. Association for Computational Linguistics.	866
814	implicit reasoning strategies. <i>Transactions of the</i>		
815	<i>Association for Computational Linguistics</i> , 9:346–	Potsawee Manakul, Adian Liusie, and Mark J. F. Gales.	867
816	361.	2023. Selfcheckgpt: Zero-resource black-box hal-	868
817	Guiven. 2023. Llm failure archive (chatgpt	lucination detection for generative large language	869
818	and beyond). https://github.com/giuven95/	models . <i>CoRR</i> , abs/2303.08896.	870
819	chatgpt-failures .		
820	Zhijiang Guo, Michael Schlichtkrull, and Andreas Vla-	Joshua Maynez, Shashi Narayan, Bernd Bohnet, and	871
821	chos. 2022. A survey on automated fact-checking .	Ryan McDonald. 2020. On faithfulness and factu-	872
822	<i>Transactions of the Association for Computational</i>	ality in abstractive summarization . In <i>Proceedings</i>	873
823	<i>Linguistics</i> , 10:178–206.	<i>of the 58th Annual Meeting of the Association for</i>	874
824	Naeemul Hassan, Chengkai Li, and Mark Tremayne.	<i>Computational Linguistics</i> , pages 1906–1919, On-	875
825	2015. Detecting check-worthy factual claims in pres-	line. Association for Computational Linguistics.	876
826	idential debates . In <i>Proceedings of the 24th ACM</i>	Jacob Menick, Maja Trebacz, Vladimir Mikulik, John	877
827	<i>International Conference on Information and Knowl-</i>	Aslanides, H. Francis Song, Martin J. Chadwick,	878
828	<i>edge Management, CIKM 2015, Melbourne, VIC,</i>	Mia Glaese, Susannah Young, Lucy Campbell-	879
829	<i>Australia, October 19 - 23, 2015</i> , pages 1835–1838.	Gillingham, Geoffrey Irving, and Nat McAleese.	880
830	ACM.	2022. Teaching language models to support answers	881
831	Robert Iv, Alexandre Passos, Sameer Singh, and Ming-	with verified quotes . <i>CoRR</i> , abs/2203.11147.	882
832	Wei Chang. 2022. FRUIT: Faithfully reflecting up-	Sewon Min, Kalpesh Krishna, Xixi Lyu, Mike	883
833	dated information in text . In <i>Proceedings of the 2022</i>	Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer,	884
		Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023.	885
		Factscore: Fine-grained atomic evaluation of fac-	886
		tual precision in long form text generation . <i>CoRR</i> ,	887
		abs/2305.14251.	888

889	Dor Muhlgay, Ori Ram, Inbal Magar, Yoav Levine, Nir Ratner, Yonatan Belinkov, Omri Abend, Kevin Leyton-Brown, Amnon Shashua, and Yoav Shoham. 2023. Generating benchmarks for factuality evaluation of language models . <i>CoRR</i> , abs/2307.06908.	944
890		945
891		946
892		947
893		948
894	Niels Mündler, Jingxuan He, Slobodan Jenko, and Martin T. Vechev. 2023. Self-contradictory hallucinations of large language models: Evaluation, detection and mitigation . <i>CoRR</i> , abs/2305.15852.	949
895		950
896		951
897		952
898	Preslav Nakov, David P. A. Corney, Maram Hasanain, Firoj Alam, Tamer Elsayed, Alberto Barrón-Cedeño, Paolo Papotti, Shaden Shaar, and Giovanni Da San Martino. 2021. Automated fact-checking for assisting human fact-checkers . In <i>Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI 2021, Virtual Event / Montreal, Canada, 19-27 August 2021</i> , pages 4551–4558. ijcai.org.	953
899		954
900		955
901		956
902		957
903		958
904		959
905		960
906		961
907	Yixin Nie, Haonan Chen, and Mohit Bansal. 2019. Combining fact extraction and verification with neural semantic matching networks . In <i>The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019</i> , pages 6859–6866. AAAI Press.	962
908		963
909		964
910		965
911		966
912		967
913		968
914		969
915		970
916		971
917	OpenAI. 2023. Gpt-4 technical report . <i>ArXiv</i> , abs/2303.08774.	972
918		973
919	Nicolas Pröllochs. 2022. Community-based fact-checking on twitter’s birdwatch platform . In <i>Proceedings of the Sixteenth International AAAI Conference on Web and Social Media, ICWSM 2022, Atlanta, Georgia, USA, June 6-9, 2022</i> , pages 794–805. AAAI Press.	974
920		975
921		976
922		977
923		978
924		979
925	Vikas Raunak, Arul Menezes, and Marcin Junczys-Dowmunt. 2021. The curious case of hallucinations in neural machine translation . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1172–1183. Online. Association for Computational Linguistics.	980
926		981
927		982
928		983
929		984
930		985
931		986
932	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing</i> . Association for Computational Linguistics.	987
933		988
934		989
935		990
936		991
937	Tal Schuster, Adam Fisch, and Regina Barzilay. 2021. Get your vitamin C! robust fact verification with contrastive evidence . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 624–643. Online. Association for Computational Linguistics.	992
938		993
939		994
940		995
941		996
942		997
943		998
		999
	James Thorne and Andreas Vlachos. 2021. Evidence-based factual error correction . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 3298–3309. Online. Association for Computational Linguistics.	994
		995
		996
		997
		998
		999
	James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for fact extraction and VERification . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 809–819. New Orleans, Louisiana. Association for Computational Linguistics.	994
		995
		996
		997
		998
		999
	Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, et al. 2023. Freshllms: Refreshing large language models with search engine augmentation . <i>arXiv preprint arXiv:2310.03214</i> .	994
		995
		996
		997
		998
		999
	David Wadden, Shanchuan Lin, Kyle Lo, Lucy Lu Wang, Madeleine van Zuylen, Arman Cohan, and Hannaneh Hajishirzi. 2020. Fact or fiction: Verifying scientific claims . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 7534–7550. Online. Association for Computational Linguistics.	994
		995
		996
		997
		998
		999
	Dustin Wright and Isabelle Augenstein. 2020. Claim Check-Worthiness Detection as Positive Unlabelled Learning . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 476–488. Online. Association for Computational Linguistics.	994
		995
		996
		997
		998
		999
	Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. 2023. Do large language models know what they don’t know? In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 8653–8665. Toronto, Canada. Association for Computational Linguistics.	994
		995
		996
		997
		998
		999
	Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. 2019. Defending against neural fake news . In <i>Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada</i> , pages 9051–9062.	994
		995
		996
		997
		998
		999
	Muru Zhang, Ofir Press, William Merrill, Alisa Liu, and Noah A. Smith. 2023a. How language model hallucinations can snowball . <i>CoRR</i> , abs/2305.13534.	994
		995
		996
		997
		998
		999
	Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. 2023b. Siren’s song in the AI ocean: A survey on hallucination in large language models . <i>CoRR</i> , abs/2309.01219.	994
		995
		996
		997
		998
		999

Appendix

A Fact-checking Background

A.1 What is Fact-checking?

Fact-checking is the task of assessing whether claims made in written or manipulated are true. This is typically broken down into the stages of claim detection, evidence retrieval, verdict prediction, and optionally justification prediction [Guo et al. \(2022\)](#); [Augenstein \(2021\)](#).

Claim detection is to identify claims that require verification, which commonly relies on the concept of check-worthiness. In the context of human-written documents, check-worthy claims are regarded as those for which the general public would be interested in knowing the truth ([Hassan et al., 2015](#); [Wright and Augenstein, 2020](#)). However, this may not be adaptable to machine-generated texts. Plausible hallucinations of LLMs make it difficult for general individuals to distinguish whether it is true or false, thus their outputs become less trustworthy than the statements made by human. Current methods tend to check all factual claims of LLM generations ([Chern et al., 2023](#)).

Evidence retrieval aims to find sources supporting or refuting the claim. Claim verification is expected to assess the veracity of the claim and produce justification based on the retrieved evidence. That is, claims are assigned truthfulness labels and explanations for verdicts are produced. A basic form of justification is to highlight the pieces of evidence used to reach a verdict ([Guo et al., 2022](#)).

Method	D	R	Granularity	Knowledge source	Datasets	Task	How_collect
Factcheck-GPT	✓	✓	claim	Google search	✓	Instruction	prompt ChatGPT and human annotation
FacTool (Chern et al., 2023)	✓	✓	article metadata	Google scholar	✓	Generate literature review	prompt ChatGPT
FacTool (Chern et al., 2023)	✓	✓	claim (gold)	Parsed Google search	RoSE/FactPrompts	Summarisation-eval/QA	human annotation: RoSE (Liu et al., 2023)
RARR (Gao et al., 2022)	✓	✓	document	Bing search	NQ.StrategyQA.QReCC	QA	human annotation
CoVe (Dhuliawala et al., 2023)	✓	✓	document	parametric knowledge	CoVe corpus	QA, instruction	human annotation
FELM (Chen et al., 2023)	✓	✗	segment	Google search	✓	Instruction	prompt ChatGPT and human annotate factuality
Self-contradictory (Mündler et al., 2023)	✓	✗	sentence	parametric knowledge	✓	Instruction	prompt ChatGPT,GPT-4 for contradictory sentence
SelfCheckGPT (Manakul et al., 2023)	✓	✗	sentence	parametric knowledge	✓	Generate Wikibio passage	prompt GPT3 and human annotate 3 factual labels
FACTOR (Muhlgay et al., 2023)	✓	✗	sentence	parametric knowledge	✓	Multichoice QA	prompt <i>davinci-003</i> for non-factual completions
HaluEval (Li et al., 2023)	✓	✗	document	parametric knowledge	✓	QA, summarise, dialogue	prompt ChatGPT to generate hallucinated answers
HaluEval (Li et al., 2023)	✓	✗	document	parametric knowledge	✓	Instruction	prompt ChatGPT, human annotate false segments
FactScore (Min et al., 2023)	✓	✗	claim	Wiki Bio Generation	✓	Instruction	prompt ChatGPT to generate biography
FRESHQA (Vu et al., 2023)	✓	✗	facts in answer	parametric knowledge	✓	QA	collect questions with time-changing answers
Snowball (Zhang et al., 2023a)	✓	✗	Yes/No answer	parametric knowledge	✓	QA	human annotation
SelfAware (Yin et al., 2023)	✓	✗	document	reference generations	✓	QA	collect unanswerable questions and prompt ChatGPT

Table 7: **Methods and benchmarks for hallucination Detection (D) and Revision (R)**. FacTool: article metadata is a tuple (paper title, year, authors). CoVe=Chain-of-Verification, CoVe corpus includes four existing datasets: Wikidata, Wiki-category, MultiSpanQA, biographic. 3 labels in SelfCheckGPT: major/minor inaccurate and accurate. Unanswerable questions: the model should express uncertainty instead of delivering conclusive responses. FRESHQA collect four types of questions: false premise, answers never change, change slowly and fast over time.

A.2 Conventional Fact-checking

Previous works either focus on hallucinations in task-specific generations with references (to detect whether the generated output contradicts the source content), such as abstractive summarization ([Maynez et al., 2020](#)), machine translation ([Raunak et al., 2021](#)) and data-to-text generation ([Liu et al., 2021](#)), or concentrate on specific topics e.g. Covid-19 ([Augenstein et al., 2019](#)), politics ([Barrera et al., 2020](#)), climate ([Diggelmann et al., 2020](#)), and specific domains such as journalism, news, social media (e.g. Twitter ([Pröllochs, 2022](#))) and Wikipedia (FEVER: [Thorne et al. \(2018\)](#)). In contrast, we set target for text generation tasks without references such as generative question answering and dialogue systems in open domain and open topic, where the source is the world knowledge.

Moreover, most early studies only perform one or two subtasks in the factual error detection and correction, instead of the whole process. For example, many models estimate a label indicating whether the statement is supported or refuted by the evidence, given a (statement, evidence) pair as input ([Thorne et al., 2018](#); [Nie et al., 2019](#); [Augenstein et al., 2019](#); [Wadden et al., 2020](#)). To adapt to situations where relevant evidence for a statement is not readily available, some works explored how to automatically retrieve evidence that may help support or refute a statement ([Fan et al., 2020](#); [Nakov et al., 2021](#); [Gao et al., 2022](#)).

More recent work has also explored how to correct claims based on retrieved evidence (Thorne and Vlachos, 2021; Schuster et al., 2021; Iv et al., 2022) and how to generate justification/explanation for verdicts on claims (Atanasova et al., 2020). However, most factual correction used human-authored edits from FEVER (Thorne et al., 2018) as both their training and automatic evaluation data. FEVER’s claims were extracted from Wikipedia. This limits the generalisability of these fact-checking models over generations of LLMs across various tasks and diverse domains.

Our goal is to automatically detect and correct factual errors end to end for open-domain factual knowledge hallucinations in a unified framework.

A.3 LLM Fact-checking

The phenomenon that LLMs produce outputs that are seemingly plausible while deviating from the user input, previously generated context, or factual knowledge, is commonly referred to as hallucination (Zhang et al., 2023b). Based on the timing of the LLM life cycle, LLM hallucinations can be addressed during pretraining by automatically selecting reliable data or filtering out noisy data to mitigate hallucinations, in supervised fine-tuning by curating small volume of high-quality training data, in reinforcement learning from human feedback (RLHF), and during inference by decoding strategies (Zhang et al., 2023b) We focus on approaches applied after inference.

Methods For post-processing approaches to alleviating LLM hallucinations, recent studies can be roughly classified into two categories depending on whether they rectify errors: (1) detecting and correcting factual errors for free-form text; and (2) only detecting whether a text contain hallucinations (*Yes* or *No*). Both of them resort to either external knowledge or parametric knowledge to identify and rectify factually-incorrect statements (Gao et al., 2022; Chern et al., 2023; Manakul et al., 2023; Dhuliawala et al., 2023). We used external knowledge retrieved from Google.

Our work falls into the first category. Though Self-contradictory (Mündler et al., 2023) involves revision, they aim to remove the conflicting information between the original sentence and the synthetically-generated contradictory sentence, instead of correcting factual errors in the original sentences. We classify it into the second category: detection only. RARR (Gao et al., 2022), FacTool (Chern et al., 2023) and CoVe (Dhuliawala et al., 2023) are three most relevant work to ours.

Given a LLM response, RARR and CoVe first generate a series of questions covering different aspects of the response, which serve as queries in the evidence retrieval, and then edit the whole response to correct factual errors. Such coarse granularity verification may miss out incorrect statements, particularly over long documents, and also makes it difficult to spot false spans precisely, thus disabling fine-grained (e.g., correct only a false number in a statement) and flexible edits (e.g., delete a completely-hallucinated sentence). Additionally, revising the whole document tends to result in poor preservation of the original input (e.g., style, vocabulary and structure), introducing irrelevant descriptions and even new hallucinations. Claim-level editing empowers precise correction and good preservation.

FacTool performs fact-checking over claims. However, gold claims are required as the input of the system. That is, users need to first decompose an output from a LLM into a list of checkable atomic claims by themselves before using FacTool to check, which complicates the fact-checking process. Moreover, it is expensive to use FacTool to check a piece of text, since the whole checking process calls APIs including OpenAI (\$0.06/1K tokens), Serper (\$1.00/1k queries) and Scraper.² This also challenges the evaluation where online API is not allowed to call with the consideration of the internal data protection.

We attempt to alleviate these issues in our framework. We decompose the fact-checking task into eight subtasks. The design of decomposing and decontextualising a long document into independent sentences and then to atomic claims allows inputs of any granularity: document, sentence or claim. The pipeline equipped with check-worthiness selection also naturally endows the flexibly-customised verification, such as skipping subjective statement, commonsense and the knowledge is well-known by the individual.

Datasets From the perspective of the evaluated benchmarks, as shown in Table 7, studies of the first category generally evaluate their methods on existing QA datasets, or revise hallucinations in a specific

²<https://www.scraperaapi.com/pricing/>

Dataset	Granularity	Factual label	Revision	Length	Size
HaluEval	document	✓	✗	82.0	4,507
FELM-WK	segment	✓	✗	51.1	184
FactPrompts	claim	✓	✗	41.8	50
Factcheck-GPT	claim	✓	✓	73.1	94

Table 8: Statistics of world-knowledge factuality evaluation benchmarks. Length=the average number of words of LLM responses.

1081 topic such as literature review and biographic generations (Chern et al., 2023; Dhuliawala et al., 2023).
1082 These topics may not be frequently requested by general users in real-world scenarios.

1083 Studies of the second category contribute a spectrum of benchmarks to detect diverse hallucinations,
1084 such as synthetically-generated contradictory sentences (Mündler et al., 2023), deliberately-generated
1085 hallucinated answers (Li et al., 2023) and non-factual completions given a prefix context (Muhlgay et al.,
1086 2023). Manakul et al. (2023) manually annotate factual labels (major/minor inaccurate and accurate)
1087 given a sentence in the generated Wikibio passage.

1088 Interestingly, Yin et al. (2023) collected 1,032 unanswerable questions from five diverse categories
1089 *no scientific consensus, imagination, completely subjective, too many variables, philosophical*, and their
1090 2,337 answerable counterparts. Unanswerable questions refer to questions where the model should express
1091 uncertainty instead of delivering conclusive responses. Zhang et al. (2023a) collected three datasets, with
1092 500 questions (all *No* or all *Yes* answers) for each. One focuses on one type of question, including whether
1093 a number is a prime, senator search (whether a US city has a specific university), and whether there is a
1094 flight from one city to another given a graph connection.

1095 However, these datasets are either only applicable in detection, or originate from a single task like
1096 biography writing (Min et al., 2023), without accounting for variations across different generations.
1097 HaluEval’s annotation over Alpaca 5K responses of various instructions, which is one of the most similar
1098 work to ours. They ask human annotators to label whether the response contains hallucinated information
1099 (*Yes* or *No*) and list the corresponding spans if there exist errors (Li et al., 2023).³ FELM with 184
1100 world-knowledge questions is labelled in the granularity of segments, while ours are over fine-grained
1101 claims to more precisely locate factual errors. Moreover, our annotations not only include factual labels
1102 of each claim, but the revised text and labels of all involved subtasks as well, e.g., decomposition of a
1103 sentence into a list of independent claims, check-worthiness of a sentence/claim, evidence stance and so
1104 on.

1105 FacTool evaluate over a knowledge-based QA dataset FactPrompts consisting of 50 (prompt, response)
1106 pairs. It is annotated by authors over atomic claims and their factual labels (true/false), but the responses
1107 tend to be short, instead of long documents (see Table 8). Overall, our dataset offers both factual labels
1108 and the revised text in three-level granularity — atomic claims, decontextualised sentences and responses,
1109 for LLM answers, with an emphasis of long documents.

1110 FELM (Chen et al., 2023) is the most relevant concurrent work with ours, but only annotated sentence-
1111 level *true or false* labels (no correction). We construct a new dataset which collects (question, ChatGPT
1112 response) pairs in real conversations. Annotators identify and edit factual errors for each atomic claim
1113 decomposed and decontextualised from the original long-form responses. This is expected to serve as a
1114 benchmark to evaluate the performance of fact-checkers.

³The hallucination is considered from the following three aspects: unverifiable, non-factual, and irrelevant.

B Dataset 1115

B.1 Sources 1116

Twitter posts and in-house brainstorming: We first collect (question, response) pairs from ChatGPT/GPT-4 failures found on social media, in Web articles and related papers.⁴ The query should satisfy criteria that the corresponding response must have factual errors, rather than failures regarding reasoning, math, coding, bias and so on; (query, response) also should be independent from a dialog. This results in 23 examples. We additionally brainstorm a spectrum of questions depending on individual usage experience of ChatGPT, and then select 22 questions whose responses contain factually-false content by manually verifying suspicious facts. 1117-1123

Dolly-15k It consists of 15,011 examples, with eight categories ranging from closed, open and general QA, to creative writing, brainstorming, information extraction, summarisation and classification.⁵ Since we pay attention to open-domain generations and responses with more factual statements, closed and open question answering pairs are chosen to be the database. 1124-1127

We first generate ChatGPT responses for 1,773 closed QA pairs without using context information (a paragraph extracted from Wikipedia relevant to the question), and 3,700 open QA pairs. After filtering questions that cannot be answered without context as well as questions ChatGPT does not answer, we further filtered responses with fewer than 200 characters. Taking human answers as the gold reference, we assume that if machine generations are semantically far from human answers, they may contain false information. So we keep the examples where the cosine similarity ≤ 0.5 between human answer and machine response based on SimCSE sentence embedding. Finally, we select 563 examples from closed QA and 528 from open QA, thus 1,136 (question, response) pairs in total with the 45 from the first source. 1128-1135

B.2 Data Selection 1136

The whole annotation process is extremely time-consuming, about 15-30 minutes for an instance even if with intermediate results from automatic methods to ease the procedure. This requests us to cherry-pick examples that highly satisfy two criteria — fact-intensive and factually-false. Therefore, we leverage FactScore to filter cases with the following four steps. 1137-1140

Sentence split and atomic claims breaking-down We first split a document into sentences using NLTK tokenizer. The most straightforward way is to prompt ChatGPT to split a sentence into claims given the response as context. However, processing sentences one by one consumes both time and API tokens. 1141-1143

Therefore, given the whole response as the context and the first sentence of the response, we ask ChatGPT to break the input sentence into independent atomic claims, and also continue the decomposition of the next sentence of the response (see the prompt in Section B.5). Specifically, ChatGPT is given three demonstration examples, so that it can follow the instruction to first break down the input sentence into atomic claims, and then sequentially find the next sentence and make the splits. Over 90% examples follow the instruction, breaking down the whole response. 105 out of 1,136 examples only decompose the first sentence, on which we process sentence by sentence based on the NLTK sentence splits. 1144-1150

Another reason why we ask ChatGPT to re-split the response into single sentences is that we observed that some sentences are incorrectly split into smaller units by NLTK, such as decomposing a paper reference into a set of meta data, while ChatGPT can remain the citation reference as a whole.⁶ A weakness of ChatGPT outputs compared with traditional models is that it is sometimes non-trivial to parse the results from the text-free responses when ChatGPT does not follow the output format as the instruction. In such cases, we have to process examples specifically. 1151-1156

Discussion: One may argue that why not directly decompose the whole response into atomic claims, but through single sentences and then to atomic claims. There are two reasons. 1157-1158

⁴<https://github.com/giujen95/chatgpt-failures>

⁵Its use is subject to the *CC BY-SA 3.0* license.

⁶In our dataset, we prioritise sentence splits by ChatGPT, using NLTK results for unsuccessfully-parsed instances. The prompt is initialised with the first sentence split by NLTK.

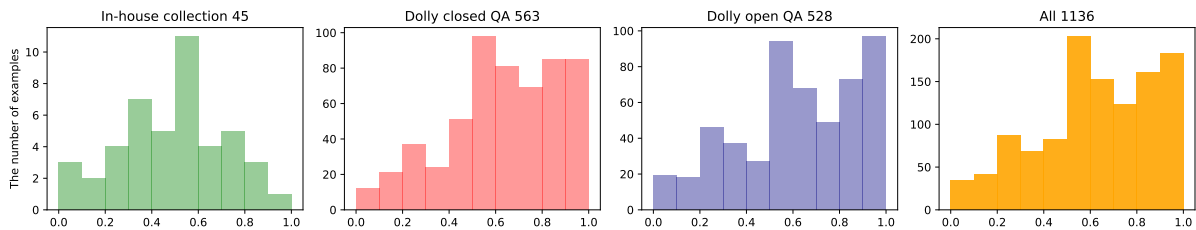


Figure 3: FactScore distribution for three component sources and their combination.

- *Avoid distortion*: atomic claims decomposed and decontextualized from a whole response by models such as ChatGPT tend to either lose or hallucinate information compared to the original response, while the quality of atomic claims of a single sentence is much better;
- *Improve annotation quality*: sentences as the intermediate state, it is easier for annotators to going through 1-5 atomic claims for a sentence as one annotation unit, instead of >5 claims for a whole response (particularly long documents), which helps annotators to pay attention and improve the annotation quality.

Evidence collection for atomic claims Given an atomic claim, following Gao et al. (2022), we first prompt ChatGPT to generate search queries for the claim, and then Google Search is used to get relevant web pages. We further split the retrieved documents into passages by sliding windows, and use a re-ranker combining lexical and semantic similarity to identify the most relevant passages for the given query, in which Sentence-BERT (Reimers and Gurevych, 2019) serves for semantic embeddings.⁷ Finally, we aggregate evidence for all queries and select the top-5 evidences per atomic claim.

FactScore calculation FactScore (Min et al., 2023) is an automatic metric for fine-grained evaluation of the factuality of long-form generations. Given a generation, FactScore is calculated as the percentage of atomic claims within the generation that are supported by a knowledge source. For verifying the claim, we use the gathered evidences as input, along with the claim, and an instruction-tuned model as the verifier.

Example selection Figure 3 shows the FactScore distribution of three component sources and the whole data set. We keep all 45 pairs from the first source, and Dolly examples whose FactScore is less than 0.2, resulting in 33 closed question-answering pairs and 37 open questions, in total of 115 examples. We remove a similar question (7 and 13 are similar), and four questions where the LLM did not provide helpful answers due to its inherent disability to access real-time data. For example, the LLM cannot browse the internet and does not have access to the latest information (“which paper got the most citations in the question generation area?” and “which large language model contains the most parameters?”), or up-to-date data and event-specific details (“who was the general chair of COLING 2023”), or individual information (“what are the awards that Prof. William Yang Wang have?”). We eventually annotate 110 examples in our first annotation stage, and more cases would be annotated in the next stage.

B.3 Annotation

As many studies illustrated, annotating a LLM factuality benchmark is a highly challenging task (Chen et al., 2023; Li et al., 2023). Our preliminary trials in which authors manually annotate labels of all subtasks empirically confirm the pain.

Preliminary Trial Based on the annotation guideline (see Appendix C), we first conduct an in-house annotation for ten examples, each example has two annotators. We annotate the whole process for all steps and manually type results into a json file as the pre-defined format. This attempt exposes three issues.

First, it is extremely time-consuming. It takes more than four hours for a fully-focused annotator to annotate a document of ~400 words with about 20 sentences, in which evidence collection takes the most time and effort, particularly for topics with which the annotator is not familiar. Second, it is ineffective

⁷cross-encoder/ms-marco-MiniLM-L-6-v2: <https://huggingface.co/cross-encoder/ms-marco-MiniLM-L-6-v2>

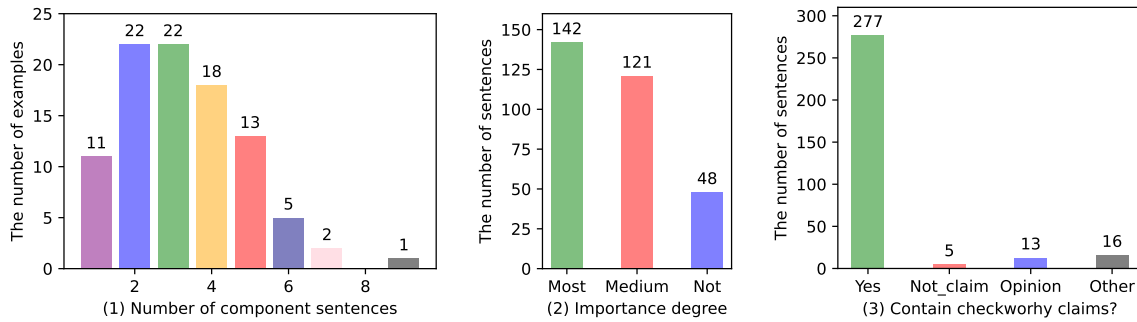


Figure 4: **Sentence analysis:** (1) Distribution of the number of sentences for each response; (2) Importance degree of sentences to answer the question (The distribution of the most importance sentences to answer the question, intermediate important and not important; (3) The number of sentences across four types in terms of whether the sentence contains statements requiring fact-checking, Not_claim refers to *not a claim*, such as a question.

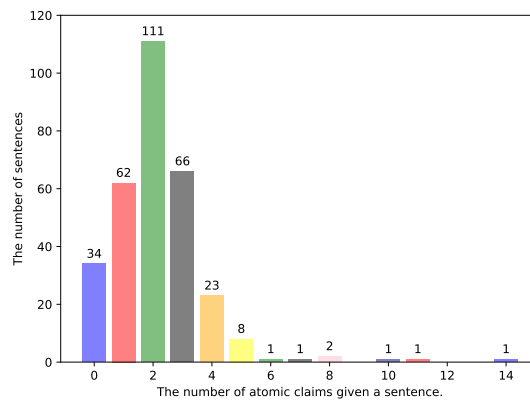


Figure 5: The distribution of component atomic claims amount given a sentence.

to extract relevant evidence passages by human eyes and basic string matching from retrieved Google search documents. This not only takes time but most importantly takes the risk of missing the most relevant evidence due to limited traversal. It is impractical for human to go through all relevant Web articles and select the most semantically-relevant and reliable ones in limited time. Humans are good at judging or making decisions while machines are good at traversing. Lastly, it is hard to reach high agreement between annotators, especially for subtasks of decomposition, evidence collection and stance identification.

B.4 Data Analysis

During annotation, we remove another 16 examples because there is no standard gold answer for these questions, such as seven involving a flow of procedures, six non-factual questions, one tricky riddle-like question, one broken generated answer and one highly-disagreed case, resulting in 94 instances.

From the perspective of LLM users, we may expect to assess any answers and identify whether they are true and reliable, including the cases deleted in our setting. It should be highlighted that the questions involving a flow of procedures, tricky riddles, or non-factual questions need to be verified, while they are just out of the verification scope of the current fact-checking pipelines that only concern facts. The causality and the global logic behind the whole answer are under-explored.

Sentences: Most responses contain 2-5 sentences, with the longest response encompassing 9 sentences as shown in Figure 4. 142 sentences are considered to be the most importance sentences, 121 and 48 fall into intermediate and not important. 278 sentences contain checkworthy statements, 16, 12 and 5 are categorised into *other*, *opinion* and *not a claim*.

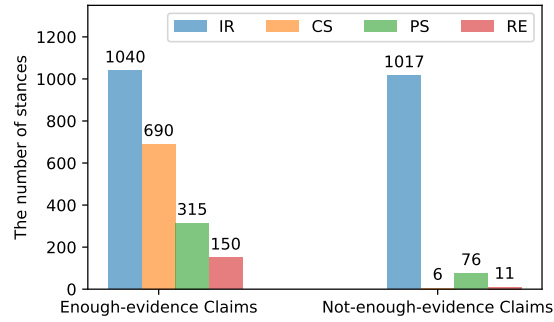


Figure 6: **Stance distribution** of claims with enough automatically-retrieved evidence to determine the factuality vs. claims without enough evidence (*CP*:completely support, *PS*: partially support, *RE*: refute, *IR*: irrelevant).

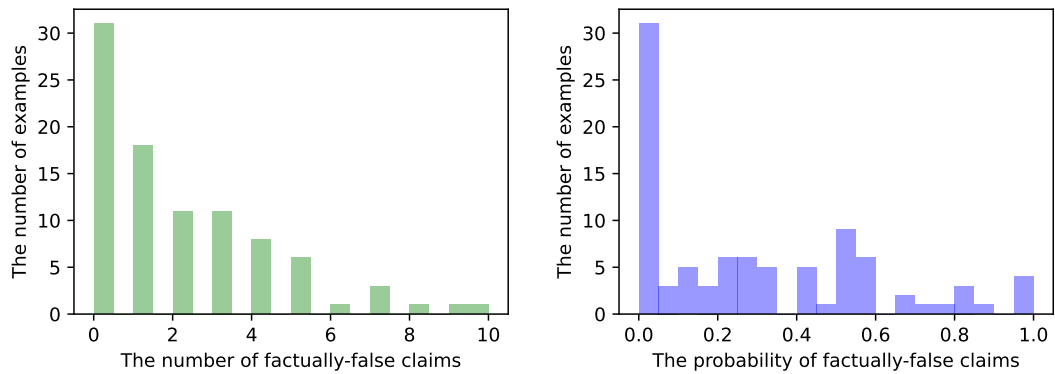


Figure 7: The number of false claims given an example.

How does the evidence support the claim? Two-thirds pieces of irrelevant evidence (2057/3305). We compare the stance distribution of claims in which automatically-retrieved evidence is enough to determine its factuality and the claims that cannot be determined by automatic evidence in Figure 6. Though the majority of evidence are irrelevant for both groups, there are only 17 strong-position stances (“completely support”: *CS* and “refute”: *RE*) in the latter, compared with 690 *CS* and 150 *RE* in the former.

B.5 Prompt to Generate Atomic Claims

Table 9: **Prompt** used to decompose and decontextualize a sentence into a set of independent atomic claims. We use three examples as demonstrations to elicit ChatGPT follow the instruction, and break the response into sentences, as well as breaking a sentence into atomic claims.

Field	Content
Prompt	<p>Depending the context, please breakdown the following sentence into independent facts.</p> <p>Context: The United States has had two black presidents: Barack Obama, who served two terms from 2009 to 2017, and Donald Trump, who served one term from 2017 to 2021. Obama was the first black president in the history of the United States. He was born in Honolulu, Hawaii, to a mother from Kansas and a father from Kenya. Trump was the second black president. He was born in New York City and previously served as a businessman and reality television personality.</p> <p>The sentence is: The United States has had two black presidents: Barack Obama, who served two terms from 2009 to 2017, and Donald Trump, who served one term from 2017 to 2021. Atomic facts for this sentence are:</p> <p>["The United States has had two black presidents: Barack Obama and Donald Trump.", "Black president Barack Obama served two terms from 2009 to 2017.", "Black president Donald Trump served one term from 2017 to 2021."]</p> <p>The sentence is: Obama was the first black president in the history of the United States. Atomic facts for this sentence are:</p> <p>["Obama was the first black president in the history of the United States."]</p> <p>The sentence is: He was born in Honolulu, Hawaii, to a mother from Kansas and a father from Kenya. Atomic facts for this sentence are:</p> <p>["Barack Obama was born in Honolulu, Hawaii.", "Barack Obama mother was from Kansas.", "Barack Obama father was from Kenya."]</p> <p>The sentence is: Trump was the second black president. Atomic facts for this sentence are:</p> <p>["Trump was the second black president."]</p> <p>The sentence is: He was born in New York City and previously served as a businessman and reality television personality. Atomic facts for this sentence are:</p> <p>["Donald Trump was born in New York City.", "Donald Trump previously served as a businessman", "Donald Trump previously served as a reality television personality."]</p>

Field	Content
	<p>Context: In 1980, the oldest justice on the United States Supreme Court was Justice William O. Douglas. He was born on October 16, 1898, and served on the Supreme Court from 1939 until his retirement in 1975. Therefore, in 1980, Justice Douglas was still alive and would have been the oldest serving justice on the Court at that time.</p> <p>The sentence is: In 1980, the oldest justice on the United States Supreme Court was Justice William O. Douglas.</p> <p>Atomic facts for this sentence are:</p> <p>["In 1980, the oldest justice on the United States Supreme Court was Justice William O. Douglas."]</p> <p>The sentence is: He was born on October 16, 1898, and served on the Supreme Court from 1939 until his retirement in 1975.</p> <p>Atomic facts for this sentence are:</p> <p>["Justice William O. Douglas was born on October 16, 1898." "Justice William O. Douglas served on the Supreme Court from 1939 until his retirement in 1975."]</p> <p>The sentence is: Therefore, in 1980, Justice Douglas was still alive and would have been the oldest serving justice on the Court at that time.</p> <p>Atomic facts for this sentence are:</p> <p>["Therefore, in 1980, Justice Douglas was still alive." "Justice William O. Douglas would have been the oldest serving justice on the Court in 1980."]</p> <p>Context: There have been only four female presidents of the United States in the country's history, so it is difficult to determine an average height for this group. The four female presidents were: 1.Abigail Adams (1797-1801) 2.Marilyn Carlson Nelson (2009-2013) 3.Luci Baines Johnson (1973-1977) 4.Hillary Clinton (2017-2021)</p> <p>The sentence is: There have been only four female presidents of the United States in the country's history, so it is difficult to determine an average height for this group.</p> <p>Atomic facts for this sentence are:</p> <p>["There have been only four female presidents of the United States in the country's history.", "It is difficult to determine an average height for four female presidents of the United States."]</p> <p>The sentence is: The four female presidents were: 1.Abigail Adams (1797-1801) 2.Marilyn Carlson Nelson (2009-2013) 3.Luci Baines Johnson (1973-1977) 4.Hillary Clinton (2017-2021)</p> <p>Atomic facts for this sentence are:</p> <p>["Abigail Adams (1797-1801) is a female president of the United States.", "Marilyn Carlson Nelson (2009-2013) is a female president of the United States.", "Luci Baines Johnson (1973-1977) is a female president of the United States.", "Hillary Clinton (2017-2021) is a female president of the United States."]</p> <p>Context: {} The sentence is: {} Atomic facts for this sentence are:</p>

Field	Content
Input	<p>Context: It is difficult to say exactly how many stars are formed in the Milky Way each year, as the rate of star formation can vary over time. However, it is generally thought that the Milky Way produces about one new star per year on average. This may seem like a small number, but keep in mind that the Milky Way is a vast galaxy with billions of stars, so even a small number of new stars added each year can add up quickly.</p> <p>First sentence of the document: It is difficult to say exactly how many stars are formed in the Milky Way each year, as the rate of star formation can vary over time.</p>
Output	<p>Atomic facts for this sentence are:</p> <p>["It is difficult to say exactly how many stars are formed in the Milky Way each year.", "The rate of star formation in the Milky Way can vary over time."]</p> <p>The sentence is: However, it is generally thought that the Milky Way produces about one new star per year on average.</p> <p>Atomic facts for this sentence are:</p> <p>["The Milky Way produces about one new star per year on average."]</p> <p>The sentence is: However, This may seem like a small number, but keep in mind that the Milky Way is a vast galaxy with billions of stars, so even a small number of new stars added each year can add up quickly.</p> <p>Atomic facts for this sentence are:</p> <p>["The Milky Way is a vast galaxy with billions of stars.", "Even a small number of new stars added each year can add up quickly in the Milky Way."]</p>

C Annotation Guidelines

Based on the annotation framework, we aim to construct a benchmark for evaluation, assessing whether a fact-checking pipeline is effective to detect and correct factual errors in texts. The benchmark comprises annotations for the whole fact-checking process spanning seven subtasks for 110 (question, response) pairs, in which most responses are generated by ChatGPT, some are by GPT-4. This section introduces the annotation guidelines, and Section 3 and 3.2 provides details of data collection and annotation.

For each example, annotators are given a pair of (question, response). A response is either an answer generated by LLMs responding to users' question, or a document returned by LLMs according to users' request. Annotators are required to give outputs of each step shown in Figure 1. We describe how to annotate for component subtasks throughout the pipeline, particularly clarify how to deal with possible ambiguous scenarios.

C.1 Decompose

It is subjective to decide the granularity of decomposition. We may aim to break down a long document into a set of atomic claims, while the definition of a atomic claim varies. Here, we practically apply the following strategy:

- Start by decomposing into single sentences.
- If the sentence contains too much information, break it into several components, but annotators do not overdo it, e.g., decomposing *Capitol Hill riots happened on January 6, 2021* to one claim for year and one for day.
- If several pieces of information are strongly dependent with each other, they are expected to co-occur in one snippet of evidence text, no more breaking-down is needed.

C.2 Decontextualise

The criteria of decontextualisation is to ensure that all separated statements fully preserve semantics presented in original context. For example, statement that *it happened on Jan 6, 2021* loses information in decomposition, which makes it uncheckable. In such cases, annotators should replace pronouns, such as *it, they, those, these, this, that*, with specific entities or events after decomposition. Decontextualisation is mostly needed over cases with coreference relation. For complex relation, such as two sentences are strongly dependent with each other, we encourage to go back to step of decomposition and keep the original text without breaking-down.

C.3 Identify checkworthy claim

We consider two aspects in check-worthiness identification:

- If a statement presents subjective opinions, then it is not check-worthy.
- If the objective facts presented in a statement are commonsense, as obvious as *sun rises from the east*, it is not worth checking.

Therefore, we regard a statement as check-worthy claim when it presents objective facts, and these facts are not apparent commonsense. There is a special case. If the objective facts presented in a statement are not publicly available information. Namely, we cannot collect any evidence over web pages related to the claim, such as personal experience. They are regarded as uncheckable claims.

Specifically, for each statement, annotators are asked to answer two questions. Which category does this claim fall into? (1) factual claim; (2) subjective opinion; (3) not a claim; and (4) other. Is this statement worth checking? (1) Yes; and (2) No.

C.4 Retrieve and collect evidence	1264
Given a checkworthy claim, annotators are asked to search and collect five most relevant snippets of text as evidence based on general web pages (including Wikipedia pages). Annotators are allowed to use any forms of queries in retrieval, e.g. questions covering some aspects of the claim, or entities in the claim, and they need to record all queries and indicate those used for searching the most relevant evidence.	1265 1266 1267 1268
Note that five pieces of evidence is not a hard criteria. If less than five (even only one) pieces of evidence are sufficient to verify the input claim, and they are from reliable sources, annotators are allowed to collect <5 results. Meanwhile, if a claim involves a controversial topic, annotators are also encouraged to collect more than five results.	1269 1270 1271 1272
For each piece of evidence, record meta-data including (1) corresponding query, (2) citation (URL) of the web page from which this piece of evidence is extracted, (3) judgement of whether the source of evidence is reliable or not, ⁸ and (4) indicator whether this individual evidence is sufficient to verify the input claim.	1273 1274 1275 1276
Aforementioned guidelines are applicable to claims for which there exist evidence over web pages. However, there are situations where there isn't any information in public web pages, e.g. personal experience. They are objective facts, but are not extensively known by the public. Put differently, they are uncheckable. Annotators can give empty list of evidence for uncheckable claims.	1277 1278 1279 1280
C.5 Identify evidence stance	1281
Given a claim and five pieces of most relevant evidence, annotators judge whether the evidence supports, partially supports, refutes or is irrelevant to the claim (see definition of stance in Section 2).	1282 1283
C.6 Determine correction	1284
For a claim, there will be K snippets of text (evidence), corresponding stance vectors $[s_1, s_2, \dots, s_K]$ and source reliability values $[r_1, r_2, \dots, r_K]$. We skip all irrelevant evidence and follow the criteria below to determine whether edits are needed for a claim.	1285 1286 1287
<ul style="list-style-type: none"> • If the claim is completely supported by evidence, no edit. • If the claim is completely refuted by evidence, check evidence and make edits accordingly one by one. • If some evidence support the claim and some refute, this means there are conflicts between evidence. In such scenario, we consider both the source reliability and the number of evidence falling into each stance. If the voice of "refute" is stronger than "support", we make editing, otherwise remain original text. • If some refute and some partially support, there are two possible situations depending on whether supported partition is the same as the refuted partition: (1) if what is supported and what is refuted are the same partition, there are conflicts between evidence, follow the steps above; and (2) if they support and refute different partition of the claim, edit the refuted partition. 	1288 1289 1290 1291 1292 1293 1294 1295 1296 1297 1298
C.7 Edit, Merge and Deduplicate	1299
In correction, we keep the principle of making minimal edits against the original text to correct factual errors. Annotators do not add extra information provided by evidence that is not directly targeted at factual errors. No extra deletion, insertion and addition. Finally, annotators merge all statements either revised or original ones in order, and deduplicate repeated information with principle of minimal edits.	1300 1301 1302 1303

⁸Source reliability can also automatically be collected from MBFC/AllSides/Politifact/, but they apply for a small number of sources.

D Conflicting Evidence Example

2. "land" can be found on Earth, Jupiter and Mars but not on Neptune and the Sun.

Evidence 1

For land on astronomical objects, see Planetary surface . For other uses, see Land (disambiguation) . Land , also known as dry land , ground , or earth , is the solid terrestrial surface of Earth not submerged by the ocean or another body of water . It makes up 29.2% of Earth's surface and includes all continents and islands . Earth's land surface is almost entirely covered by regolith , a layer of rock , soil , and minerals that forms the outer part of the crust .

Link: <https://en.wikipedia.org/wiki/Land>

Does this evidence support the claim?

Completely support Partially support Refute Irrelevant

Evidence 2

Has there been any lore revealed how exactly a city can exist on Neptune? Yes, technically Neptune has a surface, but the pressure and temperatures are so high as to make a settlement existing there (even if it be a ship that crash landed) rather... unlikely. Read more More posts you may like

Link: https://www.reddit.com/r/DestinyTheGame/comments/117tvbj/neptunes_surface/

Does this evidence support the claim?

Completely support Partially support Refute Irrelevant

Figure 8: A claim with conflicting stance evidence: *partially support* and *refute*.

E Confusion Matrix of Subtasks

1305

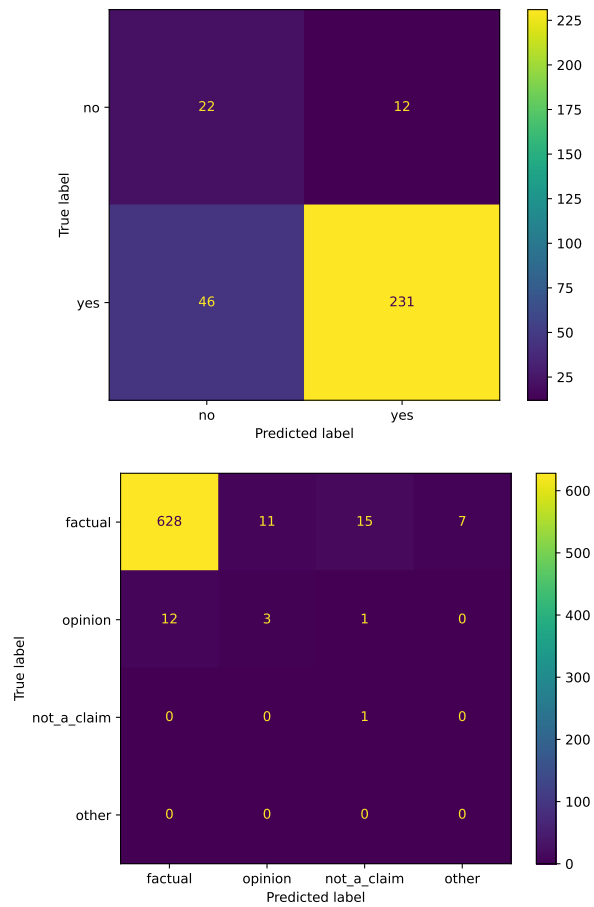


Figure 9: **ChatGPT checkworthiness** detection confusion matrix: sentence (top) and claim (bottom)

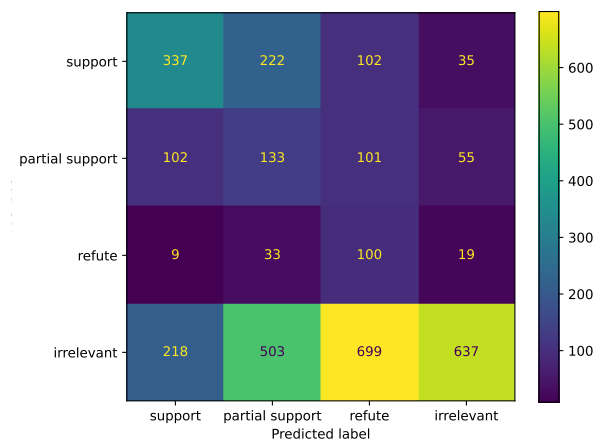


Figure 10: **(claim, evidence) Stance** detection confusion matrix based on ChatGPT with four labels: completely support (support in short in the figure), partial support, refute and irrelevant.

label	LLaMA2			ChatGPT			RoBERTa-large-mnli		
	Prec	Recall	F1	Prec	Recall	F1	Prec	Recall	F1
irrelevant	0.69	0.54	0.61	0.85	0.47	0.60	0.71	0.52	0.60
refute	0.06	0.49	0.10	0.14	0.53	0.22	0.15	0.65	0.24
support	0.48	0.12	0.19	0.53	0.76	0.63	0.75	0.65	0.69

Table 10: **Stance** detection by prompting zero-shot ChatGPT, LLaMA2 and employing a fine-tuned NLI model in three label space: support, refute and irrelevant.

F Example of Subtask 5

1306

Prompt1 without question	Given a document containing factual errors, please correct the errors in the document depending on a corresponding list of factually true claims. Note that preserve the linguistic features and style of the original document, just correct factual errors. document: {response} true claims: {correct claims} revised document:
Prompt2 with question	Given a question, and an answer containing factual errors, please correct the errors in the document depending on a corresponding list of factually true claims. Note that preserve the linguistic features and style of the original document, just correct factual errors. question: {prompt} document: {response} true claims: {claims} revised document:
Question Response	Who was the oldest justice on the US supreme court in 1980? In 1980, the oldest justice on the United States Supreme Court was Justice William O. Douglas. He was born on October 16, 1898, and served on the Supreme Court from 1939 until his retirement in 1975. Therefore, in 1980, Justice Douglas was still alive and would have been the oldest serving justice on the Court at that time.
Correct claims	In 1980, the oldest justice on the United States Supreme Court was Justice William J. Brennan Jr. Justice William O. Douglas was born on October 16, 1898. Justice William O. Douglas served on the United States Supreme Court from 1939 until his retirement in 1975. In 1980, Justice William O. Douglas died. Justice William O. Douglas was the oldest serving justice in the history United States Supreme Court, sat on the Court from 1939 to 1975 and died in 1980.
Human revision	In 1980, the oldest justice on the United States Supreme Court was Justice William J. Brennan Jr.
Prompt1, GPT-3.5-turbo	In 1980, the oldest justice on the United States Supreme Court was Justice William J. Brennan Jr. He was born on October 16, 1898, and served on the Supreme Court from 1939 until his retirement in 1975. Therefore, in 1980, Justice Douglas had already passed away and would not have been the oldest serving justice on the Court at that time.
Prompt1, GPT-4	In 1980, the oldest justice on the United States Supreme Court was not Justice William O. Douglas, but rather Justice William J. Brennan Jr. Justice Douglas, who was born on October 16, 1898, notably served on the Supreme Court from 1939 until his retirement in 1975. However, by 1980, Justice Douglas was no longer alive, having passed away that year. Previously, he had held record as the oldest serving justice in the history of the United States Supreme Court.
Prompt2, GPT-3.5-turbo	In 1980, the oldest justice on the United States Supreme Court was Justice William J. Brennan Jr. He was born on October 16, 1898 and served on the Supreme Court from 1939 until his retirement in 1975. Therefore, in 1980, Justice Douglas had already passed away and Justice Brennan would have been the oldest serving justice on the Court at that time.
Prompt2, GPT-4	In 1980, the oldest justice on the United States Supreme Court was Justice William J. Brennan Jr., not Justice William O. Douglas. While Justice Douglas was indeed born on October 16, 1898, and served on the Supreme Court from 1939 until his retirement in 1975, in 1980, he was no longer alive, contrary to the original claim. It's true that Justice Douglas was considered the oldest serving justice in the history of the United States Supreme Court, his tenure spanning from 1939 to 1975, but by 1980, he had passed away. Thus, the distinction of being the oldest justice on the Supreme Court in 1980 goes to Justice Brennan Jr.

Table 11: Outputs (revised response) comparison using different prompts and models (GPT-3.5-turbo and GPT-4). The response by Prompt1 using GPT-4 is preferred.

G Annotation Interfaces

G.1 Decomposition, Decontextualization and Check-worthiness detection

Example 1

Who was the oldest justice on the US supreme court in 1980?

1. In 1980, the oldest justice on the United States Supreme Court was Justice William O. Douglas.
2. He was born on October 16, 1898, and served on the Supreme Court from 1939 until his retirement in 1975.
3. Therefore, in 1980, Justice Douglas was still alive and would have been the oldest serving justice on the Court at that time.

Please upload a generated result json.

UPLOAD

Please select a sentence by clicking it on the left.

1. In 1980, the oldest justice on the United States Supreme Court was Justice William O. Douglas.

Decontextualized sentence

COPY

Is this sentence the most important one to answer the question?

Yes No, intermediate No, not important

Does the sentence contain statements requiring fact-checking (checkworthy)?

Yes No

Please add claims if you see any.

Statement 1

In 1980, the oldest justice on the United States Supreme Court was Justice William O. Douglas.

In 1980, the oldest justice on the United States Supreme Court was Justice William O. Douglas. COPY

What type of claim is this statement?

Factual Claim
 Opinion
 Not a Claim (e.g. a question)
 Other (e.g. I'm a language model.)

REMOVE

Which statement is the most important statement in the sentence:

1

ADD CLAIM
SAVE

Figure 11: Screenshot of the first annotation interface for Decomposition, Decontextualization and Check-worthiness detection.

G.2 Evidence stance identification and Claim correction

1309

The screenshot displays a web-based annotation interface. On the left, a yellow bar contains the claim: "In 1980, the oldest justice on the United States Supreme Court was Justice William O. Douglas." Below this, a "NEXT EVIDENCE" button is visible. A "Please upload a generated result json" section includes an "UNLOAD" button and a file named "screen2.json".

The main interface on the right is titled "Please select a claim by clicking it on the left." and lists five evidence items:

- Evidence 1:** "In 1980, the oldest justice on the United States Supreme Court was Justice William O. Douglas." This is a chronologically ordered list of Supreme Court justices from the earliest to the most recent. The date the justice took the judicial oath is used as the beginning date of service. Chief justices are denoted by italic text. John Jay (1789-95), James Wilson (1789-98), John Rutledge (1789-91), William Cushing (1790-93), John Blair (1793-94), James Iredell (1790-99), Thomas Johnson (1762-93), William Paterson (1763-1806), Samuel Chase (1766-1811), Oliver Ellsworth (1795-1800), Bushrod Washington (1799-1829), Alfred Moore (1800-04), John Marshall (1801-35), William Johnson (1804-34).
Does this evidence support the claim?
 Completely support Partially support Refute Irrelevant
- Evidence 2:** "United States Senate. This article incorporates public domain material from Justice TSB to Present. United States Supreme Court. before Justice Cort. discussed as improvionly granted Grant, waite, remond Shadow discler In chambers opinion Current members Chief Justice : John Roberts Associate justices :
Does this evidence support the claim?
 Completely support Partially support Refute Irrelevant
- Evidence 3:** "The Supreme Court of the United States is the final court of appeal and final expounder of the Constitution of the United States. The justices are appointed by the president of the United States and must be confirmed by the U.S. Senate. The chief justice of the United States is also appointed by the president and must be confirmed by the Senate. The number of justices fluctuated in the 18th and 19th centuries until 1869, when it was set at nine. A plan by Pitt...
Does this evidence support the claim?
 Completely support Partially support Refute Irrelevant
- Evidence 4:** "The Supreme Court of the United States is the highest ranking judicial body in the United States. Established by Article III of the Constitution, the detailed structure of the Court was laid down by the 1st United States Congress in 1789. Congress specified the Court's original and appellate jurisdiction, created 13 judicial districts, and fixed the number of justices at six - one chief justice and five associate justices. The number of justices on the Supreme Court changed six times before settling at the present total of nine in 1869. [1] The following tables detail the succession of justices of the Supreme Court of the United States by seat.
Does this evidence support the claim?
 Completely support Partially support Refute Irrelevant
- Evidence 5:** "www.supremecourt.gov. Supreme Court of the United States. Archived from the original on December 7, 2020. Retrieved January 21, 2018. ^
Does this evidence support the claim?
 Completely support Partially support Refute Irrelevant

Below the evidence items, there are two questions:

- "Can you determine whether the claim is true or false by evidence above?"
 Yes No
- "Please search relevant evidence and add in the box below."
A text box contains the following text: "William J. Brennan (1906 - 1997) was the oldest justice on the United States Supreme Court in 1980. From table information from Wikipedia: https://en.wikipedia.org/wiki/List_of_Justices_of_the_Supreme_Court_of_the_United_States. We can see that in 1980, William J. Brennan is the oldest alive justice."
Does this claim need to be corrected?
 Yes No

At the bottom, there is a "COPY" button, a "self retrieved" label, a "SAVE" button, and a "DOWNLOAD GENERATED JSON" button.

Figure 12: Screenshot of the second annotation interface: Evidence stance identification and Claim correction.

G.3 Claim Merge and Deduplication

Who was the oldest justice on the US supreme court in 1980?

1. In 1980, the oldest justice on the United States Supreme Court was Justice William O. Douglas.
2. He was born on October 16, 1898, and served on the Supreme Court from 1939 until his retirement in 1975.
3. Therefore, in 1980, Justice Douglas was still alive and would have been the oldest serving justice on the Court at that time.
 - a. Therefore, in 1980, Justice Douglas was still alive.
 - b. William J. Brennan would have been the oldest serving justice on the Court in 1980.

Please upload a generated result json.

0_screen3.json

Please select a sentence by clicking it on the left.

Sentence 1:

In 1980, the oldest justice on the United States Supreme Court was William J. Brennan.

Sentence 2:

Justice William O. Douglas was born on October 16, 1898 and served on the Supreme Court from 1939 until his retirement in 1975.

Sentence 3:

Therefore, in 1980, Justice Douglas was still alive, but William J. Brennan would have been the oldest serving justice on the Court in 1980.

Please Merge All Revised Sentences, duplicate and make coherent:

In 1980, the oldest justice on the United States Supreme Court was William J. Brennan. Justice William O. Douglas was born on October 16, 1898 and served on the Supreme Court from 1939 until his retirement in 1975. Therefore, in 1980, Justice Douglas was still alive, but William J. Brennan would have been the oldest serving justice on the Court in 1980.

Figure 13: Screenshot of the third annotation interface: Claim Merge and Deduplication.