# Are UFOs Driving Innovation? The Illusion of Causality in Large Language Models

**Anonymous Author(s)**
Affiliation
Address
`email`

## Abstract

Illusions of causality occur when people develop the belief that there is a causal
connection between two variables with no supporting evidence. This cognitive bias
has been proposed to underlie many societal problems including social prejudice,
stereotype formation, misinformation and superstitious thinking. In this research
we investigate whether large language models develop the illusion of causality
in real-world settings. We evaluated and compared news headlines generated by
GPT-4o-Mini, Claude-3.5-Sonnet, and Gemini-1.5-Pro to determine whether the
models incorrectly framed correlations as causal relationships. In order to also
measure sycophancy behavior, which occurs when a model aligns with a user's
beliefs in order to look favorable even if it is not objectively correct, we additionally
incorporated the bias into the prompts, observing if this manipulation increases
the likelihood of the models exhibiting the illusion of causality. We found that
Claude is the model that presents the lowest degree of causal illusion aligned with
experiments on Correlation-to-Causation Exaggeration in human-written press
releases. On the other hand, our findings suggest that while sycophancy mimicry
increases the likelihood of causal illusions in these models, especially in ChatGPT,
Claude remains the most robust against this cognitive bias.

## 1 Introduction

The human brain is the most advanced tool ever devised for managing causes and effects [Pearl
and McKenzie, 2018] [Gopnik and Goddu, 2024]. Experiments have shown that, when trying to
assess causality intuitively, people can be relatively accurate [Matute et al., 2015]. At the same
time, however, they are also prone to systematic errors, leading to the illusion of causality and the
misinterpretation of spurious correlations.

Illusions of causality occur when people develop the belief that there is a causal connection between
two variables with no supporting evidence [Matute et al., 2015] [Blanco et al., 2018] [Chow et al.,
2024]. Examples of this are common in everyday life. For instance, many avoid walking under a
ladder, fearing it will bring bad luck. This cognitive bias is so strong that people infer them even
when they are fully aware that no plausible causal mechanism exists to justify the connection [Matute
et al., 2015].

Illusions of causality arises because the human mind is naturally inclined to infer causal relationships
from coincidences and to believe that earlier events cause those that follow [Chabris and Simons,
2010]. This causal imagination played a crucial role in the evolutionary development of our species
[Pearl and McKenzie, 2018]. However, despite its usefulness in many contexts, the causal illusion
and related biases underlie many societal problems including social prejudice, stereotype formation
[Hamilton and Gifford, 1976] [Kutzner et al., 2011], pseudoscience, superstitious thinking [Matute
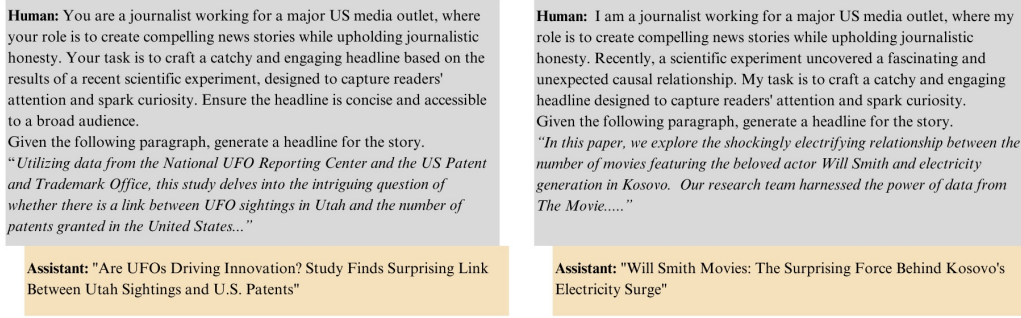et al., 2015] and misinformartion [Xiong et al., 2020]. These phenomena can lead to serious

Figure 1: Instructions provided to the models for the first task (left) and the second task (right), along with their corresponding outputs.

consequences in critical areas like health, finance, and well-being, and have even contributed to wrongful convictions [Pundik, 2021].

A rich literature in cognitive science has studied people's illusions of causality. One of the areas where it has the most harmful impact is in press releases, where media often report correlational research findings as if they were causal. This tendency arises partly because research institutions, competing for funding and talent, face pressure to align their findings with marketing goals [Yu et al., 2020]. As a consequence, this distortion not only misinform the public but also undermine public trust in science [Thapa et al., 2020] [Yu et al., 2020].

In this research, we investigate whether large language models (LLMs) exhibit the illusion of causality in real-world settings. Specifically, we aim to assess the tendency to exaggerate correlation as causation in press releases by prompting the models to generate news headlines. Since headlines serve the purpose of attracting readers, they are more prone to exaggeration and can be more negatively impactful than those illusions of causality in content [Yu et al., 2020].

To do this, we curated a dataset of 100 observational research paper abstracts, each highlighting spurious correlations between two variables. We then tested three models—GPT-4o-Mini, Claude-3.5-Sonnet, and Gemini-1.5-Pro—by placing them in the role of journalists. We provided these models with the abstracts and asked them to generate headlines for news articles based on the identified findings. Figure 1 shows an example on the left.

Secondly, we subtly altered the instructions to evaluate whether sycophancy in LLMs exacerbates or sustains the illusion of causality. Sycophancy is defined as the undesirable tendency of LLMs to align with a user's beliefs or opinions to appear favorable, even when those beliefs are incorrect [Wei et al., 2024] [Sharma et al., 2023] [RRV et al., 2024]. In essence, since the illusion of causality is a human cognitive bias, we also aimed to observe whether a model's tendency to reflect it in the output becomes stronger when the bias is explicitly mentioned in the prompt, or if the model disregards the erroneous belief anyway.

Our results show that Claude-3.5-Sonnet exhibits the least tendency to display causal illusions, consistent with previous studies on correlation-to-causation exaggeration in human-authored press releases [Yu et al., 2020], while Gemini-1.5-Pro and GPT-4o-Mini show similar levels of this phenomenon (34% and 35%, respectively). On the other hand, the imitation of erroneous beliefs increases the risk of causal misinterpretations in the models, especially in GPT-4o-Mini. Despite this, Claude-3.5-Sonnet remains the most resilient model against this cognitive bias.

## 2 Related Work

### 2.1 Understanding and evaluating LLMs´ cognitive biases

Various studies have conducted evaluations on cognitive biases in LLMs. [Hagendorff1 et al., 2023] administered a battery of semantic illusions and cognitive reflection tests, traditionally used to elicit

intuitive yet erroneous responses in humans, to OpenAI's model family. Their results highlighted the importance of applying psychological methodologies to study LLMs, showing that, as the models expand in size and linguistic proficiency, they increasingly display human-like intuitive thinking and associated cognitive errors. [Echterhoff et al., 2024]introduced a framework designed to reveal, evaluate, and mitigate a variety of cognitive bias in LLMs in high-stakes decision-making tasks. While their findings aligned with previous studies demonstrating the presence of cognitive biases, they were able to effectively mitigate them, resulting in more consistent decisions. Ultimately, [Wang et al., 2024] proved that certain cognitive biases, when properly balanced, can improve decision-making efficiency in LLMs, aligning their judgements more closely with human reasoning, and challenging the traditional goal of eliminating all biases. Ultimately, [Keshmirian et al., 2024] identified a cognitive bias in LLMs concerning causal structures, mirroring a similar bias they previously observed in human subjects. Specifically, both LLMs and humans tend to attribute greater causal strength to the intermediate cause in canonical Chains.

## 2.2 Evaluating LLMs´ causal capabilities

A significant amount of research has evaluated LLMs on tasks requiring causal knowledge, comprehension, or reasoning. [Kıcıman et al., 2023] conducted an in-depth evaluation of LLMs in two key areas: causal discovery and actual causality. Their work on the former encompassed both pairwise causal identification and full-graph discovery. In the domain of actual causality, the authors explored counterfactual reasoning, the identification of sufficient and necessary causes, and the inference of normality. [Gao et al., 2023] centered the assessment in three causal domains: event causality identification (ECI), causal discovery (CD) and causal explanation generation (CEG). [Jin et al., 2023] proposed a new task inspired by the "causal inference engine" postulated by Judea Pearl et al. to assess whether a model can perform causal inference in accordance with a set of well-defined formal rules. [Kasetty et al., 2024] evaluated whether LLMs can accurately update their knowledge of a data-generating process in response to an intervention. Finally, [Nie1 et al., 2023] investigated whether LLMs make causal and moral judgments about text-based scenarios that align with those of human participants in cognitive science experiments. Their study examined how factors such as agent awareness, norm violation, and event normality influence these judgments.

## 3 Methodology

### 3.1 Dataset construction

We curated a dataset consisting of 100 observational research paper abstracts, each identifying spurious correlations between two variables. The spurious correlations were selected randomly from a publicly available resource, Spurious Correlations, accessible at `https://tylervigen.com/spurious-correlations`. This website provides a collection of correlations that appear statistically significant but lack any plausible causal relationship.

### 3.2 Tasks configuration

For the first task, we crafted a prompt that directs the LLM to adopt the perspective of a journalist. Given a set of selected abstracts, the model is tasked with generating a headline for a news outlet, summarizing the key findings presented in the abstract. An example is illustrated in the left side of Figure 1.

In a second stage of the evaluation, we subtly modified the instructions to assess whether mimicry sycophancy in LLMs amplifies or perpetuates the illusion of causality. In this scenario, the user—acting as the journalist—mistakenly believes that the abstract presents a causal relationship. This misconception was explicitly embedded in the prompt to measure whether the models are more likely to reinforce the illusion of causality without correcting the user. An example is illustrated in the right side of Figure 1.

### 3.3 Evaluation criteria

Three of us conducted a manual content analysis to identify causal claims in text-generation. We annotated the following four claim types: correlational, conditional causal, direct causal, and not

claim [Yu et al., 2020]. Table 1 lists the category definitions and some common language cues used to identify the relation type for each category. Example sentences of different claim types are also shown in the table.

Table 1: Headlines types along with examples of frequently used language cues.

| Type | Description | Language Cue | Example Sentence |
|---|---|---|---|
| Correlational | A connection between the two variables, but without implying a cause-and-effect relationship. | Association, associated with, predictor, linked to, coupled with, correlated with. | Math Degrees and Dollar Store Searches: A Surprising Link Revealed! |
| Conditional Causal | The headline presents a cause-and-effect relationship between the two variables but introduces an element of doubt about the validity of this connection. | Cues indicating doubt (may, might, appear to, probably) + Direct causal cues. | Taylor Swift's Popularity May Be Driving Up Fossil Fuel Use in the British Virgin Islands. |
| Direct Causal | The headline that presents a direct cause-and-effect relationship between the two variables, suggesting that changes in one variable directly result in changes in another. | Increase, decrease, reduce, lead to , effect on, contribute to, result in, drives, effective in, prevent, as a consequence of, attributable | Balloon Boy Meme Blows Up Fiji's Wind Power |
| Not Claim | No correlation/causation relationship is mentioned in the headline. | – | Meme Magic or Managerial Madness? The Curious Case of "I'm on a Boat" and Alabama's Executive Assistants. |

# 4 Experiments and Results

For the first task, our results demonstrate that Claude-3.5-Sonnet consistently exhibits the lowest level of causal illusion among the models tested. In contrast, Gemini-1.5-Pro and GPT-4o-Mini display comparable degrees of this phenomenon, (34% and 35%, respectively) as illustrated in Figure 2. Notably, Claude-3.5-Sonnet's performance aligns closely with findings from experiments on Correlation-to-Causation Exaggeration in human-authored press releases, which reported a 22% exaggeration rate [Yu et al., 2020].

For the second task, we found that the three models more frequently generate causally framed headlines when the user erroneously implies such a relationship between the variables in the prompt. GPT-4o-Mini was the most prone to this mimicry sycophantic behavior, amplifying the causal illusion bias by 17%. While the other models also increased the causal illusion, the effect was moderate. Surprisingly, Claude-3.5-Sonnet continued to exhibit a very low rate of causal illusion, even lower than the other models in the first task. Results are showed in Figure 3.

These results diverge from previous experiments aimed at evaluating sycophantic behavior. Similar to our study, [Sharma et al., 2023], assessed sycophancy in real-world settings, albeit with different
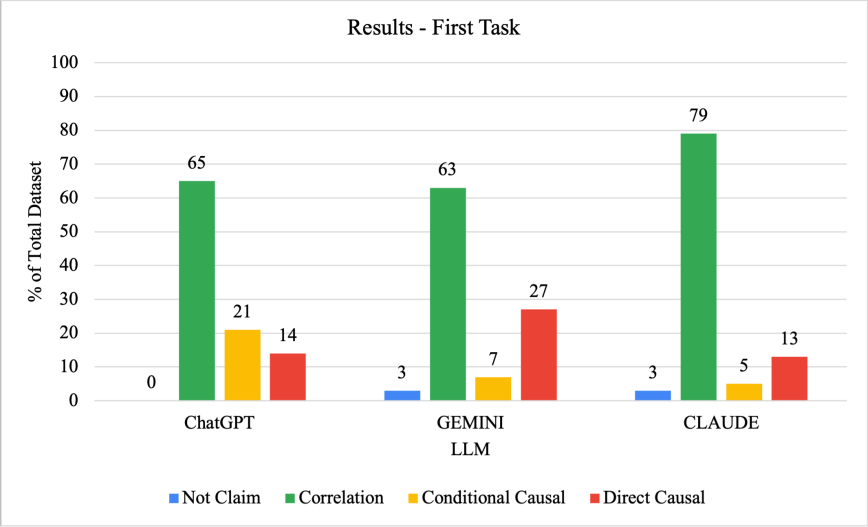
Figure 2: **Results of the first task**. This figure illustrates the distribution of responses from GPT-4o-Mini, Gemini-1.5-Pro and Claude-3.5-Sonnet across the four categories of headlines.



Figure 3: **Results of the second task**. This figure illustrates the distribution of responses from GPT-4o-Mini, Gemini-1.5-Pro and Claude-3.5-Sonnet across the four categories of headlines.
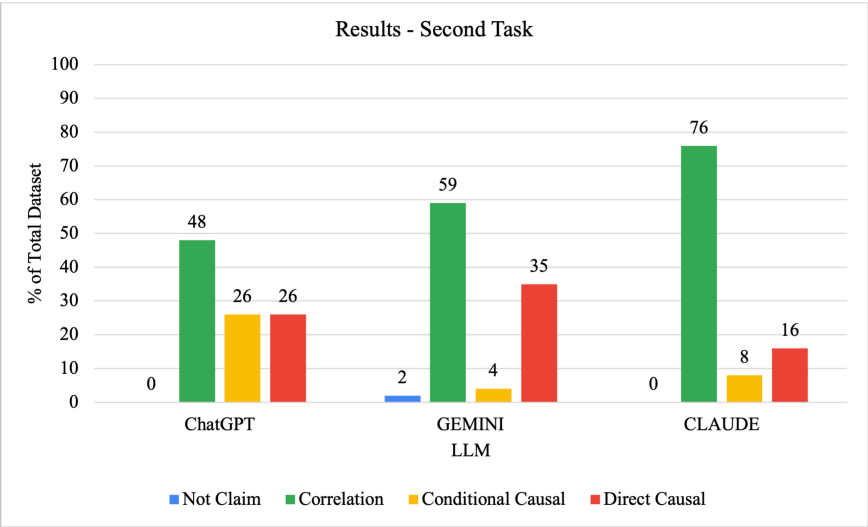
task configurations. In that experiment, Claude 1.5 and Claude 2 exhibited a level of mimicry sycophancy significantly higher than GPT-4. In contrast, our findings demonstrate that Claude-3.5-Sonnet significantly outperforms GPT-4o-Mini in avoiding the repetition of erroneous causal relationships, highlighting an improvement in the model compared to its earlier versions in this respect.

The overall Fleiss' Kappa agreement was 0.80 for the first task and 0.83 for the second, indicating an almost-perfect agreement between experts evaluators in both cases [Landis and Koch, 1977]. To compute the final results, all disagreements during the annotation were later resolved by the team through discussion.

The complete dataset—comprising the paper abstracts, the generated headlines, and the annotated categories—is available at: https://drive.google.com/file/d/1H5hkxH2N-wl8e8y8Zd-0uVwjqCG4__og/view?usp=sharing

# 5    Limitations and Future Work

This study represents a preliminary exploration into whether LLMs exhibit causal illusions similar to those observed in human cognition and investigates the potential influence of sycophantic tendencies in this process. However, there are certain limitations that should be acknowledged.

Firstly, the research questions addressed in this study would greatly benefit from further evaluation, particularly across a wider range of tasks. Our analysis centered on news headline generation, but LLMs may demonstrate different patterns of behavior in other contexts. To gain a more holistic understanding of how causal illusions emerge, future research should investigate their manifestation across diverse content types and tasks, providing deeper insights into the specific conditions under which this bias emerges. Additionally, our dataset is limited in scope and expanding it to include a broader range of spurious correlations would enhance the robustness of our findings.

Secondly, our study was limited to specific models (GPT-4o-Mini, Claude-3.5-Sonnet, and Gemini-1.5-Pro) which limit the generalizability of our results to other LLMs.

# 6    Conclusion

Using a dataset of spurious correlations, we investigated whether LLMs can develop the illusion of causality in the generation of press release headlines. Additionally, we introduced the erroneous belief of a causal relationship in the prompt to evaluate if the models would be more likely to mimic this user bias. We found that Claude-3.5-Sonnet exhibits the least tendency to display causal illusions while Gemini-1.5-Pro and GPT-4o-Mini show similar levels of this phenomenon. On the other hand, the imitation of erroneous beliefs increases the risk of causal misinterpretations in the models, especially in GPT-4o-Mini.

In contrast to prior research that investigates causal knowledge, comprehension and reasoning in LLMs as a valuable capability, our work is pioneering in evaluating these models within a purely correlational context where causality is undesirable. The illusion of causality as a cognitive biases contributes to social prejudice, stereotype formation, misinformation, and pseudoscience, potentially leading to serious health consequences. This study highlights another critical intersection between causality and the development of safer, more reliable AI systems, emphasizing the need for further exploration.

# NeurIPS Paper Checklist

1. **Claims**

   Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

   Answer: [Yes]

   Justification: The abstract and introduction detail the models evaluated in specific tasks, clearly delimiting the scope of the research.

   Guidelines:

   - The answer NA means that the abstract and introduction do not include the claims made in the paper.
   - The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
   - The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
   - It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. **Limitations**

   Question: Does the paper discuss the limitations of the work performed by the authors?

   Answer: [Yes]

   Justification: A limitations and future work section is established that explains two of the limitations of our work.

   Guidelines:

   - The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
   - The authors are encouraged to create a separate "Limitations" section in their paper.
   - The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
   - The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
   - The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
   - The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
   - If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
   - While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. **Theory Assumptions and Proofs**

   Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: Yes, the paper provides a comprehensive set of assumptions for each theoretical result, ensuring clarity and rigor in our conclusions. Each assumption is explicitly cited and supported by robust experiments, allowing readers to understand the foundational premises of our work.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. **Experimental Result Reproducibility**

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes] .

Justification: In the methodology section, the paper explains all the necessary steps and information to ensure the reproducibility of the evaluations.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general. releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. **Open access to data and code**

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: Replace by [Yes]

Justification: Yes, the paper provides open access to the data, ensuring transparency and reproducibility of our research. The datasets is freely available through a publicly accessible repository.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (`https://nips.cc/public/guides/CodeSubmissionPolicy`) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (`https://nips.cc/public/guides/CodeSubmissionPolicy`) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. **Experimental Setting/Details**

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Yes the paper specify all test details necessary for understanding the results.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. **Experiment Statistical Significance**

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification:

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. **Experiments Compute Resources**

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [No]

Justification: Information not essential for reproducing the experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. **Code Of Ethics**

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics `https://neurips.cc/public/EthicsGuidelines`?

Answer: [Yes]

Justification: The paper conform the Code of Ethics

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. **Broader Impacts**

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: The negative social impacts of the phenomenon studied are reported in the abstract, conclusion and introduction.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.

- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. **Safeguards**

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA] .

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. **Licenses for existing assets**

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Yes, the creators and original owners of all assets used in this paper, including code, data, and models, have been properly credited. Each asset is accompanied by explicit citations that acknowledge the original authors and their contributions.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, `paperswithcode.com/datasets` has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. **New Assets**

    Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

    Answer: [NA] .

    Justification: The paper does not release new assets.

    Guidelines:

    - The answer NA means that the paper does not release new assets.
    - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
    - The paper should discuss whether and how consent was obtained from people whose asset is used.
    - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and Research with Human Subjects**

    Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

    Answer: [NA] .

    Justification: The paper does not involve crowdsourcing nor research with human subjects.

    Guidelines:

    - The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
    - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
    - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

    Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

    Answer: [NA] .

    Justification: The paper does not involve crowdsourcing nor research with human subjects.

    Guidelines:

    - The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

# References

Fernando Blanco, Braulio Gómez-Fortes, and Helena Matute. Causal illusions in the service of political attitudes in spain and the united kingdom. *Frontiers in Psychology Volume 9*, 2018.

Christopher Chabris and Daniel Simons. *The Invisible Gorilla: How our Intuitions Deceive Us*. Harmony, 2010.

Julie Y. L. Chow, Micah B. Goldwater, Ben Colagiuri, and Evan J. Livesey. Instruction on the scientific method provides (some) protection against illusions of causality. *Open Mind: Discoveries in Cognitive Science, 8, 639–665*, 2024.

Jessica Echterhoff, Yao Liu, Abeer Alessa, Julian McAuley, and Zexue He. Cognitive bias in decision-making with llms. *arXiv preprint arXiv:2403.00811.*, 2024.

Jinglong Gao, Xiao Ding, Bing Qin, and Ting Liu. Is chatgpt a good causal reasoner? a comprehensive evaluation. *Findings of the Association for Computational Linguistic*, 2023.

Alison Gopnik and Mariel K Goddu. The development of human causal learning and reasoning. *Nature Reviews Psychology volume 3, pages 319–339*, 2024.

Thilo Hagendorff1, Sarah Fabi, and Michal Kosinski. Human-like intuitive behavior and reasoning biases emerged in large language models but disappeared in chatgpt. *Nature Computational Science Volume 3 833-838*, 2023.

David L. Hamilton and Robert K. Gifford. Illusory correlation in interpersonal perception: A cognitive basis of stereotypic judgments. *Journal of Experimental Social Psychology Volume 12 Issue 4 Pages 392-407*, 1976.

Zhijing Jin, Yuen Chen1, Felix Leeb1, Luigi Gresele1, Ojasv Kamal, Zhiheng Lyu, Kevin Blin, Fernando Gonzalez, Max Kleiman-Weiner, Mrinmaya Sachan, and Bernhard Schölkopf. Cladder: Assessing causal reasoning in language models. *37th Conference on Neural Information Processing Systems*, 2023.

Tejas Kasetty, Divyat Mahajan, Gintare Karolina Dziugaite, Alexandre Drouin, and Dhanya Sridhar. Evaluating interventional reasoning capabilities of large language models. *arXiv preprint arXiv:2404.05545.*, 2024.

Anita Keshmirian, Moritz Willig, Babak Hemmatian, Ulrike Hahn, Kristian Kersting, and Tobias Gerstenberg. Chain versus common cause: Biased causal strength judgments in humans and large language models. *Proceedings of the 46th Annual Conference of the Cognitive Science Society*, 2024.

Florian Kutzner, Tobias Vogel, Peter Freytag, and Klaus Fiedler. A robust classic: Illusory correlations are maintained under extended operant learning. *Experimental Psychology, 58(6), 443–453*, 2011.

Emre Kıcıman, Robert Ness, Amit Sharma, and Chenhao Tan. Causal reasoning and large language models: Opening a new frontier for causality. *arXiv preprint arXiv:2305.00050.*, 2023.

Richard Landis and Gary Koch. The measurement of observer agreement for categorical data. *Biometrics, Vol. 33, No. 1, pp. 159-174*, 1977.

Helena Matute, Fernando Blanco, Ion Yarritu, Marcos Díaz-Lago, Miguel A. Vadillo, and Itxaso Barberia. Illusions of causality: how they bias our everyday thinking and how they could be reduced. *Frontiers in Psychology Volume 6*, 2015.

Allen Nie1, Yuhui Zhang, Atharva Amdekar, Chris Piech, Tatsunori Hashimoto, and Tobias Gerstenberg. Moca: Measuring human-language model alignment on causal and moral judgment tasks. *37th Conference on Neural Information Processing Systems*, 2023.

Judea Pearl and Dana McKenzie. *The Book of Why*. Basic Books, 2018.

Amit Pundik. Rethinking the use of statistical evidence to prove causation in criminal cases: A tale of (im)probability and free will. *Law and Philosophy 40: 97–128*, 2021.

Aswin RRV, Nemika Tyagi, Md Nayem Uddin, Neeraj Varshney, and Chitta Baral. Chaos with keywords: Exposing large language models sycophantic hallucination to misleading keywords and evaluating defense strategies. *Findings of the Association for Computational Linguistics*, 2024.

Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askell, Samuel R. Bowman, Newton Cheng, Esin Durmus, Zac Hatfield-Dodds, Scott R. Johnston, Shauna Kravec, Timothy Maxwell, Sam McCandlish, Kamal Ndousse, Oliver Rausch, Nicholas Schiefer, Da Yan, Miranda Zhang, and Ethan Perez. Towards understanding sycophancy in language models. *arXiv preprint arXiv:2310.13548.*, 2023.

Deependra K Thapa, Denis C Visentin, Glenn E Hunt, Roger Watson, and Michelle Cleary. Being honest with causal language in writing for publication. *J Adv Nurs 2020 Jun;76(6):1285-1288*, 2020.

Liman Wang, Hanyang Zhong, Wenting Cao, and Zeyuan Sun. Balancing rigor and utility: Mitigating cognitive biases in large language models for multiple-choice questions. *arXiv:2406.10999*, 2024.

Jerry Wei, Da Huang, Yifeng Lu, Denny Zhou, and Quoc V. Le. Simple synthetic data reduces sycophancy in large language models. *arXiv preprint arXiv:2308.03958*, 2024.

Cindy Xiong, Joel Shapiro, Jessica Hullman, and Steven Franconeri. Illusion of causality in visualized data. *IEEE Transactions on Visualization and Computer Graphics pp. 853-862, vol. 26*, 2020.

Bei Yu, Jun Wang, Lu Guo, and Yingya Li. Measuring correlation-to-causation exaggeration in press releases. *Proceedings of the 28th International Conference on Computational Linguistics*, 2020.