

Facilitating Human-LLM Collaboration through Factuality Scores and Source Attributions

Hyo Jin Do*
hjdo@ibm.com
IBM Research
USA

Rachel Ostrand*
rachel.ostrand@ibm.com
IBM Research
USA

Justin D. Weisz
jweisz@us.ibm.com
IBM Research
USA

Casey Dugan
cadugan@us.ibm.com
IBM Research
USA

Prasanna Sattigeri
psattig@us.ibm.com
IBM Research
USA

Dennis Wei
dwei@us.ibm.com
IBM Research
USA

Keerthiram Murugesan
keerthiram.murugesan@ibm.com
IBM Research
USA

Werner Geyer
werner.geyer@us.ibm.com
IBM Research
USA

ABSTRACT

While humans increasingly rely on large language models (LLMs), they are susceptible to generating inaccurate or false information, also known as “hallucinations”. Technical advancements have been made in algorithms that detect hallucinated content by assessing the factuality of the model’s responses and attributing sections of those responses to specific source documents. However, there is limited research on how to effectively *communicate* this information to users in ways that will help them appropriately calibrate their trust toward LLMs. To address this issue, we conducted a scenario-based study (N=104) to systematically compare the impact of various design strategies for communicating factuality and source attribution on participants’ ratings of trust, preferences, and ease in validating response accuracy. Our findings reveal that participants preferred a design in which phrases within a response were color-coded based on the computed factuality scores the most. Participants found it easy to validate the accuracy of an LLM’s response and increased their trust in this style compared to a baseline in which no style was applied. Additionally, participants increased their trust ratings when relevant sections of the source material were highlighted or responses were annotated with reference numbers corresponding to those sources, compared to when they received no annotation in the source material. Our study offers practical design guidelines to facilitate human-LLM collaboration and it promotes a new human role to carefully evaluate and take responsibility for their use of LLM outputs.

KEYWORDS

Large Language Models, Hallucinations, Factuality, Source Attribution, Human-AI Collaboration

1 INTRODUCTION

The rapid advancement of natural language generation technologies has led to the widespread use of large language models (LLMs) such as GPT, Bard, and LLaMA, in various tasks and contexts. However, LLMs are prone to presenting factually incorrect information as

if it were true, a phenomenon known as “*hallucination*” [15]. The presence of these hallucinations in LLM outputs, coupled with users’ inability to easily detect them and the tendency to over-trust LLMs, has resulted in several high-profile incidents. These include lawyers being reprimanded by judges for presenting nonexistent case law that had been hallucinated [40], new products being rapidly shelved due to hallucinated scientific references [37], news outlets issuing corrections to articles written with AI assistance [38], and company share prices falling after a hallucination leads to a blunder during a new product demo [12]. Upon realizing the hallucinations, users may lose trust in LLMs, which further impedes technology adoption [45]. Researchers are actively investigating methods to mitigate hallucinations, such as refining and improving datasets and models [15] using techniques such as reinforcement learning with human feedback [15, 34] and retrieval-augmented generation [4, 23]. However, technical advancements alone cannot completely resolve the issue; ultimately, it falls upon end-users to build an appropriate level of trust, be trained in how to carefully evaluate LLM outputs, and be accountable for their use.

Human-centered evaluation approaches offer a promising solution to address the hallucination issue. Several new methods have been developed that aid users in assessing the factual accuracy of a model’s response. Notably, these include *factuality scoring*, which evaluates the extent to which a model’s response is truthful to a source document [6, 19, 20, 30, 31, 48] and *source attribution*, which links the generated response to its source material [1, 7, 24, 33, 43]. However, it is currently unclear how to effectively *communicate* this factuality information to users. Should it be conveyed numerically or visually? At what level of linguistic granularity should such information be presented (e.g. word, phrase)? Recently, Leiser et al. [22] conducted participatory workshops where people brainstormed design strategies to identify hallucinations in the LLM outputs. However, no studies have systematically compared the effectiveness of these strategies in helping users assess the accuracy of the model’s response and calibrate their trust. Our study aims to identify the most effective and preferred strategy for communicating two pieces of information about an LLM’s response: (1) the *factuality score*: the automated assessment of how factual the

*Both authors contributed equally to this research.

response is, and (2) *source attribution*: identification of the sections within a source document from which the response was generated.

We explore various design strategies for communicating factuality information in a question-and-answer scenario. We first developed a factuality score color scale ranging from 0 (red) to 1 (green). We presented participants with three styles for visualizing factuality scores within an LLM’s response: (1) *highlight-all*, which annotates all of the linguistic content in the LLM response with varying background colors according to the scale, (2) *highlight-threshold*, which annotates only those parts of the LLM response where the factuality score is below a given threshold, and (3) *score*, which shows the numeric factuality score associated with each part of the response. Factuality scores were evaluated at two levels of linguistic granularity – *phrase* and *word* – and the three factuality styles were presented at each level of granularity. We additionally investigated two styles for presenting source attribution: (1) *highlight gradients*, in which linguistic components of the source that were used to generate the model’s response are highlighted, and (2) *reference numbers*, which displays in-line citations within the model’s response to specific, numbered parts of the source.

We conducted a scenario-based survey study (N=104) to compare the effects of these design strategies on participants’ ratings of trust, preference, and ease of evaluating response accuracy. Based on our findings, this paper makes three contributions to the literature on human-AI collaboration: 1) We explore the design space for presenting factuality and source attribution information to users and identify a set of promising approaches for deep evaluation based on user feedback; 2) We show that our design strategies have significant effects on ratings of trust and ease of accuracy validation; 3) We offer practical guidance on how to communicate factuality scores and source attribution within the user interface of LLM-based applications and thereby facilitate human-LLM interactions.

2 RELATED WORK

2.1 Calibrating End-User Trust for Human-AI Collaboration

Successful human-AI collaboration requires a user to modulate their level of trust in concert with the true reliability of the AI system, a process known as trust calibration [21, 44]. Miscalibrated trust can lead to overreliance – by accepting an incorrect AI recommendation – or underreliance – by not accepting a correct AI recommendation [44]. Researchers identified three main factors that influence trust: 1) AI-related factors of performance (e.g., reliability, failure rate) and attributes (e.g., anthropomorphism), 2) user-related factors such as ability (e.g., expertise, prior experiences) and characteristics (e.g., demographic), and 3) environment-related factors such as team collaboration (e.g., culture, communication, reputation) and tasks (e.g., task complexity and type) [10, 18]. Kim et al. argued that participants with limited domain expertise had difficulty in assessing the accuracy of AI-generated outputs [18]. Thus, we might expect that a user’s assessment of the accuracy of the model and their trust in the model would be related, due to the interrelationship between accuracy assessment and expertise, and separately, trust and expertise. In particular, these prior works suggest that a user’s trust in AI may be *contingent* on their accuracy assessment of the AI response. To get a full picture of trust, it is

also crucial to examine both general trustworthiness perceptions such as user expectations in the use of LLMs, and instance-specific trust-related behaviors such as users’ accuracy assessments of the LLM-generated output in the case of specific model responses [18]. In our work, we measured users’ ratings of the ease of validating the model response’s accuracy, their trust in the model, and their overall preferences among the different designs of factuality. We provided empirical evidence of trust calibration depending on their initial accuracy assessment of the model’s output.

2.2 Hallucination and Factuality Detection in Large Language Models

The widespread usage of LLMs in society has drawn attention to their risks and limitations. Notably, LLMs have the potential to generate text that seems plausible at first glance, but in reality, it is factually incorrect, a phenomenon known as *hallucination*. In the context of natural language generation technologies, hallucination refers to the generation of content that is factually inconsistent or unfaithful to the provided source. The counterpart of hallucination is *factuality*, defined as “truthfulness or the quality of being based on fact” [15]. The *source* document is essential in determining the factuality of an LLM output. If the model’s response aligns with the information from a reliable source, it is likely to be factually correct. On the other hand, *faithfulness* means the LLM-generated response stays consistent with the source. In this study, we assume a reliable source as our basis for “fact” so that faithfulness has the same meaning as factuality [30].

Hallucination in LLMs can stem from various factors, such as noisy, biased, and/or erroneous training data, as well as the model itself. As summarized in survey papers [13, 15], researchers have addressed data-related issues by establishing ground truth data through human annotators and enhancing model inputs with external knowledge [11, 13, 15]. However, it is impossible to completely resolve the hallucination issue inherent in AI technologies. Ultimately, it is the responsibility of end-users to carefully evaluate and be accountable for the use of LLM responses. As part of the effort to assist end-users in evaluating LLM responses, there has been ongoing research to develop methods for scoring the factuality of LLM outputs [6, 20, 29, 30, 48]. These methods can either use lexical matching-based metrics [2, 27, 35] or model-based metrics using neural networks [31, 39, 47]. This increasing body of research raises new questions for LLM developers and designers on how to effectively *communicate* factuality information to end-users. Specifically, there are no guidelines on which parts of the LLM response (e.g., correct, incorrect, or both) should be annotated, in what visual style (e.g., numeric, color-coded), and at what level of linguistic granularity accuracy should be assessed (e.g., word, phrase, entire response). Furthermore, we have little understanding of how communicating the factuality of LLM outputs mitigates the effects of hallucination and calibrates end-users’ trust. In addressing this gap, the present work identifies the most preferred and effective design to communicate the factuality and source attribution of LLM outputs, and presents practical guidelines based on our findings.

2.3 Source Attributions as Explanations

Source attribution is connected to research on explainable AI (XAI), providing explanations to support appropriate human understanding of AI-generated outputs [5, 16, 17, 25, 29, 32, 36, 42, 46]. Liao and Vaughan emphasized the importance of communicating and being transparent about information during interactions with LLMs, given that LLMs raise unique challenges in XAI including new and complex model capabilities, behaviors, and applications of LLMs, massive and opaque architectures, and organizational pressure to move fast and deploy at scale [26]. We investigated two design strategies – highlight gradients and reference numbers – to explain which parts of a source are most relevant to the LLM-generated response. These design strategies are widely used in LLMs [5, 8, 9, 16, 29, 36, 43, 46] and other similar contexts that incorporate sources (e.g., news media [41] and academic research citations [14]). Despite the common use of both design strategies in real-world applications and research, there is no research we are aware of that investigates the effects of these strategies on end-users’ trust and ease of validating the response accuracy within the context of LLMs, both of which we address in the present work.

3 METHODS

The goal of the current study was to understand how to present information about factuality and source attributions to an end-user in a way that is easy to understand and helps them calibrate their trust in LLM-generated text. To achieve this goal, we first reviewed designs of commercial generative AI systems and prior research to understand how other researchers and designers have presented factuality and source attribution information, and we ideated additional ways to present this information to a user. For the controlled study, we selected six different designs for representing a factuality score and two different designs for representing source attributions. This selection was based on a pilot study in which we interviewed ten participants about their preferences on numerous design options.

3.1 Participants

There were 104 participants in the controlled study, who were employees of a large, multinational technology company. Participant recruitment was advertised widely within the company on 25 internal Slack channels spanning multiple divisions and geographic regions in order to recruit a diverse sample across multiple professional, demographic, and experiential characteristics. Participants’ work locations consisted of 20 unique countries, and job roles spanned a wide array of disciplines, including design, customer service, engineering, sales, research, and HR, among others. Participants had a range of experiences with AI as a technology, with some having heard about it from the news, work, friends or family (N=14), others reporting that they “closely follow” AI news (N=26), many reporting some work or educational experience regarding AI (N=49), and a small number with “significant” work experience with AI (N=15). Participants also reported a wide range of experience with LLMs, with daily usage (N=10), a few times a week (N=21), once a week (N=10), once a month (N=15), a couple of times a year (N=15), a couple of times in life (N=13), never (N=19), and not sure (N=1). Participants self-rated their English proficiency

on a 7-point Likert scale, with 68% rating themselves at 7 (*native or native-like proficiency*), 19% as 6, 8% as 5, 4% as 4 (*medium*), and 1% as 3. All participants provided written informed consent and were treated in accordance with the guidelines for the ethical treatment of human participants.

3.2 Procedure

Participants were told to put themselves in the place of a user of an AI-powered language model. Their task was to evaluate different designs for presenting the factuality of the model’s response, and the source that the response was drawn from. On each trial, participants were shown a snapshot of a supposed interaction with an LLM, with three components: a *Question*, a *Response*, and a *Source*. The Question was “What movies did Beyonce star in and with whom?” and the AI-generated Response was “Beyonce starred in the musical comedy *The Fighting Temptations* in 2002 and in the documentary film *Austin Powers in Goldmember* in 2003, alongside Missy Elliott and Foxy Cleopatra, respectively.” The source was a paragraph from Wikipedia, and we asked participants to assume that the source was factually accurate. The response was written to have some factually inaccurate propositions (and importantly, contradictory of the source document) and other accurate ones; overall, it was approximately half accurate and half inaccurate. The question and the response were carefully selected to test our design strategies, and covered a topic that was not technical, to enable non-expert participants to make their own assessment as to its accuracy.

During the study, participants were shown a series of different design strategies to evaluate. Each design was demonstrated using the same Question, Response, and Source text, to hold constant the content and level of accuracy across different designs. This allowed us to reduce the number of variables tested and ensure a more targeted exploration of our design strategies. Participants were always shown the *Baseline* design first, which had no markup, and displayed only the text of the Question, Response, and Source. Next, participants were presented with six different design strategies for displaying the model response’s factuality, and finally, two different design strategies for showing source attribution in the response. The study was conducted as a within-subjects experiment, and thus all participants viewed and rated the same designs.

As a first step, participants were asked to rate their perceptions about the model and its response on a 7-point Likert scale for the *Baseline* design along three dimensions:

- (1) *Perceived accuracy*: How accurate do you think this AI-generated response is?
- (2) *Ease of validation*: With the information presented in this way, how easy is it for you to determine the accuracy of this AI-generated response?
- (3) *Trust*: With the information presented in this way, how much do you trust the AI system that generated the response?

3.2.1 Factuality Score. Following the baseline design, participants were introduced to the concept of a *factuality score* – a feature that compares linguistic components of the response against the source – and that a high factuality score indicated that the response aligned with the information in the source and thus is likely to be correct.

Three factuality design strategies were presented to each user, each at two levels of granularity. The designs were *highlights-all*, in which every part of the response text was highlighted with a color-coding to show its level of factuality on a red (0) to green (1) scale; *highlights-threshold*, in which only the sections of the response text with a factuality score below 0.5 were highlighted, to signal inaccuracies; and *score*, in which all parts of the response text were tagged with their factuality score, but instead of highlights, with color-coded underlines and the numerical factuality score value.

In addition, the designs were presented at two levels of *granularity* – either *word-level* or *phrase-level*. This refers to the amount of text over which the factuality was evaluated. At phrase-level granularity, if there was an inaccuracy in one word in a phrase, then the entire phrase would be tagged with a lower factuality score. In contrast, at word-level granularity, only that word or a group of words would be tagged with a lower factuality score, while the other words in the sentence (if correct), would individually be tagged with a higher factuality score. Table 1 shows the six designs that users evaluated.

Participants were shown each factuality design strategy one at a time, and asked to rate their perceptions on two dimensions: ease of validation and trust (questions (2) and (3) as listed in Section 3.2), using a 7-point Likert scale. Note that we did not ask users about perceived accuracy (question (1)) for any designs except for the Baseline, because the wording of the text was identical in the Baseline and every presented design.

Participants performed this rating task for the three designs at one granularity (word-level vs. phrase-level), and then were asked to rank-order them, along with the baseline, in their order of preference. They then performed the same rating and preference-ranking task for the three designs at the other granularity. The three designs within each granularity were presented in a randomized order to each participant, and the order of the two granularities was also randomized across participants, to reduce possible order effects in the aggregated data. Following the ratings, we asked a supplemental question regarding preference between the two types of granularity designs.

3.2.2 Source Attribution. Next, participants were introduced to the source attribution feature, in which the source was annotated to show which parts were used to generate the model’s response. There were two designs that were presented to users: *reference numbers*, in which each sentence of the source document was numbered, and propositions in the response were tagged with the number corresponding to the source sentence from which it was derived; and *highlight gradients*, in which sections of the source that provided the information for parts of the response were highlighted. The order of presentation of the two source attribution designs was randomized between participants. Fig. 1 presents the two designs that users evaluated.

After each source attribution design, participants were asked the same two questions as for the factuality score designs. After rating the two source attribution designs, participants were then asked to rank-order their preference among the two source attribution designs and the baseline. At the end of the survey, participants responded to some demographic and professional questions as reported in Section 3.3.

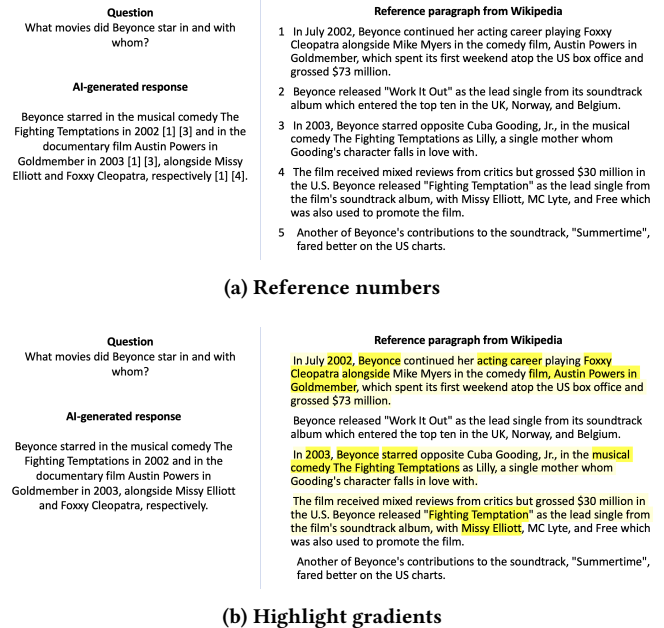


Figure 1: The set of designs presented to each participant for displaying the source attribution by the model. Each participant saw and rated both designs, in a randomized order.

3.3 Participants

There were 104 participants in the study, who were employees of a large, multinational technology company. As such, participant recruitment was advertised widely within the company on 25 internal Slack channels spanning multiple divisions and geographic regions. All participants provided written informed consent and were treated in accordance with the guidelines for the ethical treatment of human participants. Participants’ work locations consisted of 20 unique countries. Job roles spanned a wide array of disciplines, including design, customer service, engineering, sales, research, HR, among others. Participants had a range of experiences with AI as a technology, with some having heard about it from the news, work, friends or family (N=14), others reporting that they "closely follow" AI news (N=26), the largest subset reporting some work or educational experience regarding AI (N=49), and others with "significant" work experience with AI (N=15). Participants reported a wide range of experience with LLMs, and varying degrees of English exposure and proficiency.

4 RESULTS

The analyses were conducted using generalized linear mixed-effects models, with one model for each dependent variable: participants’ ratings of (1) trust and (2) ease of validating the response accuracy. In both models, the categorical independent variable Design Strategy was treatment-coded with the Baseline design set as the reference level, such that each design’s rating was statistically compared against the baseline. Participant ID was a random variable. Following the omnibus models, pairwise contrasts were conducted

Table 1: The set of designs presented to each participant for displaying factuality scores on the model’s response. Each participant saw and rated all six designs, in a randomized order but grouped by granularity.

| Granularity | Highlight-all | Highlight-threshold | Score |
|-------------|--|---|--|
| Word | Beyonce starred in the musical comedy The Fighting Temptations in 2002 and in the documentary film Austin Powers in Goldmember in 2003, alongside Missy Elliott and Foxy Cleopatra , respectively. | Beyonce starred in the musical comedy The Fighting Temptations in 2002 and in the documentary film Austin Powers in Goldmember in 2003, alongside Missy Elliott and Foxy Cleopatra , respectively. | Beyonce ⁵ starred ⁵ in the musical comedy ⁷ The Fighting Temptations , ⁶ in 2002 ³ and in the documentary film ⁴ Austin Powers in Goldmember ⁶ in 2003 ³ , alongside Missy Elliott ⁴ and Foxy Cleopatra ⁴ , respectively. |
| Phrase | Beyonce starred in the musical comedy The Fighting Temptations in 2002 and in the documentary film Austin Powers in Goldmember in 2003, alongside Missy Elliott and Foxy Cleopatra , respectively. | Beyonce starred in the musical comedy The Fighting Temptations in 2002 and in the documentary film Austin Powers in Goldmember in 2003, alongside Missy Elliott and Foxy Cleopatra , respectively. | Beyonce ⁵ starred in the musical comedy The Fighting Temptations in 2002 ⁵ and in the documentary film Austin Powers in Goldmember in 2003, ³ alongside Missy Elliott and Foxy Cleopatra , respectively. ⁶ |

to explore comparisons between each pair of Design Strategy levels, with p -values adjusted for multiple comparisons using the Tukey correction.

4.1 Factuality Score

4.1.1 Trust. First, we compared users’ ratings of their trust of the model that produced the response, for each of the designs compared against the baseline. All six of the designs were rated as significantly more trustworthy compared to the baseline, suggesting that all of the designs that presented response factuality increased users’ trust in the model. The mean and the standard error (SE) of the ratings for each design, along with results from the statistical model comparing each design to the baseline, are displayed in Table 2a. Post-hoc pairwise comparisons between all pairs of designs revealed no additional significant differences after correction for multiple comparisons.

As an exploratory analysis, we investigated whether participants’ baseline rating of the model response’s accuracy affected how much they trusted the model when it subsequently provided factuality scores. To do so, we ran another linear mixed-effects model with the same structure as above, with the addition of *perceived accuracy of the response at Baseline* as an independent variable. Perceived baseline accuracy significantly affected participants’ subsequent ratings of trust ($t = 4.00, p < .001$). For visualization purposes, we categorized participants into two groups: The *low baseline accuracy* group, which was defined as those participants who rated the baseline response accuracy at or below 4 ($N=87$), and the *high baseline accuracy* group, who rated the baseline response accuracy as 5 or higher ($N=17$). As illustrated in Fig. 2(a), participants in the low baseline accuracy group initially rated the baseline with low trust, but subsequently increased their trust ratings after reviewing the factuality scores. In contrast, participants in the high baseline accuracy group (Fig. 2(b)) initially had higher trust in the model’s response but subsequently decreased their trust ratings after examining the factuality information presented in each design style, particularly at the phrase-level.

4.1.2 Ease of validation. We next compared users’ ratings about the ease of assessing the model’s accuracy for each of the designs, compared against the baseline. Of the six design strategies, three were rated as significantly easier to assess the response accuracy

compared to the baseline: *highlights-all* at phrase-level granularity, and *highlights-all* and *highlights-threshold* at word-level granularity. The other three designs were not rated significantly different from the baseline. Post-hoc pairwise comparisons between each pair of design strategies revealed no significant differences after correction for multiple comparisons. The mean and standard errors of the ratings for each design, along with results from the statistical model comparing each design to the baseline, are displayed in Table 2b.

4.1.3 Preference. Participants next rank-ordered each of the three designs plus the baseline, within each granularity level. Thus, each ranking compared four designs. The results are presented with the ranking scores reversed (i.e., 4 - ranking score) such that a higher score corresponds to a more preferred design for comparability with the trust and validation ratings. At phrase-level granularity, the *highlights-all* design was the most preferred ($M = 1.94$), *score* was second ($M = 1.69$), *highlights-threshold* was third ($M = 1.67$) and baseline was the least preferred ($M = 0.69$). At word-level granularity, rankings were similar: the *highlights-all* design was again the most preferred ($M = 1.87$), *highlights-threshold* was second ($M = 1.79$), *score* was third ($M = 1.56$), and the baseline was the least preferred ($M = 0.78$).

Participants were also asked their preference between the two types of granularities: 52.9% of participants preferred *phrase-level* granularity, while 26.9% of the participants preferred *word-level* granularity, with 10.6% of the participants responding with “don’t know” and 9.6% of the participants selecting “other”.

4.2 Source Attribution

4.2.1 Trust. Both *reference numbers* and *highlight gradients* caused the model to be rated as significantly more trustworthy compared to the baseline, as shown in Table 3a. Post-hoc pairwise comparisons between the two design strategies revealed no additional significant differences.

4.2.2 Ease of validation. Both *reference numbers* and *highlight gradients* were rated significantly *lower* (i.e., worse) compared to the baseline, as shown in Table 3b. Post-hoc pairwise comparisons between the two design strategies showed no significant differences between them.

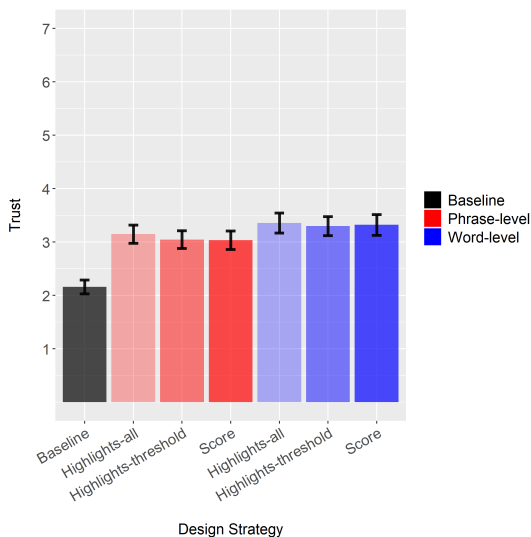
Table 2: Users’ ratings of (a) their trust in the model and (b) the ease of assessing the accuracy of the model’s response, for the baseline and each of the factuality score designs. t and p values were calculated within the omnibus statistical model, and represent the comparison of each Design Strategy against the Baseline as the reference level. Bolded Design Strategy names indicate a significant difference from the Baseline.

| Design Strategy | Mean | SE | t | p |
|-----------------------------|------|------|------|--------|
| Baseline | 2.63 | 0.17 | - | - |
| Phrase-level granularity | | | | |
| Highlights-all | 3.31 | 0.17 | 4.59 | < .001 |
| Highlights-threshold | 3.22 | 0.16 | 4.00 | < .001 |
| Score | 3.21 | 0.17 | 3.93 | < .001 |
| Word-level granularity | | | | |
| Highlights-all | 3.62 | 0.18 | 6.57 | < .001 |
| Highlights-threshold | 3.49 | 0.18 | 5.80 | < .001 |
| Score | 3.54 | 0.18 | 5.85 | < .001 |

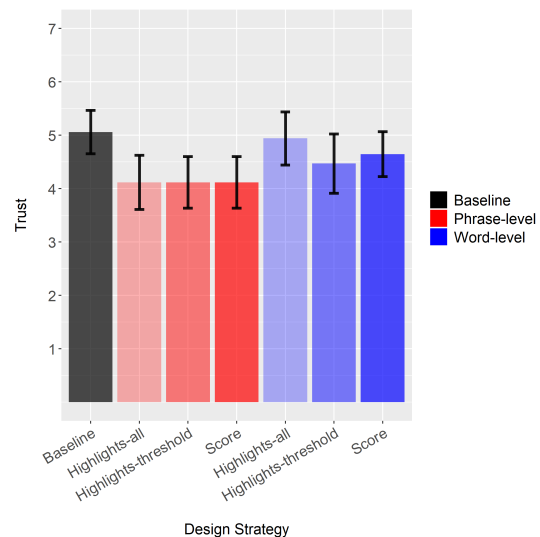
(a) User trust of the model

| Design Strategy | Mean | SE | t | p |
|-----------------------------|------|------|------|------|
| Baseline | 4.29 | 0.20 | - | - |
| Phrase-level granularity | | | | |
| Highlights-all | 4.74 | 0.17 | 2.24 | .03 |
| Highlights-threshold | 4.53 | 0.17 | 1.19 | .23 |
| Score | 4.38 | 0.17 | 0.48 | .63 |
| Word-level granularity | | | | |
| Highlights-all | 4.72 | 0.19 | 2.10 | .04 |
| Highlights-threshold | 4.88 | 0.16 | 2.96 | .003 |
| Score | 4.32 | 0.20 | 0.13 | .89 |

(b) Ease of assessing the accuracy of the response



(a) Participants with low baseline accuracy rating



(b) Participants with high baseline accuracy rating

Figure 2: Participants showed different levels of trust in the model as a function of their perceived accuracy at baseline. Participants who identified errors in the baseline response (a) reported higher levels of trust after reviewing factuality designs. Participants who initially missed errors in the baseline response (b) reported lower levels of trust after reviewing the factuality designs.

4.2.3 Preference. Participants were asked to rank-order their preference for the two designs and the baseline as a way to present source attribution information. As with the factuality scores, the data here are presented with the ranking scores reversed (i.e., 3 - ranking score) such that a higher score corresponds to a more preferred design. Participants preferred both of the two designs over the baseline, with *reference numbers* the most preferred by a small margin ($M = 1.21$), followed by *highlight gradient* ($M = 1.17$), and the least preferred was the baseline ($M = 0.62$).

5 DISCUSSION

In this study, we offer insight into ways of presenting information to an end-user during human-AI collaboration, allowing the user to assess the factuality of LLM responses that have the potential to be hallucinated. This user study presented multiple design strategies for displaying factuality and source attribution information from an LLM. Here, we abstract across the different results to present design recommendations and discuss limitations.

The most common granularity preference among participants for presenting factuality information within a model’s response was *phrase-level granularity*. Among the phrase-level granularity styles,

Table 3: Users’ ratings of (a) their trust in the model and (b) the ease of assessing the accuracy of the model’s response, for the baseline and each of the designs for source attribution. Bolded Design Strategy names indicate significant differences from baseline.

| Design Strategy | Mean | SE | <i>t</i> | <i>p</i> |
|----------------------------|------|------|----------|----------|
| Baseline | 2.63 | 0.17 | - | - |
| Reference numbers | 3.40 | 0.18 | 4.71 | < .001 |
| Highlight gradients | 3.23 | 0.18 | 3.74 | < .001 |

(a) User trust of the model

| Design Strategy | Mean | SE | <i>t</i> | <i>p</i> |
|----------------------------|------|------|----------|----------|
| Baseline | 4.29 | 0.20 | - | - |
| Reference numbers | 3.81 | 0.19 | -2.18 | .03 |
| Highlight gradients | 3.81 | 0.19 | -2.19 | .03 |

(b) Ease of assessing the accuracy of the response

the *highlights-all* style, which color-codes each of the phrases in the response with the color associated with the computed factuality score, was the most preferred. The statistical results showed that the *highlights-all* style also led to significantly higher trust and was significantly easier to validate response accuracy than the baseline. Thus, we make the design recommendation to present factuality information using a design similar to our **highlights-all at phrase-level granularity** design.

For source attribution, we recommend adopting either the **reference numbers** or **highlight gradients** design to enhance trust and accommodate the preferences of end-users, as both design strategies were effective in increasing trust compared to the baseline. The finding supports existing XAI research that explanations can increase user trust. The procedural justice theory [28] suggests that people’s trust is strongly impacted by procedural explanations and not only the outcome. Our design strategies provided a procedural explanation of how the model’s response was generated, even where the response was partially inaccurate. In contrast, participants commented that they had difficulty assessing the accuracy of the model’s response using the design strategies because they felt overwhelmed and distracted by them. Therefore, we recommend incorporating a feature that enables users to turn off or filter styles to reduce distractions, or even remove the styles if prioritizing ease of validation over building high trust.

How participants perceived the accuracy of the model’s response in the baseline design had a substantial impact on their trust after they viewed the factuality scores. Participants who initially rated the model accuracy as high, despite the presence of multiple errors in its response, decreased their trust upon seeing the errors called out through the factuality scores. In contrast, participants who initially rated the model accuracy as low increased their trust when they observed that the factuality scores accurately flagged those errors. Expectancy violations theory supports the finding that positive violations (i.e., when perceived performance exceeds rather than meets the expected level of performance) have a stronger positive effect on satisfaction, while negative violations produce a negative effect [3]. Therefore, to calibrate the level of end-users’ trust, these results support incorporating factuality information in the LLM response.

While our study assumed that the algorithm generating factuality scores for the model’s response is reliable, the algorithm itself may be imperfect or erroneous. A similar situation exists for source attribution. We assumed the existence of a reliable source attribution algorithm, but in practice, source attribution and AI

model explanation are ongoing research topics, especially for generative language models that output text rather than numerical predictions. These issues are heightened when there are multiple source documents, some of which may be irrelevant or unreliable. Therefore, end-users should be aware of the limitations of these AI technologies, particularly in the context of LLMs, which may appear factual at face value. It is crucial for users to always verify the model’s responses across multiple sources to ensure the reliability of information and prevent themselves from placing over-trust in LLMs.

The current study had a few limitations that should be areas of future research. It focused on a single question-and-answer task, as there were multiple interventions (i.e., design strategies) to compare within this task. Additionally, the model’s response and the assigned factuality scores were handcrafted, rather than generating an actual LLM response and scores from existing factuality or source attribution algorithms. This allowed us to easily create designs that effectively tested our research questions, but real-world LLM responses and algorithms may be different. While we made efforts to recruit participants with diverse backgrounds, skills, and locations, our recruitment was restricted to individuals within our company. Finally, it is important to note that our research did not aim to exhaustively explore all potential design strategies. Instead, our study should be viewed as a starting point, encouraging researchers to delve deeper into diverse design strategies and expand the discussion.

6 CONCLUSION

Large language models have known problems with so-called hallucinations. To address these challenges, a growing area of research is the development of algorithms to assess the model response’s factuality and attribute it to sources. However, how to effectively communicate factuality and source attribution to end-users is an open question. In this study, we designed and compared six design strategies for communicating factuality scores and two design strategies for conveying source attribution. Conducting a scenario-based study through an online survey, we discovered that highlighting every phrase in the model’s response based on the factuality score is the most preferred strategy and leads to higher trust than the baseline without any markups. Our findings also revealed that participants calibrated their trust in the model based on their initial accuracy assessment of the response. Regarding source attribution, reference numbers and highlight gradients enhanced trust in the model, but did not alleviate the challenge of assessing the response

accuracy. We provide practical guidance on communicating source attribution and factuality scores to facilitate successful human-LLM collaboration.

REFERENCES

- [1] [n. d.]. Perplexity. <https://www.perplexity.ai/> Accessed: 2023-12-12.
- [2] Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*. 65–72.
- [3] Judee K Burgoon, Joseph A Bonito, Paul Benjamin Lowry, Sean L Humpherys, Gregory D Moody, James E Gaskin, and Justin Scott Giboney. 2016. Application of expectancy violations theory to communication with and judgments about embodied agents during a decision-making task. *International Journal of Human-Computer Studies* 91 (2016), 24–36.
- [4] Deng Cai, Yan Wang, Lemaoy Liu, and Shuming Shi. 2022. Recent advances in retrieval-augmented text generation. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 3417–3419.
- [5] Hanjie Chen, Guangtao Zheng, and Yangfeng Ji. 2020. Generating Hierarchical Explanations on Text Classification via Feature Interaction Detection. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 5578–5593. <https://doi.org/10.18653/v1/2020.acl-main.494>
- [6] I Chern, Steffi Chern, Shiqi Chen, Weizhe Yuan, Kehua Feng, Chunting Zhou, Junxian He, Graham Neubig, Pengfei Liu, et al. 2023. FacTool: Factuality Detection in Generative AI—A Tool Augmented Framework for Multi-Task and Multi-Domain Scenarios. *arXiv preprint arXiv:2307.13528* (2023).
- [7] Kundan Krishna et al. [n. d.]. Evidence Inspector. <https://evinspector.site/> Accessed: 2023-12-12.
- [8] Luyi Gao, Zhuyuan Dai, Panupong Pasupat, Anthony Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vincent Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan, et al. 2023. Rarr: Researching and revising what language models say, using language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 16477–16508.
- [9] Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. 2023. Enabling Large Language Models to Generate Text with Citations. *arXiv preprint arXiv:2305.14627* (2023).
- [10] Peter A Hancock, Deborah R Billings, Kristin E Schaefer, Jessie YC Chen, Ewart J De Visser, and Raja Parasuraman. 2011. A meta-analysis of factors affecting trust in human-robot interaction. *Human factors* 53, 5 (2011), 517–527.
- [11] Or Honovich, Roei Aharoni, Jonathan Herzig, Hagai Taitelbaum, Doron Kukliansky, Vered Cohen, Thomas Scialom, Idan Szpektor, Avinatan Hassidim, and Yossi Matias. 2022. TRUE: Re-evaluating Factual Consistency Evaluation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 3905–3920.
- [12] Elizabeth Howell. 2023. James Webb Telescope question costs Google \$100 billion – here’s why. <https://www.space.com/james-webb-space-telescope-google-100-billion>
- [13] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2023. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232* (2023).
- [14] Ken Hyland. 1999. Academic attribution: Citation and the construction of disciplinary knowledge. *Applied linguistics* 20, 3 (1999), 341–367.
- [15] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *Comput. Surveys* 55, 12 (2023), 1–38.
- [16] Yiming Ju, Yuanzhe Zhang, Kang Liu, and Jun Zhao. 2023. A Hierarchical Explanation Generation Method Based on Feature Interaction Detection. In *Findings of the Association for Computational Linguistics: ACL 2023*. 12600–12611. <https://doi.org/10.18653/v1/2023.findings-acl.798>
- [17] Siwon Kim, Jihun Yi, Eunji Kim, and Sungroh Yoon. 2020. Interpretation of NLP models through input marginalization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 3154–3167. <https://doi.org/10.18653/v1/2020.emnlp-main.255>
- [18] Sunnie SY Kim, Elizabeth Anne Watkins, Olga Russakovsky, Ruth Fong, and Andrés Monroy-Hernández. 2023. Humans, AI, and Context: Understanding End-Users’ Trust in a Real-World Computer Vision Application. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*. 77–88.
- [19] Wojciech Kryściński, Bryan McCann, Caiming Xiong, and Richard Socher. 2020. Evaluating the Factual Consistency of Abstractive Text Summarization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 9332–9346.
- [20] Philippe Laban, Tobias Schnabel, Paul N Bennett, and Marti A Hearst. 2022. SummaC: Re-visiting NLI-based models for inconsistency detection in summarization. *Transactions of the Association for Computational Linguistics* 10 (2022), 163–177.
- [21] John D Lee and Katrina A See. 2004. Trust in automation: Designing for appropriate reliance. *Human factors* 46, 1 (2004), 50–80.
- [22] Florian Leiser, Sven Eckhardt, Merlin Knaeble, Alexander Maedche, Gerhard Schwabe, and Ali Sunyaev. 2023. From ChatGPT to FactGPT: A Participatory Design Study to Mitigate the Effects of Large Language Model Hallucinations on Users. In *Proceedings of Mensch und Computer 2023*. 81–90.
- [23] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Kuttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems* 33 (2020), 9459–9474.
- [24] Daniel Li, Thomas Chen, Alec Zadikian, Albert Tung, and Lydia B Chilton. 2023. Improving Automatic Summarization for Browsing Longform Spoken Dialog. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–20.
- [25] Q Vera Liao, Daniel Gruen, and Sarah Miller. 2020. Questioning the AI: informing design practices for explainable AI user experiences. In *Proceedings of the 2020 CHI conference on human factors in computing systems*. 1–15.
- [26] Q Vera Liao and Jennifer Wortman Vaughan. 2023. AI Transparency in the Age of LLMs: A Human-Centered Research Roadmap. *arXiv preprint arXiv:2306.01941* (2023).
- [27] Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*. 74–81.
- [28] E Allan Lind and Tom R Tyler. 1988. *The social psychology of procedural justice*. Springer Science & Business Media.
- [29] Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. *Advances in neural information processing systems* 30 (2017).
- [30] Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization. *arXiv preprint arXiv:2005.00661* (2020).
- [31] Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. FActScore: Fine-grained Atomic Evaluation of Factual Precision in Long Form Text Generation. *arXiv preprint arXiv:2305.14251* (2023).
- [32] Edoardo Mosca, Ferenc Szigeti, Stella Tragianni, Daniel Gallagher, and Georg Groh. 2022. SHAP-Based Explanation Methods: A Review for NLP Interpretability. In *Proceedings of the 29th International Conference on Computational Linguistics*. 4593–4603. <https://aclanthology.org/2022.coling-1.406>
- [33] Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. 2021. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332* (2021).
- [34] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems* 35 (2022), 27730–27744.
- [35] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*. 311–318.
- [36] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. “Why should I trust you?” Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 1135–1144.
- [37] Jackson Ryan. 2022. Meta Trained an AI on 48M Science Papers. It Was Shut Down After 2 Days. *CNET* (20 Nov 2022). Retrieved 04-Jan-2024 from <https://www.cnet.com/science/meta-trained-ai-on-48-million-science-papers-it-was-shut-down-after-2-days/>
- [38] Mia Sato and Emma Roth. 2023. CNET found errors in more than half of its AI-written stories. *The Verge* (25 Jan 2023). Retrieved 04-Jan-2024 from <https://www.theverge.com/2023/1/25/23571082/cnet-ai-written-stories-errors-corrections-red-ventures>
- [39] Thibault Sellam, Dipanjan Das, and Ankur P Parikh. 2020. BLEURT: Learning robust metrics for text generation. *arXiv preprint arXiv:2004.04696* (2020).
- [40] Karen Sloan. 2023. A lawyer used ChatGPT to cite bogus cases. What are the ethics? <https://www.reuters.com/legal/transactional/lawyer-used-chatgpt-cite-bogus-cases-what-are-ethics-2023-05-30/>
- [41] S Shyam Sundar. 1998. Effect of source attribution on perception of online news stories. *Journalism & Mass Communication Quarterly* 75, 1 (1998), 55–68.
- [42] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic Attribution for Deep Networks. In *Proceedings of the 34th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 70)*. 3319–3328. <https://proceedings.mlr.press/v70/sundararajan17a.html>
- [43] Shahbaz Syed, Tariq Yousef, Khalid Al-Khatib, Stefan Jänicke, and Martin Potthast. 2021. Summary explorer: Visualizing the state of the art in text summarization. *arXiv preprint arXiv:2108.01879* (2021).
- [44] Magdalena Wischnewski, Nicole Krämer, and Emmanuel Müller. 2023. Measuring and Understanding Trust Calibrations for Automated Systems: A Survey of the State-Of-The-Art and Future Directions. In *Proceedings of the 2023 CHI Conference*

- on *Human Factors in Computing Systems*. 1–16.
- [45] Kewen Wu, Yuxiang Zhao, Qinghua Zhu, Xiaojie Tan, and Hua Zheng. 2011. A meta-analysis of the impact of trust on technology acceptance model: Investigation of moderating influence of subject and context type. *International Journal of Information Management* 31, 6 (2011), 572–581.
- [46] Kayo Yin and Graham Neubig. 2022. Interpreting Language Models with Contrastive Explanations. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. 184–198. <https://doi.org/10.18653/v1/2022>.
- emnlp-main.14
- [47] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675* (2019).
- [48] Jiawei Zhou, Yixuan Zhang, Qianni Luo, Andrea G Parker, and Munmun De Choudhury. 2023. Synthetic lies: Understanding ai-generated misinformation and evaluating algorithmic and human solutions. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–20.