

Large Language Models Help Humans Verify Truthfulness— Except When They Are Convincingly Wrong

Anonymous ACL submission

Abstract

Large Language Models (LLMs) are increasingly used for accessing information on the web. Their truthfulness and factuality are thus of great interest. To help users make the right decisions about the information they get, LLMs should not only provide information but also help users fact-check it. We conduct human experiments with 80 crowdworkers to compare language models with search engines (information retrieval systems) at facilitating fact-checking. We prompt LLMs to validate a given claim and provide corresponding explanations. Users reading LLM explanations are significantly more efficient than those using search engines while achieving similar accuracy. However, they over-rely on the LLMs when the explanation is wrong. To reduce over-reliance on LLMs, we ask LLMs to provide contrastive information—explain both why the claim is true and false, and then we present both sides of the explanation to users. This contrastive explanation mitigates users’ over-reliance on LLMs, but cannot significantly outperform search engines. Further, showing both search engine results and LLM explanations offers no complementary benefits compared to search engines alone. Taken together, our study highlights that natural language explanations by LLMs may not be a reliable replacement for reading the retrieved passages, especially in high-stakes settings where over-relying on wrong AI explanations could lead to critical consequences.

1 Introduction

Imagine you are told a claim about Neptune: “*Only one spacecraft has visited Neptune and it has more than 13 moons.*” and you want to verify whether it is factual. What would you do—look up relevant pages from search engines or ask ChatGPT for its take? This is not just a question of checking a piece of trivia; our information ecosystem depends on

people being able to check the veracity of information online. Misinformation, whether accidental or deliberate, has the potential to sway public opinion, influence decisions, and erode trust in credible sources (Faris et al., 2017; Mendes, 2017). Moreover, the wide adoption of large language models like ChatGPT increases the danger of misinformation, both by malicious actors and models generating inadvertent hallucinations (Pan et al., 2023).

Consequently, verifying the accuracy of information has taken on great importance. Fact-checking claims is a well-established task in NLP (Thorne et al., 2018; Guo et al., 2021). However, automated fact-checkers are far from perfect, and they are only useful when users trust their predictions (Nakov et al., 2021). Building that trust and providing effective help is crucial: a team without trust leads to suboptimal human-AI team performance while over-trusting wrong AI predictions could lead to catastrophic failures in high-stakes applications. Therefore, in real-life applications, we care about the AI-assisted human accuracy of fact-checking, rather than evaluating and improving automated fact-checkers alone (Shneiderman, 2022).

The two major types of tools for helping human users (many of which are non-expert fact-checkers) are *retrieval* and *explanation* (Nakov et al., 2021), exemplified by the widely-used web search engines (*e.g.*, Google) and generative language models (*e.g.*, ChatGPT) respectively. Showing retrieved passages to users has long been established as an effective information-seeking tool (Vlachos and Riedel, 2014). In contrast, the usefulness of generative explanations on fact-checking remains understudied. On the one hand, competent generative models (especially LLMs) can generate fluent and convincing-looking natural language explanations that not only provide an answer (*i.e.*, whether the claim is true or false), but also elucidate the context and basis of its judgment. On the other hand, these models are prone to hallucinations (Min et al., 2023; Liu et al.,

041
042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081

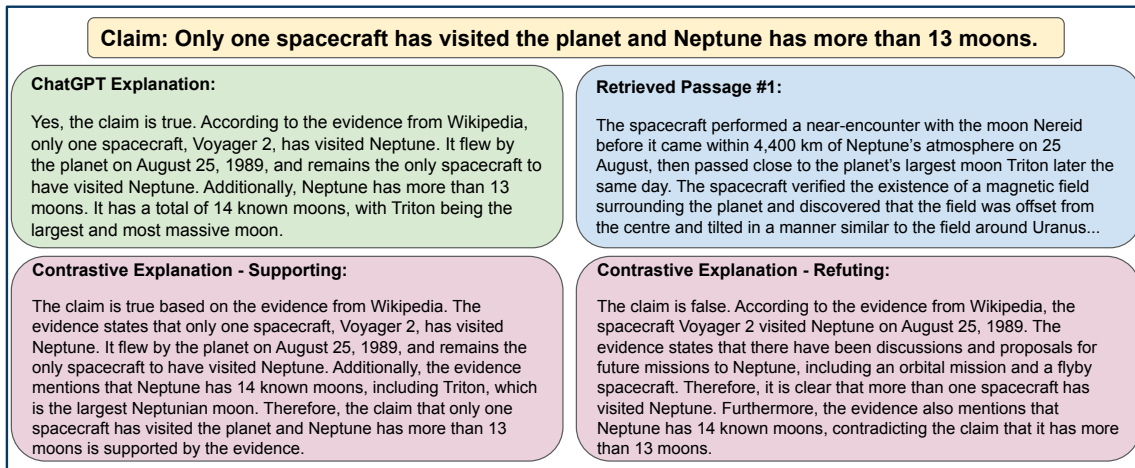


Figure 1: An example claim and the corresponding ChatGPT explanation, retrieved passages (abridged), and contrastive explanation. The claim is true and the refuting explanation has factual errors and reasoning contradiction.

2023), so the users are frequently left to their own devices.

In this work, we conduct human experiments to study *whether language models can assist fact-checking*. To contextualize the effectiveness of explanations against search engines, we compare them with retrieval models mimicking a search engine experience and experiment with ways where retrieved passages can be paired with explanations, aiming to provide a practical guide to users on what is the most helpful tool. We base our evaluation on FoolMeTwice (Eisenschlos et al., 2021), an adversarial dataset with interesting claims crowdsourced and gold evidences from Wikipedia (Eisenschlos et al., 2021). Our participants verify whether the claim is factually true or false: Figure 1 shows an example to illustrate the explanation and retrieved passages that participants see.

Our study reveals that showing explanation and retrieved passages lead to similar human accuracy (74% and 73% respectively) on difficult-to-verify claims (59% without AI assistant), but reading natural language explanations is significantly faster (1.01 min/claim vs 2.53 min/claim). However, humans over-trust ChatGPT explanations where they agree with the explanation most of the time, even when the explanation gives a wrong answer.

To combat the issue of over-reliance on natural language explanations, we explore two improvements: 1) contrastive explanations—present both supporting and refuting arguments generated by ChatGPT to the user and 2) combining retrieval and explanation (showing both to users). Both methods significantly reduce over-reliance on wrong AI explanations, however, they do not show a significant gain in user fact-checking accuracy compared to just showing users the retrieved passages. Overall,

our work underscores the potential benefit and danger of natural language explanations as a tool in the battle against misinformation. They can save time, but at the same time the difficulty of combatting over-reliance and the redundancy when combining retrieval and explanation remains. Turning back to the question of what users should do to verify factuality: taking longer time to read the retrieved passages is still more reliable!

2 Related Work

2.1 Fact Checking

Abundant datasets have been collected for training and evaluating automatic fact-checking models, such as FEVER (Thorne et al., 2018; Schuster et al., 2021; Guo et al., 2021) and SciFact (Wadden et al., 2020). Various techniques have been proposed to improve the fact-checking pipeline, such as jointly reasoning across evidence articles and claims (Popat et al., 2018b), and breaking complex claims into atomic sub-claims (Min et al., 2023; Kamoi et al., 2023). Instead of improving automatic fact-checking, we focus on how to improve human fact-checking via user studies.

Compared to automated approaches, there are relatively few prior user studies. Notably, Fan et al. (2020) synthesized summaries for retrieved passages to improve efficiency for users and Robbmond et al. (2022) compared showing explanations in different modalities to users. However, the advent of LLMs such as ChatGPT made it possible to generate plausible natural language explanations, and we are the first work to systematically evaluate such explanations in comparison to conventional retrieval methods.

153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183

184

185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201

2.2 Explainable AI

A thread of work in explainable AI (XAI) attempts to generate useful explanations in various formats (Wiegrefe and Marasović, 2021), such as highlighting (Schuff et al., 2022), feature importance (Ribeiro et al., 2016), free-text rationales (Ehsan et al., 2018), and structured explanations (Lamm et al., 2020). As the end goal of explanations is to aid human verification of AI predictions and inform decision-making (Vasconcelos et al., 2022; Fok and Weld, 2023), several work in XAI literature has focused on human-centered evaluation of explanations (Poursabzi-Sangdeh et al., 2021). Closest to our work, Feng and Boyd-Graber (2018) evaluated human-AI collaborative Quizbowl question answering and compared the effectiveness of showing retrieved passages, highlighting, and showing multiple guesses made by the system. This previous work used only a retrieval component, while our new approach allows us to directly compare ChatGPT-generated explanations (in the form of free-text rationales) with retrieved passages for aiding claim verification and explore whether natural language explanations and retrieved evidence yield complementary benefits. Joshi et al. (2023) studied free-text explanations in question-answering setting: their rationales do not help users much, especially when the rationales are misleading. In contrast to their work, we *contrast* model-generated explanations with passages retrieved from external sources (Wikipedia).

2.3 Trust Calibration and Over-Reliance

Existing work has identified the issue of human over-reliance on AI predictions, where humans tend to trust AI predictions *even when they are wrong* (Bussone et al., 2015b; Lai et al., 2021). A growing line of work attempts to mitigate such over-reliance, for example by providing explanations (Bansal et al., 2020; Zhang et al., 2020; Vasconcelos et al., 2022), communicating model uncertainty (Prabhudesai et al., 2023; Si et al., 2022), showing AI model accuracy (Yin et al., 2019), and prompting slow thinking (Buçinca et al., 2021) to help users calibrate their trust. Our work also contributes to this line of work by revealing the over-reliance issue in fact-checking. We propose new ways of potentially combatting over-reliance including contrastive explanation and combining explanation with retrieval.

3 Research Questions

To understand the comparative advantages of retrieval and explanation in human fact verification, we pose the following research questions:

- **RQ1:** Are natural language explanations more effective than retrieved passages for human fact-checking?
- **RQ2:** Can contrastive explanations—arguing for or against a fact being true—mitigate over-reliance and be more effective than non-contrastive explanations?
- **RQ3:** Are there complementary benefits in presenting both natural language explanations and retrieved passages?

We investigate these questions through a series of human studies: we show participants claims that need to be verified, potentially aiding them with different pieces of evidence (Figure 1). This is a between-subjects study; thus, we vary the evidence presented to participants in different conditions:

- **Baseline:** We show users only the claims without any additional evidence.
- **Retrieval:** We show the top 10 paragraphs retrieved from Wikipedia along with the claim to be verified.
- **Explanation:** We show the ChatGPT¹ explanation along with the claim.
- **Contrastive Explanation:** We present users ChatGPT’s supporting and refuting arguments side by side.
- **Retrieval + Explanation:** We present both the retrieved passages as well as the (non-contrastive) natural language explanations to users.

In the **Explanation** and **Retrieval + Explanation** conditions, the ChatGPT prediction on whether the claim is true or false is part of the explanation, while in the other conditions, users only see the evidence but not the prediction.

4 Study Design Overview

4.1 Task, Data, and Variables

We ask human annotators to look at claims and decide whether it is true or false. We use the FoolMeTwice dataset (Eisenschlos et al., 2021) over other claim-verification datasets because FoolMeTwice is adversarial: crowdworkers write claims based on Wikipedia to maximally fool another set of annotators whose task is to verify these claims. This ensures that all the claims are hard to

¹We use gpt-3.5-turbo in all experiments.

verify, mimicking potential real-world fake news arms race. For our human studies, we create a test set by randomly sampling 200 claims where half are true and half are false. To ensure that the selected claims are sufficiently complex, we only sample claims requiring at least two different sentences from Wikipedia to verify.

We sample 20 claims (half true and half false) for each participant to verify and randomize their order. For each claim, we ask for the participant’s binary decision of whether they think the claim is true or false. We measure the accuracy of human decisions given that we know the gold labels of these claims. We also ask for the participant’s confidence in their judgment on a scale of 1 to 5, and record the time used for verifying each claim. We also ask for a free-form response of how the annotator makes their judgments. Appendix A.1 and Figure 6 illustrate the interface setup.

4.2 Retriever

For the **Retrieval** and **Retrieval + Explanation** conditions, we show users the most relevant passages from Wikipedia. Specifically, we adopt a similar retrieval setup as Min et al. (2023), where we use the state-of-the-art Generalizable T5-based Retriever (GTR-XXL), an unsupervised dense passage retriever (Ni et al., 2021). We retrieve the top 10 most relevant paragraphs from Wikipedia, where each paragraph has an average length of 188 words. To measure the retrieval quality, we report two metrics on our test set. The full recall measures how often the top 10 retrieved passages contain all evidence sentences required to verify the claim, which is 81.5%; the partial recall measures how often the top 10 retrieved passages contain at least one evidence sentence required to verify the claim, which is 93.0%.

4.3 Explanation Generation

We study two types of natural language explanations with ChatGPT: non-contrastive explanation and contrastive explanation. In the **Explanation** and **Retrieval + Explanation** conditions, we generate **non-contrastive** explanations, where we construct the prompt by concatenating the top 10 retrieved passages, followed by the claim to be verified, then appending the question “*Based on the evidence from Wikipedia, is the claim true? Explain in a short paragraph.*” We measure the accuracy of these explanations by manually extracting the answer (true or false) from the explanations and com-

paring with the gold labels. ChatGPT-generated explanations achieve an accuracy of 78.0% (judged based on the AI predictions only, not the reasoning processes). In the **Contrastive Explanation** condition, we prompt ChatGPT to generate both a supporting answer and a refuting answer. Specifically, after concatenating the retrieved passages and the claim, we append two different questions: 1) “*Based on the evidence from Wikipedia, explain in a short paragraph why the claim is true.*” and 2) “*Based on the evidence from Wikipedia, explain in a short paragraph why the claim is false.*” We then show both of these generated explanations to annotators, which functions similarly to a single-turn debate (Parrish et al., 2022; Michael et al., 2023).

Additionally, in **Retrieval + Explanation**, we automatically insert citations to the explanation text to attribute the arguments to corresponding retrieved passages. This is implemented by prompting ChatGPT where we provide a manually crafted example of inserting citations into the explanations based on the retrieved passages, which has been shown to be an effective method for enabling citations in language model generations (Gao et al., 2023). For all cases, we ground the explanation generation on the retrieved passages. This is because grounding significantly improves the accuracy of explanations. For example, for non-contrastive explanations, grounding improves the accuracy from 59.5% to 78.0%. For all cases, we use a temperature value of 0 for ChatGPT generation to minimize randomness.

4.4 Users

We recruit participants from Prolific. We recruit 16 annotators for each condition and each annotator verifies 20 claims, resulting in $20 \times 16 \times 5 = 1500$ annotations across all 5 conditions. We compensate all annotators at least \$14 per hour, as well as additional bonuses to users who perform particularly well on the task or who have left very insightful comments as an additional incentive. Our experiment is approved by an IRB.

5 Experiment Results

5.1 RQ1: Are natural language explanations more effective than retrieved passages for human fact checking?

We compare three conditions: the Baseline condition (showing users only the claims); the Retrieval condition (showing the top 10 para-

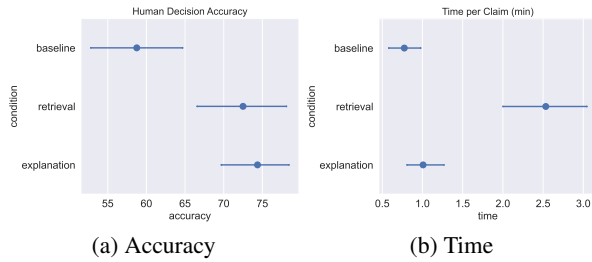
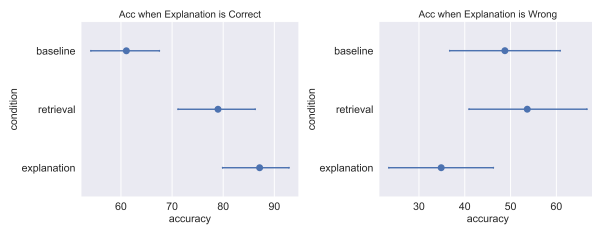


Figure 2: Human decision accuracy and average time spent on verifying a claim. Both retrieval and explanation significantly improve human verification accuracy, while explanation takes a significantly shorter time.



(a) Human decision accuracy on examples where the explanation is correct. (b) Human decision accuracy on examples where the explanation is wrong.

Figure 3: Human verification accuracy broken down into two subsets: examples on which the explanation gives the correct labels, and examples on which the explanation gives the wrong labels. Humans over-rely on explanations so that they achieve significantly lower accuracy than the baseline when the explanation is wrong.

graphs retrieved from Wikipedia); and the Explanation condition (showing the ChatGPT explanation along with the claim). We do not set a time limit but record the time taken for each claim.

Figure 2a shows the AI-assisted human verification accuracy across conditions. We test the significance of our results using Student’s t-tests with Bonferroni correction.² We start with examining whether ChatGPT explanations and retrieved passages are indeed helpful for humans.

Showing ChatGPT explanation improves human accuracy. When showing explanations to users, the accuracy is $\mu = 0.74 \pm \sigma = 0.09$ compared to the baseline condition where claims are shown without any additional evidence (0.59 ± 0.12). The improvement in accuracy is significant ($z = -4.08, p = 0.00015$).

Showing retrieved passages improves human accuracy. When showing retrieved passages to users, they achieve the accuracy of (0.73 ± 0.12) as compared to the baseline condition where claims

²We inspected all data with histograms and Q-Q plots to verify that the data approximate normality before applying t-tests.

are shown without any additional evidence (0.59 ± 0.12). The improvement in accuracy is significant ($z = -3.15, p = 0.0018$). Now that both ChatGPT explanation and retrieved passages help humans more accurately verify claims, we examine their comparative advantages in both accuracy and time.

Showing ChatGPT explanation does not achieve significantly higher accuracy than showing retrieved passages. Comparing the accuracy in the explanation condition (0.74 ± 0.09) and the retrieval condition (0.73 ± 0.12), the improvement in accuracy is not significant ($z = -0.48, p = 0.32$).

However, **reading ChatGPT explanation is significantly faster than reading retrieved passages.** We compare the time taken to verify claims in Figure 2b. When verifying with retrieved passages, the time taken to verify each claim is (2.53 ± 1.07) minutes while for the explanation condition, it takes (1.01 ± 0.45) minutes. Showing explanations allows significantly faster decision time than showing retrieved passages ($z = -5.09, p = 9.1e - 6$).

5.2 Breakdown Analysis: The Danger of Over-Reliance

While ChatGPT explanations show promise in aiding human fact verification, the aggregate results obscure the danger when the explanation gives wrong answers. To examine what happens in those cases, we break down the analysis, manually annotating the ChatGPT explanation for each claim based on whether it gives the correct answer (whether the claim is true or false). We then split all user responses into two subsets: ones with correct answers from ChatGPT and ones where the ChatGPT explanation is wrong (Figure 3a and Figure 3b, respectively).

Users achieve the highest accuracy when the explanations are correct, but below-random accuracy when the explanations are wrong. When the explanation is correct, users’ accuracy is (0.87 ± 0.13), higher than the baseline of having no evidence (0.61 ± 0.13) as well as the retrieval condition (0.79 ± 0.15). However, when the explanation is wrong, users tend to over-trust the explanations and only achieve an accuracy of (0.35 ± 0.22) as compared to the baseline condition (0.49 ± 0.24) and the retrieval condition (0.54 ± 0.26). Moreover, **users spend similar time on claims with correct and wrong explanations**, further indicating that they are not deliberately differentiating correct and wrong explanations and instead tend to trust most of the explanations. We also look at the

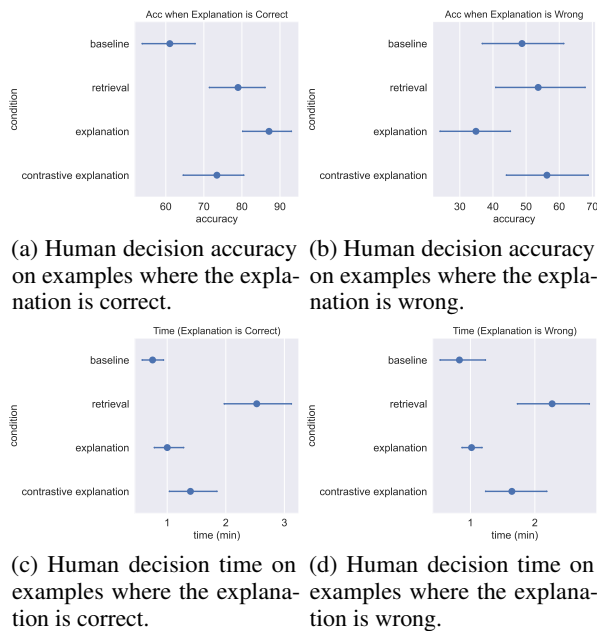


Figure 4: Verification accuracy and time broken down by whether the (non-contrastive) explanation is correct. Contrastive explanation achieves significantly higher accuracy than non-contrastive explanation on examples where the non-contrastive explanation is wrong, with some drop in accuracy on examples where the non-contrastive explanation is correct.

free-form responses from users for their decision rationales, the most common responses include: (1) ChatGPT’s explanation looks convincing, especially with quotes from the retrieved passages (even when the quotes or reasoning are wrong); (2) They do not have any prior knowledge on the topic so would just trust ChatGPT.

In comparison, retrieved passages suffer less from over-reliance. On examples where the ChatGPT explanations are correct, the retrieval condition achieves the accuracy of (0.79 ± 0.15) , surpassing the baseline condition (0.61 ± 0.13) . On examples where the ChatGPT explanations are wrong, the retrieval condition achieves the accuracy of (0.54 ± 0.26) compared to the baseline (0.49 ± 0.24) . While there is still an accuracy drop in these examples, possibly because they are harder to verify, the performance discrepancy between the two cases (ChatGPT explanation correct vs wrong) is much less severe in the retrieval condition. This highlights the pitfall of using ChatGPT explanation to aid helpful verification: users over-rely on the explanations, even when they are wrong and misleading. To combat this problem, we next explore two strategies for mitigation: contrastive explanation and combining retrieval and explanation.

5.3 RQ2: Can contrastive explanations mitigate over-reliance and be more effective than non-contrastive explanations?

In addition to the three conditions from the previous section (Baseline, Retrieval, and Explanation), we additionally compare the Contrastive Explanation condition where we present users ChatGPT’s supporting and refuting arguments side by side. The experiment results are in Figure 4. We first compare contrastive explanation with non-contrastive explanation.

Contrastive explanation achieves higher human accuracy than non-contrastive explanation when the non-contrastive explanation is wrong. When the non-contrastive explanation is wrong, humans only achieve an accuracy of (0.35 ± 0.22) due to over-reliance, but when switching to contrastive explanation improves the accuracy to (0.56 ± 0.24) , which is significantly higher ($z = -2.52, p = 0.009$). When analyzing the free-response rationales of human judgment, the most common patterns of how people make correct judgments based on contrastive explanations are: (1) The correct side of the explanation is more compelling or thorough; (2) The wrong side of the explanation contains factual errors and wrong reasoning; (3) Both sides of the explanations give the same answer (even though ChatGPT was prompted to explain why the claim is true and false in the two sides of explanations).

However, **contrastive explanation achieves lower human accuracy than non-contrastive explanation when the non-contrastive explanation is correct.** When the non-contrastive explanation is correct, humans achieve an accuracy of (0.87 ± 0.13) , higher than contrastive explanation (0.73 ± 0.15) , indicating a significant drop ($z = -2.56, p = 0.008$). Unlike the case in non-contrastive explanations where users can just take the AI prediction as the answer, for contrastive explanations they have to decide between the two sides of the explanation and choose one, which can sometimes be difficult since LLMs can generate convincing explanations even for the wrong statements. For example, given the false claim “Joe Torre was the manager of the New York Yankees and guided the team to four World Series championships, and ranks third all-time in MLB history with 2,326 wins as a manager.”, ChatGPT generates the supporting explanation “Yes, the claim is

498 *true. According to the evidence from Wikipedia, Joe*
 499 *Torre was the manager of the New York Yankees*
 500 *from 1996 to 2007. He also ranks third all-time in*
 501 *MLB history with 2,326 wins as a manager.” and*
 502 *generates the refuting explanation “The claim is*
 503 *false. According to the evidence from Wikipedia,*
 504 *Joe Torre was the manager of the New York Yan-*
 505 *kees and guided the team to six pennants and four*
 506 *World Series championships. He ranks fifth all-time*
 507 *in MLB history with 2,326 wins as a manager, not*
 508 *third.” The fact is that Torre ranks fifth all-time*
 509 *in MLB history with 2,326 wins as a manager but*
 510 *ChatGPT still generated a convincingly looking*
 511 *explanation for the wrong side by hallucinating he*
 512 *ranks third all-time rather than fifth. As a result,*
 513 *some users were misled into the wrong judgment.*
 514 *Overall, contrastive explanation shows promise in*
 515 *reducing over-reliance but incurs a trade-off in ac-*
 516 *curacy when the non-contrastive explanation is cor-*
 517 *rect. Next, we also compare contrastive explana-*
 518 *tion with retrieval.*

519 **Contrastive explanation does not significantly**
 520 **improve human accuracy over retrieval.** On ex-
 521 amples where the non-contrastive explanation is
 522 correct, providing contrastive explanation achieves
 523 the human accuracy of (0.73 ± 0.15) , lower than
 524 the accuracy in the retrieval condition (0.79 ± 0.15) .
 525 On examples where the non-contrastive explana-
 526 tion is wrong, contrastive explanation has com-
 527 parable human accuracy of (0.56 ± 0.24) com-
 528 pared to retrieval (0.54 ± 0.26) , and the difference
 529 is not significant ($z = 0.29, p = 0.61$). There-
 530 fore, in both cases, contrastive explanations do not
 531 achieve significantly better human accuracy than
 532 retrieval, despite the results that contrastive expla-
 533 nations can mitigate over-reliance as compared to
 534 non-contrastive explanations.

535 Apart from the above quantitative results, we
 536 also manually analyze the free-form responses of
 537 user decision rationales to understand how users
 538 leverage contrastive explanations to make deci-
 539 sions. **Users mostly base their judgment on the**
 540 **relative strength of the two sides of the expla-**
 541 **nations** (*i.e.*, is the supporting or refuting expla-
 542 nation more convincing) (41.8%). Example user
 543 rationales include: *“The refutation seems more*
 544 *logically sound.”* and *“The support explanation*
 545 *seems like it’s trying too hard to make the claim*
 546 *true, but the refute puts it more plain and simple*
 547 *and makes more sense.”* **Sometimes both sides**
 548 **converge on the same answer** (26.9%) and users
 549 would just agree with that. For example, for the

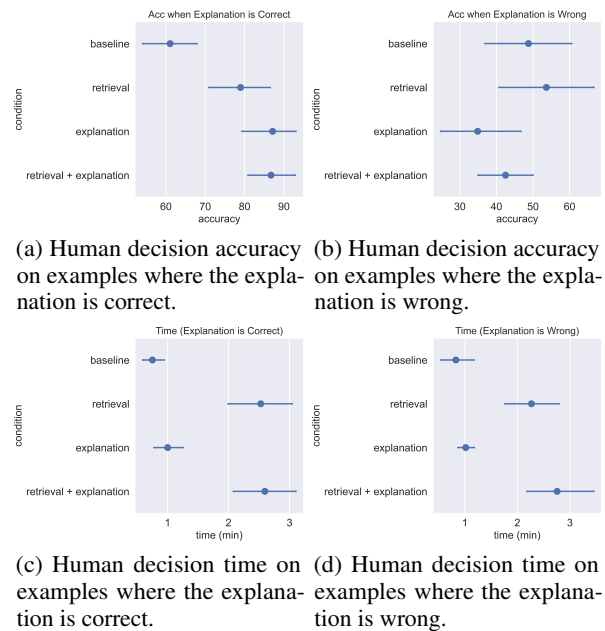


Figure 5: Verification accuracy and time breakdown. Combining retrieval and explanation is not significantly better than just showing retrieved passages alone.

550 false claim *“The only verified original sled prop*
 551 *from Citizen Kane was sold at a price of over a*
 552 *hundred thousand dollars.”*, users found that *“Both*
 553 *sides acknowledge that there were more than 1 sled*
 554 *prop, therefore refuting the claim.”*, even though the
 555 ChatGPT supporting explanation said *“The claim*
 556 *is true.”* In several cases, ChatGPT would simply
 557 say the claim is true even though we prompt it for a
 558 refuting explanation (and vice versa), giving users
 559 a clear cue that the model could not make a strong
 560 argument for the wrong side.

5.4 RQ3: Are there complementary benefits in presenting both natural language explanations and retrieved passages?

561 Apart from the Baseline, Retrieval, and
 562 Explanation conditions from earlier, we also com-
 563 pare with the (Retrieval + Explanation) condi-
 564 tion where we present both to users.
 565

5.5 Results

566 Results are plotted in Figure 5 and we start by
 567 comparing whether combining explanation with
 568 retrieval is better than explanation alone.
 569

570 **Combining retrieval and explanation does not**
 571 **achieve significantly higher accuracy than expla-**
 572 **nation alone in cases where the explanation is**
 573 **correct.** When the explanation is correct, users
 574 achieve the accuracy of (0.87 ± 0.13) relying on
 575 explanations, as compared to combining both re-
 576

retrieval and explanation (0.87 ± 0.12). We do not observe a significant advantage of combining retrieval and explanation in this case ($z = 0.084, p = 0.53$).

Combining retrieval and explanation does not achieve significantly higher accuracy than explanation alone in cases where the explanation is wrong either. When the explanation is wrong, users’ accuracy (0.35 ± 0.22) in the explanation condition is slightly lower than combining retrieval and explanation (0.43 ± 0.16). The advantage of combining retrieval and explanation is not significant ($z = -1.06, p = 0.15$). Taken together, combining explanation and retrieval is not better than explanation alone. Next, we compare whether combining explanation with retrieval is better than retrieval alone.

Combining retrieval and explanation does not achieve significantly higher accuracy than retrieval alone in cases where the explanation is correct. When the explanation is correct, users achieve the accuracy of (0.79 ± 0.15) in the retrieval alone condition as compared to combining both retrieval and explanation (0.87 ± 0.12). There is a slight advantage of combining retrieval and explanation in this setting but the advantage is not significant ($z = -1.48, p = 0.07$).

Combining retrieval and explanation does not achieve significantly higher accuracy than retrieval alone in cases where the explanation is wrong. When the explanation is wrong, users’ accuracy of (0.54 ± 0.26) in the retrieval alone condition beats combining both retrieval and explanation (0.43 ± 0.16), indicating a drop in accuracy in this case when combining retrieval and explanation. This means that combining retrieval and explanation offers no complementary benefits compared to retrieval alone. To understand whether users indeed read both the explanation and retrieved passages, we compare their reading time.

Combining retrieval and explanation takes a longer time. In the retrieval alone condition, users take (2.5 ± 1.1) minutes to verify a claim; in the explanation condition, users take (1.0 ± 0.4) minutes to verify a claim; in the retrieval + explanation condition, users take (2.7 ± 1.0) minutes to verify a claim, indicating that combining retrieval and explanation increases the verification time, so users indeed spend time reading the explanation and retrieved passages in most cases. Moreover, in analyzing the free-form responses, the majority of the users base their judgment on the retrieved passages since the ChatGPT explanations are not

always credible, further indicating that presenting ChatGPT explanations grounded on the retrieved passages does not really offer additional benefits than just presenting the retrieved passages themselves. Overall, our results suggest that combining retrieval and explanation might be redundant and inefficient.

5.6 Meta-Analysis

We also conduct a series of meta-analyses and summarize the main findings below. We refer readers to Appendix A.2 for more details.

- **Confidence Calibration:** Users are overconfident on wrong judgments across all experiment conditions, with average confidence above 0.6 (Appendix A.2.1).
- **Impact of Retrieval:** The explanation accuracy is much lower when the retrieval recall is low, and the human decision accuracy is also much lower when the retrieval recall is low (Appendix A.2.2).
- **Correlation between Accuracy and Time:** The correlation is weak in all conditions (Appendix A.2.3).
- **Analysis of Free-form Responses:** We categorize and qualitatively analyze when users disagree with ChatGPT explanations, which mostly happens when they detect ChatGPT’s self-contradictions, identify evidence from retrieved passages, or just rely on their own knowledge (Appendix A.2.4).
- **Additional Related Work:** We discuss additional related works from NLP and HCI in Appendix A.3.

6 Conclusion

Our human studies highlight the false promise of using natural language explanation produced by ChatGPT to aid human fact-checking. Humans over-rely on explanations even when they are wrong, making human accuracy worse than showing retrieval or the baseline of not showing any evidence. In attempts to combat over-reliance, contrastive explanation mitigates users’ over-reliance on wrong explanations, while combining retrieval and explanation does not achieve significant complementary improvement. Overall, neither of these two approaches significantly outperforms the retrieval baseline. highlighting the need for better methods for combatting over-reliance on AI.

678 Limitations

679 We acknowledge several limitations of this work:
680 (1) Our experiments are at a limited scale with
681 participants recruited from Prolific. It is possible
682 that other factors such as knowledge of the top-
683 ics, familiarity with language models, and trust in
684 automation in general, could impact our conclu-
685 sions and future work should consider scaling up
686 the study with diverse populations to capture such
687 nuances.

688 (2) We only experimented with a limited set of
689 explanation methods and our explanations are all
690 static (*i.e.*, not personalized for different partic-
691 ipants). Future work could explore how to cus-
692 tomize the best sets of evidence for different users
693 in different conditions (Feng and Boyd-Graber,
694 2022; Bansal et al., 2020).

695 (3) We observed little benefit from combining re-
696 trieval and explanation, future work could further
697 explore how to strategically combine retrieval and
698 explanation so that they can actually complement
699 each other in both accuracy and efficiency. For in-
700 stance, when the explanation is likely to be correct,
701 we can show users the explanation; but when the
702 explanation is likely to be wrong, we should priori-
703 tize showing users the raw retrieved passages. This
704 might also require better uncertainty estimation or
705 calibration to help users identify AI mistakes.

706 Ethical Considerations

707 In our human studies, we made sure to compen-
708 sate all participants fairly, with a minimum rate
709 of \$14 per hour. We do not expect any potential
710 mental stress or harm to the participants from the
711 experiments. Our work highlights and explores so-
712 lutions for combatting human over-reliance on AI,
713 which has important societal implications given
714 that LLMs like ChatGPT are being widely used.
715 We hope our results can contribute positively to
716 society by reducing catastrophic harms caused by
717 such over-reliance and also offering practical guid-
718 ance for how to effectively verify potential fake
719 information on the Internet.

References

- 720 Yasmeeen Alufaisan, Laura R. Marusich, Jonathan Z.
721 Bakdash, Yan Zhou, and Murat Kantarcioglu. 2021.
722 [Does explainable artificial intelligence improve hu-
723 man decision-making?](#) *Proceedings of the AAAI Con-
724 ference on Artificial Intelligence*, 35(8):6618–6626.
725
- Gagan Bansal, Tongshuang Sherry Wu, Joyce Zhou,
726 Raymond Fok, Besmira Nushi, Ece Kamar,
727 Marco Tulio Ribeiro, and Daniel S. Weld. 2020.
728 [Does the whole exceed its parts? the effect of ai
729 explanations on complementary team performance.](#)
730 *Proceedings of the 2021 CHI Conference on Human
731 Factors in Computing Systems*.
732
- Andrew Beers. 2019. [Examining the digital toolsets of
733 journalists reporting on disinformation.](#)
734
- Samuel R Bowman, Jeeyoon Hyun, Ethan Perez, Edwin
735 Chen, Craig Pettit, Scott Heiner, Kamile Lukosuite,
736 Amanda Askeell, Andy Jones, Anna Chen, et al. 2022.
737 [Measuring progress on scalable oversight for large
738 language models.](#) *arXiv preprint arXiv:2211.03540*.
739
- Zana Bućinca, Maja Barbara Malaya, and Krzysztof Z.
740 Gajos. 2021. [To trust or to think: Cognitive forcing
741 functions can reduce overreliance on ai in ai-assisted
742 decision-making.](#) *Proc. ACM Hum.-Comput. Inter-
743 act.*, 5(CSCW1).
744
- Adrian Bussone, Simone Stumpf, and Dympna
745 O’Sullivan. 2015a. [The role of explanations on trust
746 and reliance in clinical decision support systems.](#) In
747 *2015 International Conference on Healthcare Infor-
748 matics*, pages 160–169.
749
- Adrian Bussone, Simone Stumpf, and Dympna
750 O’Sullivan. 2015b. [The role of explanations on
751 trust and reliance in clinical decision support sys-
752 tems.](#) *2015 International Conference on Healthcare
753 Informatics*, pages 160–169.
754
- Jifan Chen, Grace Kim, Aniruddh Sriram, Greg Dur-
755 rett, and Eunsol Choi. 2023. [Complex claim veri-
756 fication with evidence retrieved in the wild.](#) *ArXiv*,
757 abs/2305.11859.
758
- Chun-Wei Chiang and Ming Yin. 2021. [You’d bet-
759 ter stop! understanding human reliance on machine
760 learning models under covariate shift.](#) In *Proceed-
761 ings of the 13th ACM Web Science Conference 2021*,
762 WebSci ’21, page 120–129, New York, NY, USA.
763 Association for Computing Machinery.
764
- Upol Ehsan, Brent Harrison, Larry Chan, and Mark O.
765 Riedl. 2018. [Rationalization: A neural machine trans-
766 lation approach to generating natural language ex-
767 planations.](#) In *Proceedings of the 2018 AAAI/ACM
768 Conference on AI, Ethics, and Society*, AIES ’18,
769 page 81–87, New York, NY, USA. Association for
770 Computing Machinery.
771
- Julian Martin Eisenschlos, Bhuwan Dhingra, Jannis Bu-
772 lian, Benjamin Borschinger, and Jordan L. Boyd-
773 Graber. 2021. [Fool me twice: Entailment from
774 wikipedia gamification.](#) *ArXiv*, abs/2104.04725.
775

776	Angela Fan, Aleksandra Piktus, Fabio Petroni, Guillaume Wenzek, Marzieh Saeidi, Andreas Vlachos, Antoine Bordes, and Sebastian Riedel. 2020. Generating fact checking briefs . In <i>Conference on Empirical Methods in Natural Language Processing</i> .	Vivian Lai, Chacha Chen, Qingzi Vera Liao, Alison Smith-Renner, and Chenhao Tan. 2021. Towards a science of human-ai decision making: A survey of empirical studies . <i>ArXiv</i> , abs/2112.11471.	829
777			830
778			831
779			832
780			
781	Robert Faris, Hal Roberts, Bruce Etling, Nikki Bourassa, Ethan Zuckerman, and Yochai Benkler. 2017. Partisanship, propaganda, and disinformation: Online media and the 2016 u.s. presidential election . <i>Social Science Research Network</i> .	Matthew Lamm, Jennimaria Palomaki, Chris Alberti, Daniel Andor, Eunsol Choi, Livio Baldini Soares, and Michael Collins. 2020. Qed: A framework and dataset for explanations in question answering . <i>Transactions of the Association for Computational Linguistics</i> , 9:790–806.	833
782			834
783			835
784			836
785			837
786			838
787	Shi Feng and Jordan L. Boyd-Graber. 2018. What can ai do for me?: evaluating machine learning interpretations in cooperative play . <i>Proceedings of the 24th International Conference on Intelligent User Interfaces</i> .	John D. Lee and Katrina A. See. 2004. Trust in automation: Designing for appropriate reliance . <i>Human Factors</i> , 46(1):50–80. PMID: 15151155.	839
788			840
789			841
790			
791	Shi Feng and Jordan L. Boyd-Graber. 2022. Learning to explain selectively: A case study on question answering . In <i>Conference on Empirical Methods in Natural Language Processing</i> .	Nelson F. Liu, Tianyi Zhang, and Percy Liang. 2023. Evaluating verifiability in generative search engines . <i>ArXiv</i> , abs/2304.09848.	842
792			843
793			844
794			
795	Claudia Flores-Saviaga, Shangbin Feng, and Saiph Savage. 2022. Datavoidant: An ai system for addressing political data voids on social media . <i>Proceedings of the ACM on Human-Computer Interaction</i> , 6:1 – 29.	Ricardo Mendes. 2017. Troops, trolls and troublemakers: A global inventory of organized social media manipulation .	845
796			846
797			847
798			
799	Raymond Fok and Daniel S Weld. 2023. In search of verifiability: Explanations rarely enable complementary performance in ai-advised decision making . <i>arXiv preprint arXiv:2305.07722</i> .	Julian Michael, Salsabila Mahdi, David Rein, Jackson Petty, Julien Dirani, Vishakh Padmakumar, and Samuel R. Bowman. 2023. Debate helps supervise unreliable experts . <i>ArXiv</i> , abs/2311.08702.	848
800			849
801			850
802			851
803	Tianyu Gao, Ho-Ching Yen, Jiatong Yu, and Danqi Chen. 2023. Enabling large language models to generate text with citations . <i>ArXiv</i> , abs/2305.14627.	Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hanna Hajishirzi. 2023. Factscore: Fine-grained atomic evaluation of factual precision in long form text generation . <i>ArXiv</i> , abs/2305.14251.	852
804			853
805			854
806	Navita Goyal, Eleftheria Briakou, Amanda Liu, Connor Baumler, Claire Bonial, Jeffrey Micher, Clare R Voss, Marine Carpuat, and Hal Daumé III. 2023. What else do i need to know? the effect of background information on users’ reliance on ai systems . <i>arXiv preprint arXiv:2305.14331</i> .	Preslav Nakov, David P. A. Corney, Maram Hasanain, Firoj Alam, Tamer Elsayed, Alberto Barr’on-Cedeno, Paolo Papotti, Shaden Shaar, and Giovanni Da San Martino. 2021. Automated fact-checking for assisting human fact-checkers . <i>ArXiv</i> , abs/2103.07769.	855
807			856
808			857
809			858
810			859
811			860
812	Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. 2017. On calibration of modern neural networks . In <i>International Conference on Machine Learning</i> .	Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernandez Abrego, Ji Ma, Vincent Zhao, Yi Luan, Keith B. Hall, Ming-Wei Chang, and Yinfei Yang. 2021. Large dual encoders are generalizable retrievers . <i>ArXiv</i> , abs/2112.07899.	861
813			862
814			863
815			864
816	Zhijiang Guo, M. Schlichtkrull, and Andreas Vlachos. 2021. A survey on automated fact-checking . <i>Transactions of the Association for Computational Linguistics</i> , 10:178–206.	Yikang Pan, Liangming Pan, Wenhua Chen, Preslav Nakov, Min-Yen Kan, and William Yang Wang. 2023. On the risk of misinformation pollution with large language models . <i>ArXiv</i> , abs/2305.13661.	865
817			866
818			867
819			868
820	Brihi Joshi, Ziyi Liu, Sahana Ramnath, Aaron Chan, Zhewei Tong, Shaoliang Nie, Qifan Wang, Yejin Choi, and Xiang Ren. 2023. Are machine rationales (not) useful to humans? measuring and improving human utility of free-text rationales . <i>ArXiv</i> , abs/2305.07095.	Alicia Parrish, H. Trivedi, Ethan Perez, Angelica Chen, Nikita Nangia, Jason Phang, and Sam Bowman. 2022. Single-turn debate does not help humans answer hard reading-comprehension questions . <i>ArXiv</i> , abs/2204.05212.	869
821			870
822			871
823			872
824			873
825			874
826	Ryo Kamoi, Tanya Goyal, Juan Diego Rodriguez, and Greg Durrett. 2023. Wice: Real-world entailment for claims in wikipedia . <i>ArXiv</i> , abs/2303.01432.	Kashyap Popat, Subhabrata Mukherjee, Jannik Strotgen, and Gerhard Weikum. 2018a. Credeye: A credibility lens for analyzing and explaining misinformation . <i>Companion Proceedings of the The Web Conference 2018</i> .	875
827			876
828			877
			878
			879
			880

881	Kashyap Popat, Subhabrata Mukherjee, Andrew Yates, and Gerhard Weikum. 2018b. Declare: Debunking fake news and false claims using evidence-aware deep learning . <i>ArXiv</i> , abs/1809.06416.	overreliance on ai systems during decision-making . <i>Proceedings of the ACM on Human-Computer Interaction</i> , 7:1 – 38.	936 937 938
885	Forough Poursabzi-Sangdeh, Daniel G Goldstein, Jake M Hofman, Jennifer Wortman	Andreas Vlachos and Sebastian Riedel. 2014. Fact checking: Task definition and dataset construction . In <i>LTCSS@ACL</i> .	939 940 941
886	Wortman Vaughan, and Hanna Wallach. 2021. Manipulating and measuring model interpretability . In <i>Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems</i> , CHI '21, New York, NY, USA. Association for Computing Machinery.	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.	942 943 944 945 946 947 948
887			
888			
889			
890			
891			
892			
893	Snehal Prabhudesai, Leyao Yang, Sumit Asthana, Xun Huan, Qingzi Vera Liao, and Nikola Banovic. 2023. Understanding uncertainty: How lay decision-makers perceive and interpret uncertainty in human-ai decision making . <i>Proceedings of the 28th International Conference on Intelligent User Interfaces</i> .	Xinru Wang and Ming Yin. 2021. Are explanations helpful? a comparative study of the effects of explanations in ai-assisted decision-making . In <i>26th International Conference on Intelligent User Interfaces</i> , IUI '21, page 318–328, New York, NY, USA. Association for Computing Machinery.	949 950 951 952 953 954
894			
895			
896			
897			
898			
899	Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "why should i trust you?": Explaining the predictions of any classifier . In <i>Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining</i> , KDD '16, page 1135–1144, New York, NY, USA. Association for Computing Machinery.	Sarah Wiegrefe and Ana Marasović. 2021. Teach me to explain: A review of datasets for explainable natural language processing . In <i>NeurIPS Datasets and Benchmarks</i> .	955 956 957 958
900			
901			
902			
903			
904			
905			
906	Vincent Robbmond, Oana Inel, and Ujwal Gadiraju. 2022. Understanding the role of explanation modality in ai-assisted decision-making . <i>Proceedings of the 30th ACM Conference on User Modeling, Adaptation and Personalization</i> .	Ming Yin, Jennifer Wortman Vaughan, and Hanna M. Wallach. 2019. Understanding the effect of accuracy on trust in machine learning models . <i>Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems</i> .	959 960 961 962 963
907			
908			
909			
910			
911	Hendrik Schuff, Alon Jacovi, Heike Adel, Yoav Goldberg, and Ngoc Thang Vu. 2022. Human interpretation of saliency-based explanation over text . <i>Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency</i> .	Yunfeng Zhang, Qingzi Vera Liao, and Rachel K. E. Bellamy. 2020. Effect of confidence and explanation on accuracy and trust calibration in ai-assisted decision making . <i>Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency</i> .	964 965 966 967 968
912			
913			
914			
915			
916	Tal Schuster, Adam Fisch, and Regina Barzilay. 2021. Get your vitamin c! robust fact verification with contrastive evidence . In <i>North American Chapter of the Association for Computational Linguistics</i> .		
917			
918			
919			
920	Ben Shneiderman. 2022. <i>Human-Centered Artificial Intelligence</i> .		
921			
922	Chenglei Si, Chen Zhao, Sewon Min, and Jordan L. Boyd-Graber. 2022. Re-examining calibration: The case of question answering . In <i>Conference on Empirical Methods in Natural Language Processing</i> .		
923			
924			
925			
926	James Thorne and Andreas Vlachos. 2018. Automated fact checking: Task formulations, methods and future directions . <i>ArXiv</i> , abs/1806.07687.		
927			
928			
929	James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. Fever: a large-scale dataset for fact extraction and verification . <i>ArXiv</i> , abs/1803.05355.		
930			
931			
932			
933	Helena Vasconcelos, Matthew Jörke, Madeleine Grundle-McLaughlin, Tobias Gerstenberg, Michael Bernstein, and Ranjay Krishna. 2022. Explanations can reduce		
934			
935			

969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989

990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011

1012
1013
1014
1015
1016
1017

A Appendix

A.1 Interface Design

Figure 6 shows an example user interface for the Contrastive Explanation condition. We identify keywords as the non-stopwords in the claim and highlight them in the claims and explanations to aid reading (we also do keyword highlighting in the retrieved passages in the retrieval conditions). For the retrieved paragraphs, we rank them by relevance and only show the first paragraph in full by default and annotators can click to expand the other paragraphs.

In the task instructions, we explicitly discourage participants from searching the claims on the internet. Each participant verifies 20 claims one by one. We provide a tutorial at the beginning of the study. We include two attention check questions at different points in the study asking participants’ selection from the most recent claim and rejecting the responses from users who fail both attention checks.

A.2 Meta-Analysis

A.2.1 Confidence Calibration

We convert users’ confidence levels into discrete values $\mathcal{C} = \{0, 0.25, 0.5, 0.75, 1.0\}$. Our goal is for users to have high confidence in their correct judgments and low confidence in their wrong judgments. We plot their average confidence on correct and wrong judgments in Figure 7. User confidence is always low in the Baseline condition, which is reasonable since they do not have additional supporting evidence and are mostly making educated guesses. On correct judgments, users generally have high confidence (above 0.6). However, **users are over-confident on wrong judgments**, with average confidence above 0.6 as well. The Explanation and Contrastive Explanation conditions incur lower user confidence on both correct and wrong judgments as compared to the Retrieval condition, as well as the Retrieval + Explanation condition. Overall, these results highlight the difficulty of achieving appropriate calibration in users’ judgments.

A.2.2 Impact of Retrieval Recall

In previous sections, we performed a breakdown analysis based on the correctness of the explanations. In this section, we analyze another important dimension—the retrieval recall. We split examples into two categories: the first group where the

top-10 retrieved passages contain all the necessary evidence to verify the claim (*i.e.*, full recall = 1), and the second group where not all evidence is retrieved within the top-10 passages (*i.e.*, full recall = 0). We analyze how the retrieval recall affects both the explanation accuracy as well as the human decision accuracy.

The explanation accuracy is much lower when the retrieval recall is low. Over the entire test set of 200 examples, when the full recall = 1, the explanation accuracy is 80.4%; when the full recall = 0, the explanation accuracy is 67.6%. This indicates that retrieval quality has a high impact on explanation accuracy, which in turn affects human decision accuracy.

Human decision accuracy is much lower when the retrieval recall is low. Human decision accuracy broken down by retrieval recall is plotted in Figure 8. In all cases (apart from the Baseline condition where users do not see any evidence), the human decision accuracy is lower when the full retrieval recall is 0, sometimes it is lower than the case of full recall = 1 by large margins, *e.g.*, in the Retrieval condition and the Retrieval + Explanation condition.

A.2.3 Correlation Between Accuracy and Time

We plot the correlation between human decision accuracy and time in Figure 9. Overall, there is little correlation between accuracy and time ($r = 0.099$). We also report the correlation broken down for each condition:

- Baseline: $r = -0.180$
- Retrieval: $r = 0.089$
- Explanation: $r = -0.539$
- Contrastive Explanation: $r = -0.008$
- Retrieval + Explanation: $r = 0.148$

The correlation is weak for all conditions.

A.2.4 When Do Users Disagree with Explanations

We manually analyze the free-form responses provided by the users as their decision rationales. We first categorize these rationales in each condition and perform manual coding to quantitatively measure the frequency of each category of rationales. We analyze in what cases would users choose to disagree with the answers given by ChatGPT, and we break it down into the following two questions:

- (1) **How do users make the correct decision when the explanations are wrong?**

1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042

1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055

1056
1057
1058
1059
1060
1061
1062
1063
1064
1065

1066
1067

Task
1/20

Claim (Topic: Queen Victoria): Queen Victoria was heralded throughout the island North of England for starving the Irish.

Supporting Explanation

Refuting Explanation

ChatGPT Explanation

The claim that Queen Victoria was heralded throughout the island North of England for starving the Irish is not supported by the evidence from Wikipedia. In fact, during the Great Famine in Ireland, Victoria was labeled "The Famine Queen" and faced criticism for her response to the crisis. While she did personally donate £2,000 to famine relief efforts and supported the Maynooth Grant to a Roman Catholic seminary in Ireland, her actions were not widely seen as sufficient or effective in addressing the scale of the famine. The claim that she was heralded for starving the Irish is therefore not true.

ChatGPT Explanation

The claim that Queen Victoria was heralded throughout the island North of England for starving the Irish is false. The evidence from Wikipedia shows that during the Great Famine in Ireland, Victoria was actually labeled "The Famine Queen" and faced criticism for her response to the crisis. While she did personally donate a significant amount of money to famine relief efforts and supported the Maynooth Grant to a Roman Catholic seminary in Ireland, she was still criticized for not doing enough to alleviate the suffering of the Irish people. Therefore, the claim that she was heralded for starving the Irish is not supported by the evidence.

Do you think the claim is true?

No Yes

How confident are you about your judgment?

Very Uncertain Uncertain Neutral Certain Very Certain

How did you make that judgment?

Figure 6: Interface for the contrastive explanation condition. We present ChatGPT’s explanations for both sides together to encourage more careful thinking. We also highlight all the keywords to ease reading.

1068 • In the Explanation condition, most users rely
 1069 on **self-contradiction** in the ChatGPT explanations
 1070 (40.7%). For example, given the true claim
 1071 “Charles Evans Hughes shuffled off this mortal
 1072 coil in Massachusetts, and then was taken to
 1073 New York to be submerged in soil.”, ChatGPT
 1074 generates the explanation “The claim is false. Ac-
 1075 cording to the information provided, Hughes died
 1076 in Osterville, Massachusetts, and was interred at
 1077 Woodlawn Cemetery in the Bronx, New York City.”
 1078 where the explanation actually supports the claim
 1079 despite it saying the claim is false. Users did man-
 1080 age to catch this: “The explanation sounds like
 1081 it’s actually agreeing with the claim.” and made
 1082 the correct judgment.

1083 • In the Retrieval + Explanation condi-
 1084 tion, users mostly rely on information from re-
 1085 trieved passages (63.5%) and occasionally based
 1086 on ChatGPT’s self-contradiction (15.9%), e.g.,
 1087 users responded “I made the judgment by sum-
 1088 marizing the highlighted areas in the passages.”
 1089 The remaining less common cases are mostly based
 1090 on personal knowledge or guesses.

1091 **(2) When do users make wrong judgments even**
 1092 **when the explanations give correct answers?**

1093 • In the Explanation condition, most users rely
 1094 on personal knowledge or guess (46.4%) and
 1095 sometimes because the explanations have wrong
 1096 or poor-quality reasoning (25.0%). For example,
 1097 one user responded “I once took a tour of Alca-
 1098 traz and I believe I remember this as being true.”
 1099 to the claim “Within Alcatraz was a music room
 1100 where inmates could be rewarded for positive
 1101 behaviors with playing time.” which is in fact
 1102 false.

1103 • In the Retrieval + Explanation condition,
 1104 users mostly misinterpreted the evidence (38.1%)
 1105 or there was just insufficient evidence (28.6%)
 1106 and they had to make educated guesses. For ex-
 1107 ample, to the false claim “The Bee Gees went on
 1108 tour eleven times.”, one user responded: “The
 1109 articles mention mostly their songs and a couple
 1110 tours. I didn’t find anything about 11 tours, just a
 1111 couple of them. It was mainly songs I saw.” and
 1112 they judged the claim to be true even though the
 1113 ChatGPT explanation is correct: “Based on the
 1114 evidence from Wikipedia, the claim that the Bee
 1115 Gees went on tour eleven times is not supported.
 1116 The evidence mentions several tours that the Bee
 1117 Gees went on, including the 2 Years On Tour,

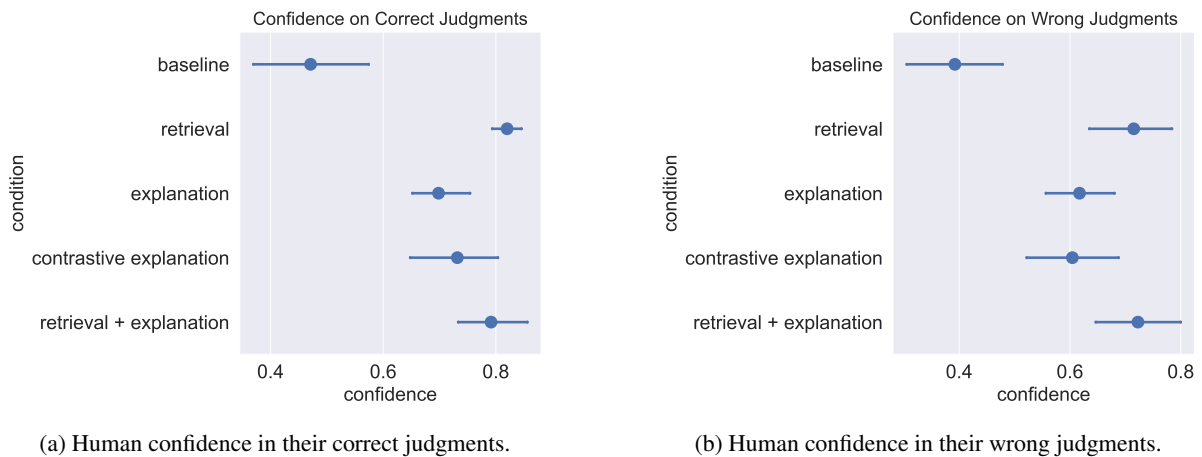


Figure 7: Human confidence broken down by their correct and wrong judgments. Users are over-confident are wrong judgments.

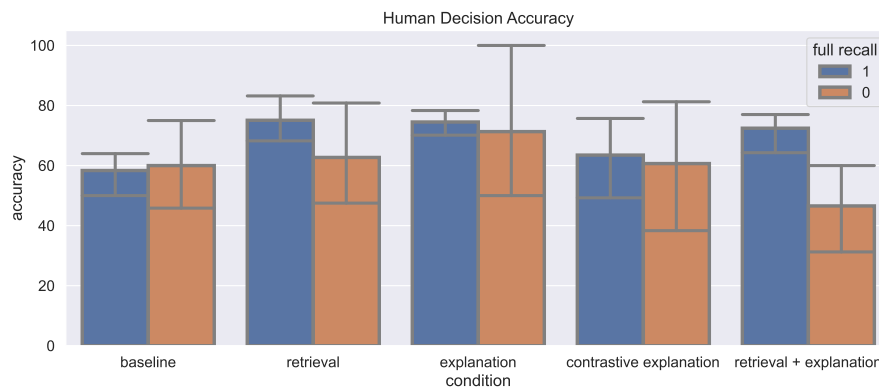


Figure 8: Human accuracy broken down by retrieval recall. Human accuracy is lower when the retrieval recall is low.

Trafalgar Tour, Mr. Natural Tour, Main Course Tour, Children of the World Tour, Spirits Having Flown Tour, One for All World Tour, High Civilization World Tour, and One Night Only World Tour. However, the evidence does not provide information about any additional tours beyond these mentioned. Therefore, it cannot be concluded that the Bee Gees went on tour eleven times based on the evidence provided.”

A.3 Additional Related Work

On Fact-Checking: Fact-checking is a well-established task in NLP where the typical task format is to input evidence text (*e.g.*, retrieved from Wikipedia) and the claim to the model and output a label of support or refute (or sometimes a third class of not enough information) (Vlachos and Riedel, 2014; Thorne and Vlachos, 2018). Automated fact-checking systems often make use of multiple evidence pieces for making predictions, and

optionally generating generations along with predictions (Popat et al., 2018a,b; Chen et al., 2023). Various HCI studies have also analyzed how fact-checking systems are used by domain experts such as journalists (Flores-Saviaga et al., 2022; Beers, 2019).

On Explanations: Explanations have been long sought as a useful tool to help users, not only in understanding AI predictions (Lee and See, 2004) but also aiding them in calibrating their reliance on these predictions (Bussone et al., 2015a). Some works find that explanations can support human-AI decision-making by exceeding both human-alone or AI-alone performance (Feng and Boyd-Graber, 2018; Bowman et al., 2022), whereas some other works find that explanations lead to worse human-AI performance (Alufaisan et al., 2021; Bansal et al., 2020; Wang and Yin, 2021). Vasconcelos et al. (2022) and Fok and Weld (2023) argue that to facilitate complementary human-AI decision-

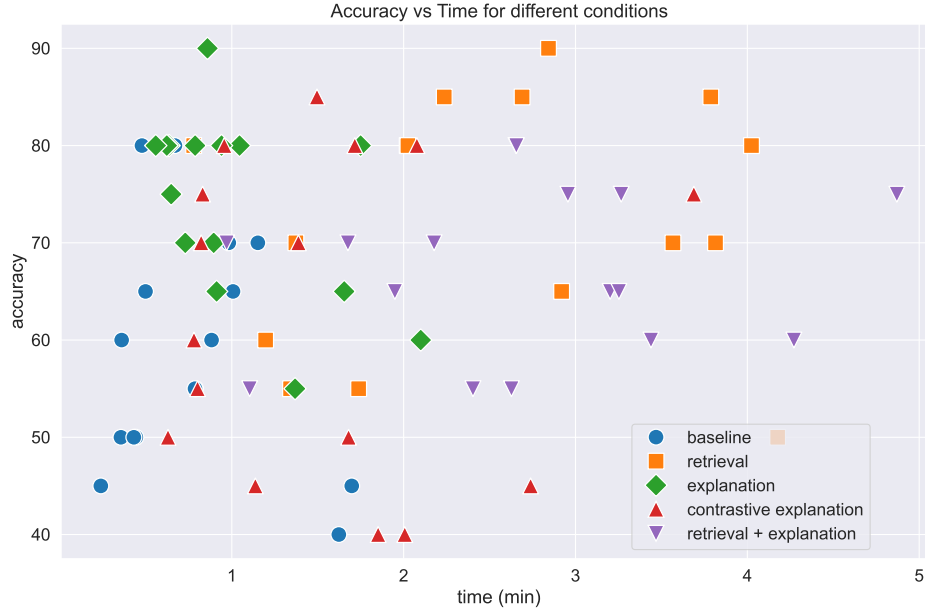


Figure 9: Correlation between each participant’s average decision accuracy (y-axis) and time (x-axis). We do not observe a strong correlation overall.

1157 making, explanations must aid users in verify-
 1158 ing the AI prediction to yield truly complemen-
 1159 tary human-AI performance. Explanations target-
 1160 ing verifiability have indeed shown promising av-
 1161 enues in human-AI collaborations (Feng and Boyd-
 1162 Graber, 2018; Vasconcelos et al., 2022; Goyal et al.,
 1163 2023).

1164 **On Explanations for Mitigating Over-Reliance:**

1165 In line with explanations, model indicators such as
 1166 confidence (Zhang et al., 2020) and accuracy (Yin
 1167 et al., 2019) have been found to yield mixed ben-
 1168 efits. On the one hand, uncertainty indicators can
 1169 promote slow thinking (Prabhudesai et al., 2023),
 1170 helping users calibrate trust in AI prediction. On
 1171 the other hand, humans find it difficult to interpret
 1172 numbers, leading to limited utility of such indica-
 1173 tors Zhang et al. (2020). Further, these indicators
 1174 can be unreliable as models’ accuracy in-the-wild
 1175 may differ from the reported accuracy on the eval-
 1176 uation set (Chiang and Yin, 2021) and models’ con-
 1177 fidence tend to be uncalibrated (Guo et al., 2017).
 1178 To resolve these limitations, Bussone et al. (2015b)
 1179 find that detailed explanations exacerbates the over-
 1180 reliance on the model predictions, whereas less
 1181 detailed explanations lead to distrust in the model,
 1182 but increases users’ self-reliance.