

# **Can LLMs serve in identifying fake Health Information: it depends on how and who you ask.**

Francois Bolduc, Ashwani singla, Manpreet Kaur, Mohammad Reza Taesiri, Keneizha Rubanarayana, Abhishek Dhankar, Osmar Zaiane, Marek Reformat

Submitted to: Journal of Medical Internet Research  
on: October 08, 2024

**Disclaimer:** © The authors. All rights reserved. This is a privileged document currently under peer-review/community review. Authors have provided JMIR Publications with an exclusive license to publish this preprint on its website for review purposes only. While the final peer-reviewed paper may be licensed under a CC BY license on publication, at this stage authors and publisher expressly prohibit redistribution of this draft paper other than for review purposes.

Table of Contents

Original Manuscript..... 4

Supplementary Files..... 25

    Multimedia Appendixes ..... 26

        Multimedia Appendix 1..... 26

        Multimedia Appendix 2..... 26

        Multimedia Appendix 3..... 26

        Multimedia Appendix 4..... 26

        Multimedia Appendix 5..... 26

        Multimedia Appendix 6..... 26

# Can LLMs serve in identifying fake Health Information: it depends on how and who you ask.

Francois Bolduc<sup>1</sup> md, phd; Ashwani singla<sup>1</sup> msc; Manpreet Kaur<sup>1</sup> Msc; Mohammad Reza Taesiri<sup>1</sup> phd; Keneizha Rubanarayana<sup>1</sup>; Abhishek Dhankar<sup>1</sup> Mac; Osmar Zaiane<sup>1</sup>; Marek Reformat<sup>1</sup> PhD

<sup>1</sup>University of Alberta Edmonton CA

## Corresponding Author:

Francois Bolduc md, phd  
University of Alberta  
3020 Katz Group Centre  
Edmonton  
CA

## Abstract

Misleading information has significant implications for society but can have disastrous impact for health matters. Transformative artificial intelligence (AI) tools such as large language models (LLMs) have the potential for limitless content generation (including fake), soon making internet information impossible to assess using traditional human approaches. We asked if the same LLMs (GPT4 and Gemini1-5-Pro) could be part of a more scalable solution. We tested 2 publicly available LLMs for their ability to identify misinformation in HealthReleases previously labeled by human experts. We found that simple prompts lead to overall low accuracy (F1 Macro 0,45 (GPT4) and 0,49 (Gemini1-5Pro)), but very different profiles for each LLM. Adding specific criteria used by experts to critically assess the Releases enhanced Gemini (0.66) but surprisingly reduced GPT4 (0,37) performances. We therefore developed a novel approach incorporating summaries of expert feedback into prompts and then observed major improvements in performance for both LLMs(GPT4;0.63 and Gemini1-5Pro; 0.96). Our study provides the first use case of LLMs as high throughput proofing of medical text, but more importantly provides insights into LLMs' "truth biases". We provide a novel paradigm integrating knowledge into the prompts which may reduce the need for LLM training, and the requirement for ever larger datasets and compute power. Importantly, we show how experts could and need to be involved in LLMs used to enhance their performance and potentially minimize the data wall issue.

(JMIR Preprints 08/10/2024:67329)

DOI: <https://doi.org/10.2196/preprints.67329>

## Preprint Settings

1) Would you like to publish your submitted manuscript as preprint?

✓ **Please make my preprint PDF available to anyone at any time (recommended).**

Please make my preprint PDF available only to logged-in users; I understand that my title and abstract will remain visible to all users.

Only make the preprint title and abstract visible.

No, I do not wish to publish my submitted manuscript as a preprint.

2) If accepted for publication in a JMIR journal, would you like the PDF to be visible to the public?

✓ **Yes, please make my accepted manuscript PDF available to anyone at any time (Recommended).**

Yes, but please make my accepted manuscript PDF available only to logged-in users; I understand that the title and abstract will remain visible to all users.

Yes, but only make the title and abstract visible (see Important note, above). I understand that if I later pay to participate in [http://www.jmir.org/](#)

## Original Manuscript

Can LLMs serve in identifying fake Health Information: it depends on how and who you ask.

Ashwani Singla<sup>1</sup>, Manpreet Kaur<sup>1</sup>, Mohammed Reza Taesiri<sup>3</sup>, Keneizha Rubanarayana<sup>1</sup>, Abhishek Dhankar<sup>2</sup>, Osmar Zaiane<sup>2</sup>, Marek Reformat<sup>3,4</sup>, Francois V Bolduc<sup>1,5,6</sup>

1 Department of Pediatrics, University of Alberta, Edmonton, Alberta, Canada

2 Department of Computing Science, University of Alberta, Edmonton, Alberta, Canada

3 Electrical and Computer Engineering Department, University of Alberta, Canada

4 Information Technology Institute, University of Social Science, Poland

5 Neuroscience and Mental Health Institute, University of Alberta,

6 Women and Children Health Research Institute, University of Alberta

Misleading information has significant implications for society but can have disastrous impact for health matters. Transformative artificial intelligence (AI) tools such as large language models (LLMs) have the potential for limitless content generation (including fake), soon making internet information impossible to assess using traditional human approaches. We asked if the same LLMs (GPT4 and Gemini1-5-Pro) could be part of a more scalable solution. We tested 2 publicly available LLMs for their ability to identify misinformation in Health Releases previously labeled by human experts. We found that simple prompts lead to overall low accuracy (F1 Macro 0.45 (GPT4) and 0.49 (Gemini1-5Pro)), but very different profiles for each LLM. Adding specific criteria used by experts to critically assess the Releases enhanced Gemini (0.66) but surprisingly reduced GPT4 (0.37) performances. We therefore developed a novel approach incorporating summaries of expert feedback into prompts and then observed major improvements in performance for both LLMs (GPT4; 0.63 and Gemini1-5Pro; 0.96). Our study provides the first use case of LLMs as high throughput proofing of medical text, but more importantly provides insights into LLMs' "truth biases". We provide a novel paradigm integrating knowledge into the prompts which may reduce the need for LLM training, and the requirement for ever larger datasets and compute power. Importantly, we show how experts could and need to be involved in LLMs used to enhance their performance and potentially minimize the data wall issue.

**Keywords:** Large language model, GPT4, Gemini1-5Pro, artificial intelligence, natural language processing, LLM tuning, health, medicine, question answering, summarization, prompt engineering, news releases, misleading or fake news.

## Introduction

Online platforms have greatly democratized access to medical information<sup>1</sup>, but have also highlighted the need of evaluating the quality of information on unprecedented scales, in order to minimize the sharing of misleading information<sup>2,3</sup>. The rapid developments in generative artificial intelligence (AI), especially large language models (LLMs), make scalable content evaluation even more needed, while leaving health experts and institutions to redefine not only how information is critically appraised but also how to engage with AI generated content. But, misinformation, especially when using exaggeration, satire, parody, fabrication, mistakes<sup>4</sup> or emotionally charged messages<sup>5</sup> remains very challenging to detect for machines.

Nonetheless, there have been several attempts at using computers to detect misinformation. Pre-LLM, natural language processing (NLP)<sup>6-8</sup> consisted of linguistic feature extraction<sup>9,10</sup> combined with machine learning (ML) based multi-modal features-based misinformation detection<sup>11,12</sup>, and evidence-based fact-checking with deep learning algorithms (BERT, SciBERT, and RoBERTa)<sup>13,14,15</sup>. Using a dataset of news releases related to the health (and annotated by experts as true or misleading)<sup>16</sup>, fine tuning<sup>17</sup> and inclusion of linguistic features<sup>18</sup> to the ML model allowed for detection of misleading information but with relatively low accuracy (0.616- 0.658 MacroF1), as F1 above 0.7 is generally accepted (REF). But, integrating user and content-related information (number of followers, verified status, number of retweets and likes) enhanced accuracy (0.843 MacroF1)<sup>19</sup>. While promising, these required dedicated teams, large computing power, are domain specific and remain not accessible to most health professionals or the lay public.

LLMs have the potential to transform text generation<sup>20</sup> for medical queries, medical examinations, and medical assistants.<sup>21-23</sup> Combined with chatbots (as seen with OpenAI ChatGPT<sup>24</sup>) LLM can be intuitively accessed by both the general public<sup>25,26</sup> and health professionals.<sup>27</sup> Moreover, ChatGPT has recently been used for fact-checking general information.<sup>28,29</sup> On the other hand, LLMs have been shown to confabulate<sup>30-32</sup>, harbor gender and race biases<sup>33,34</sup>, or respond based on outdated, unreliable, or domain-nonspecific training data<sup>35,36</sup>. Moreover, LLMs' answers are influenced by the way the question (prompt) is designed.<sup>37,38,39</sup>

Here we show that publicly available LLM-chatbots, GPT4<sup>40</sup> and Gemini can detect misleading information from health-related news releases. We also elucidate key aspects in prompt engineering, including a novel approach integrating expert input in order to improve performances significantly. Our work shows how integrating expert knowledge in prompts could mitigate the data wall and thus help reduce compute needs.

## Method

**Dataset used for medical texts.** We leveraged the same published dataset<sup>16</sup> previously used in other publications<sup>17,18,19</sup> including 606 texts (HealthRelease) annotated by domain experts (journalism, medicine, health services research, and public health) with a ground truth value (315 rated by experts as true and 291 as fake) as well as individual labels on 10 standardized questions (criteria)<sup>41</sup> (**Supp. Table 1**). 14 resources (news\_reviews\_00018, 00136, 00168, 00235, 00284, 00395, 00422, 00490, 00499, 00507, 00589, 00590, 00601, 00605) out of 606 are not considered (no description text or no news\_source was provided).

**Prompt engineering.** We first tested a zero-shot prompting approach: (**Prompt1**) where GPT4 or Gemini-1.5-Pro are instructed to detect misleading information in the NewsRelease without providing any context. To restrict the output, it was asked to show only those claims which it can provide a reference for (details in **Supp. Method file**). For the second prompt (**Prompt2**), we included in the prompt a specific criteria (quality of evidence, sensational language, novelty of approach, disease mongering and benefits of treatment, existing alternatives, harms of intervention, availability of treatment, cost of intervention, and conflict of interest) which was used also previously by experts to rate the NewsReleases<sup>41</sup> (**Supp. Table 1**) as a context to make the model perform guided decisions along certain domain-specific dimensions. Additionally, the output instruction was modified to assign satisfactory or not satisfactory labels to each criteria and the overall text. For the third prompt (**Prompt3**), we first selected, for each criteria analyzed (quality of evidence, sensational language, novelty of approach, disease mongering and benefits of treatment), 50 NewsReleases labeled fake by experts and noted as not satisfactory for the criteria. We then asked GPT4 or Gemini-1.5-Pro to summarize into bullet point the justification provided by experts for marking a release as misleading (**Supp. Method File, Supp Table 2**).

We then incorporated the summarization into the prompt. The remaining five criteria (existing



alternatives, harms of intervention, availability of treatment, costs of intervention, and conflict of interest) were kept the same as prompt 2.

**Query process.** We queried GPT4 and Gemini1-5-Pro using the requests library in Python following the request/response protocol. Initially, all the claim text description was stored in an Excel sheet which was taken as input by the Python script. A list of claim texts was iterated and a prompt was created by embedding claim text descriptions into prompt templates to obtain the prompts for both experiment settings. Each prompt was passed as a message to Open AI or Gemini API calls along with other required API parameters with default values (role = user, frequency\_penalty = 0.0, and temperature 0.0/0.2 for GPT4 and default and temperature 0/1 for Gemini1-5Pro). The API returned a list of responses for a given prompt, and only the first response in the list was considered for further analysis. For prompt 3, we utilized LangChain framework<sup>42</sup> for GPT to create chat templates and the json output parser. Code is publicly available at GitHub<sup>43</sup>.

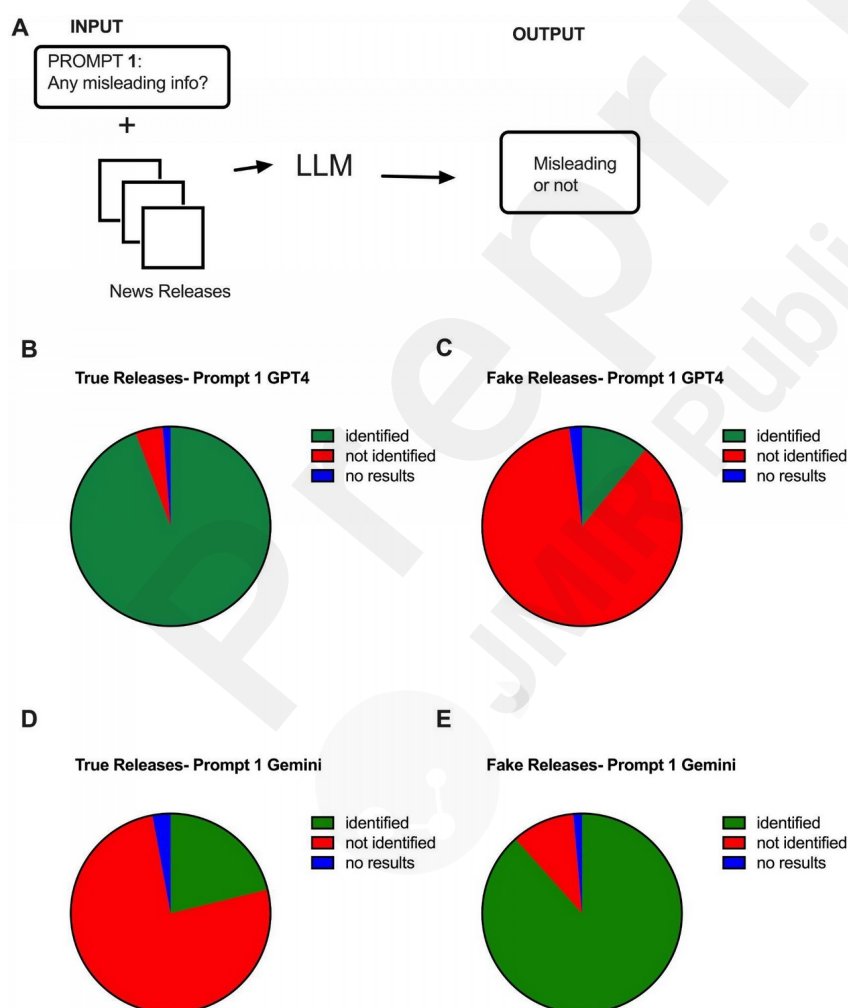
**Analysis:** for prompt 1 and 2, GPT-4 and Gemini responses were further analyzed (using a Python script followed by manual result validation) to create a specific label (satisfactory or not satisfactory) for each output. For evaluating the prompts, the analysis of Identifying Fakes we used the following definitions: True Positive (TP): The Number of fakes identified correctly. False Positive (FP): The number of true items identified as fakes. True Negative (TN): The number of true items identified correctly as true. False Negative (FN): The number of fakes identified incorrectly as true. For the analysis of Identifying Trues: True Positive (TP): The number of true items identified correctly as true. False Positive (FP): The number of fakes identified incorrectly as true. True Negative (TN): The number of fakes identified correctly. False Negative (FN): The number of true items identified as fakes. We used those to calculate the Precision =  $TP / (TP + FP)$  and the Recall =  $TP / (TP + FN)$ . The F1 score was calculated as before  $2 * ((precision * recall) / (precision + recall))$ . Macro average of the classification metrics (precision, recall and f1) is calculated by taking the unweighted mean of the metrics for each label using scikit-learn library<sup>44</sup>.

## Results

We started by asking if GPT4 or Gemini could label a NewsRelease as containing or not misleading information (Prompt 1, zeroshot) (**Fig. 1A**). Importantly, those NewsReleases were labeled for ground truth by experts and used in ML based experiments for fake detection<sup>16</sup>.

Importantly the dataset contained both true (315) and fake (291) resources. We found that GPT-4 matched expert annotation for true news in 90.6% of cases but was only able to identify 11.2% of the Fake ones, resulting in a low overall performance (as indicated by MacroF1 of 0.45) (**Fig.**

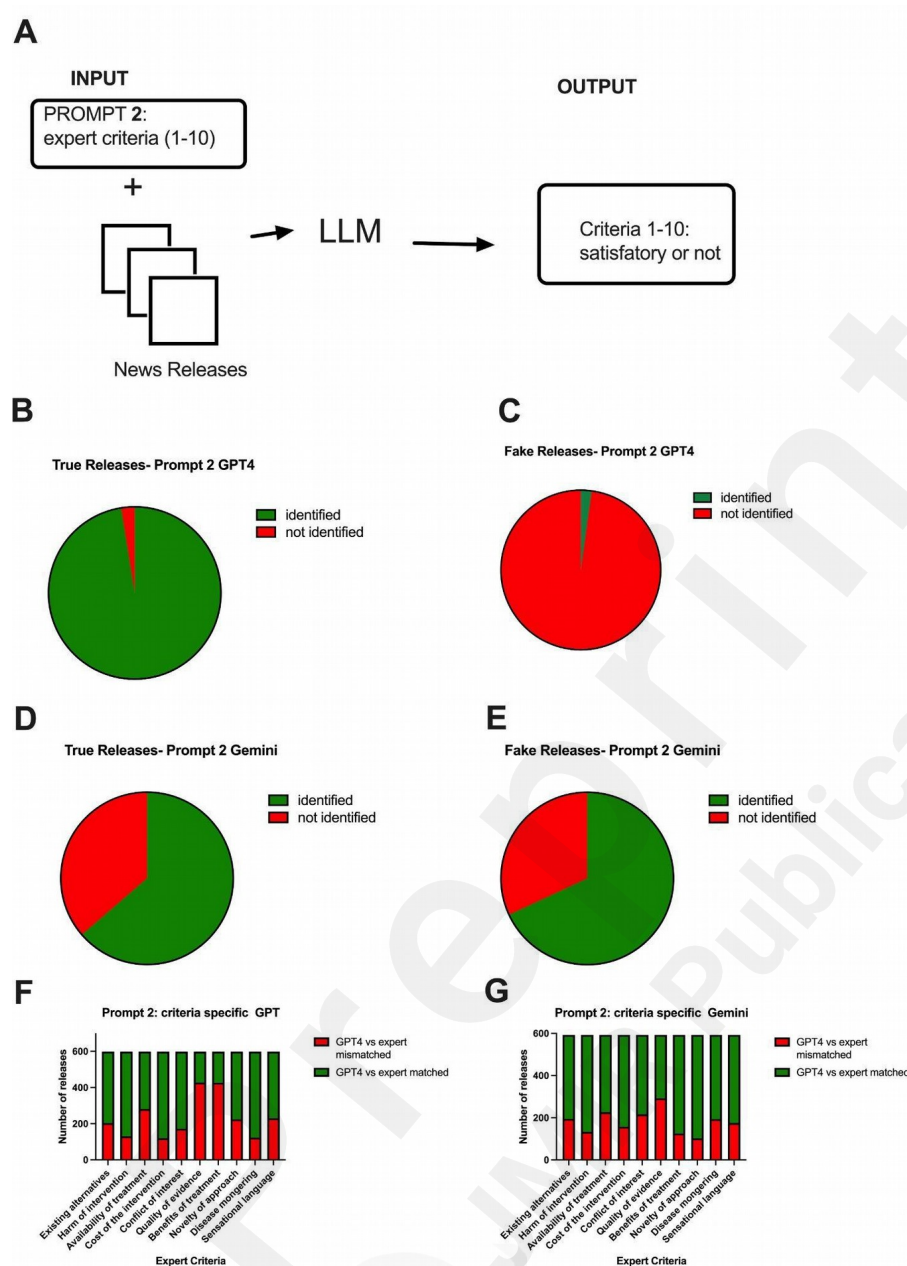
**1B-C; Supp. Table 3-4**). Surprisingly, Gemini1-5 Pro presented with an almost mirror image, low performance for True (21.2%) and high performance (88.4%) for Fake, leading to a similar Macro F1 of 0.49 (**Fig. 1D-E; Supp. Table 3-5**). For both models, we found no significant differences by varying temperature from 0 to default settings (**Supp. Table 4-5**).



**Figure 1. LLMs true/fake detection. Prompt 1: Identify misleading information. A)** Schematic representation of the workflow. GPT4 or Gemini are asked to assess the NewsReleases for misleading information. **B)** For Releases labeled by experts as true, GPT4 identified 90.6% as accurate, while reporting 4.2% of them as inaccurate. **C)** But, GPT4 only detected 11.2% of the Releases labeled as Fake by experts. **D)** In comparison, Gemini1-5Pro identified only 21.2% of True Releases as such but was able to recognize misleading information in 88.4% of Fake Releases.

Next, we postulated that using prompts similar to criterias used by experts to evaluate the NewsRelease would enhance LLMs performance (**Fig. 2 A**). We therefore use the same criteria and ground truth scoring (more than 6/10 criteria scored as misleading is labeled Fake) to compare LLM and expert labels. Surprisingly, we found that GPT4 performance increased for its ability to match experts on true news (97.5%) but worsened Fake detection (2.2%), resulting in low MacroF1 (0.37) (**Fig. 2B-C**); **Supp. Table 3-4**). For Gemini1-5 Pro, we noticed a major improvement of performance for True (63.6%) but reduced for Fake (68.1%) nonetheless resulting in a more balanced performance and improved MacroF1 of 0.66 (**Fig. 2D-E**); **Supp.**

**Table 3-5**).



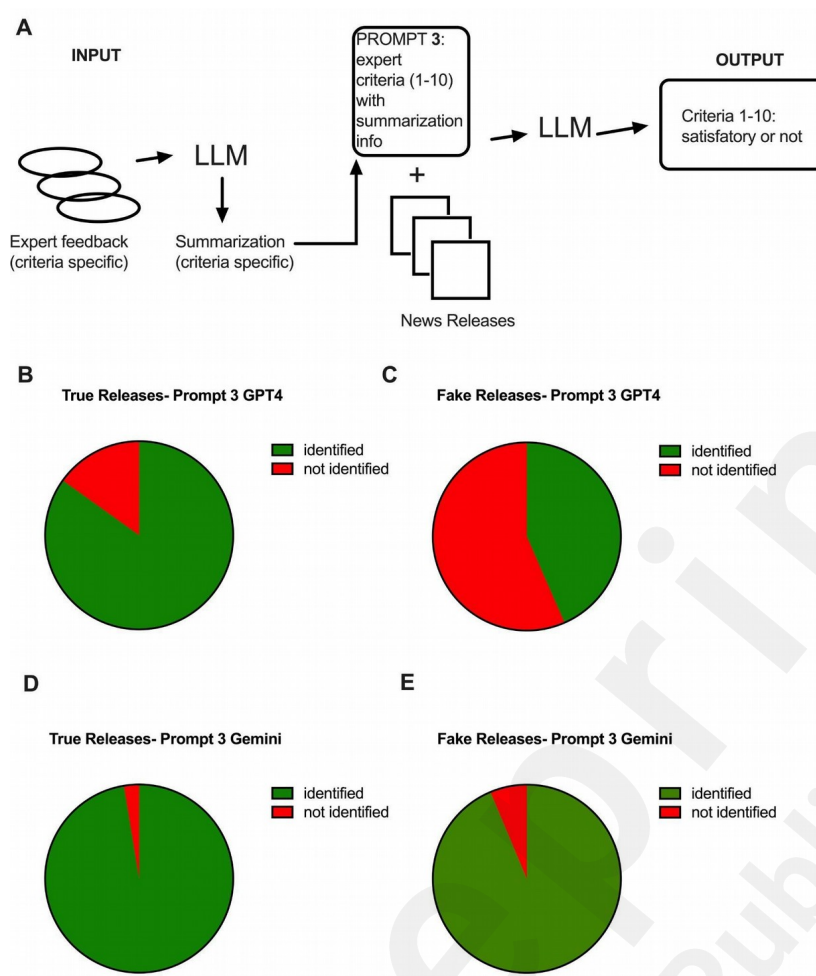
**Figure 2. Criteria guided queries reveal uneven performance of LLMs.** A) Schematic representation of the workflow. The LLM has to mark the release as satisfactory (or not) on specific criteria (asserted by experts on each release). 6 or more not satisfactory results in the Release being marked as Fake. GPT was able to identify **B) True releases in 97.5%** but **C) Fake in only 2.2%**. On the other hand, Gemini identified **D) true in 63.6 %** and **E) fake in 68.1 %**. Performance of **F) GPT4** and **G) Gemini** reveals heterogeneity in LLM responses based on individual criteria.

Examining performance on individual criterias revealed heterogeneous performance (**Fig. 2F-G**) in GPT4 and to a lesser extent in Gemini1-5Pro. Interestingly, examining justifications provided by both models (GPT4 for Fake and Gemini1-5Pro for True) revealed differences in bias toward critical thinking. For GPT4, when assessing disease mongering, experts raised the issue of the prevalence of conditions either lacking, not being supported by references, or cherry-picking the highest prevalence, when GPT4 thought the information provided was satisfactory. Also, controversies about disorders or treatments were not flagged by GPT4 but were noted by experts.

Similarly, in release fails related to the quality of evidence, GPT4 did not achieve enough critical analysis and did not pick on limitations (study design, power, data quality, outcome measures), methodological flaws (statistical analysis), or exaggeration of the significance of the results (translating findings, lack of peer review). In many cases, GPT4 marked the releases as true if the authors stated those limitations somewhere in the text instead of reporting the NewsRelease as misleading. Some issues in GPT4 understanding nuances about language became more explicit by looking at the question about Emotional or inflammatory language. Here, words like “breakthrough, cure, or revolutionary” were felt to be unsubstantiated by experts but marked as quotes and therefore not judged misleading by GPT4. Interestingly, GPT4 commented that those quotes were “optimistic or exciting” or justified considering the promising research instead of flagging them as inappropriate.

On the other hand, Gemini1.5-Pro struggled with True news by being very critical: reporting News as overly optimistic based on the evidence provided and reporting that authors were using language which was not conservative enough for health related topics.

So we hypothesized that by fine tuning the prompt (instead of the LLM itself) with input from the experts (through summarization), we could enhance (rebalance) performance for both LLMs (**Fig. 3 A**). Iterative testing revealed that summarization of experts' feedback from misleading statements led to best performance (**Supp. Table 3**). For GPT4, prompt 3 boosted overall detection of fake releases to 43.4% but deflated the performance for true news (84.9%) (**Fig. 3 B-C**) (**Macro F10.63**) (**Supp. Table 3-4**). For Gemini1-5Pro, we saw improvement in both True (97.4%) and Fake detection (93.7%) resulting in overall F1Macro 0.96 (**Supp. Table 3-5**).



**Figure 3. Integrating critical thinking into the prompt enhances LLMs performance.** A) Schematic representation of workflow. The feedback provided by experts as to why they rated a News Release as not satisfactory for an individual criteria (N=50 releases summarized per criteria) are summarized and integrated into the criteria specific prompt (3). GPT4 performance on B) True (84.9%) and C) Fake (43.4%) News Releases. Gemini1-5Pro performance on D) True (97.4%) and E) Fake (93.7%) releases.

## Discussion

Finding ways of scaling the identification of misleading information remains challenging but is becoming a crucial aspect related to the value of online information. While we better appreciate now how misleading information ranging from “Fake News” to misinformation campaigns can be impactful on society and public health, we can only imagine how disastrous it could be when it comes to health. Moreover, generative AI opens a totally new level of capabilities with almost unlimited voluntary (or not) machine based production of misleading information.

So, we asked if LLM could be used to evaluate information quality but also sought to better understand their inner workings. This is not trivial considering the limitations inherent to LLM: being predictive models, prone to confabulation, grown on different sets of data, released at different stages of validation and rapidly evolving. Recent work using LLM (ChatGPT) to assert misinformation in political claims<sup>28,29,45,46</sup>, education<sup>47,48</sup>, social media post on COVID<sup>49,50</sup> shows promising results but their usefulness for health related news remains untested. Using a dataset of News Releases labeled by experts as either True or Fake, we tested two publicly available LLMs (GPT4 and Gemini1-5Pro) for their ability at detecting misleading information. Not surprisingly considering recent work on prompt contextualization<sup>51,52,39</sup>, the first zero shot approach led to relatively low accuracies but revealed how different LLMs would respond in very different ways: GPT4 labeled most articles as true, and therefore missed most of the Fake ones whereas Gemini1-5Pro proved overly critical and therefore labeled most true articles as containing misleading information. This could be due to differences in training data for each model, but recent work on another LLM, CLAUDE (Golden Gate version)<sup>53</sup> showed how individual entities could be assigned more weight and completely change the LLM’s output. Another surprise was the relatively poor performance of both models when including specific scientific criteria as prompt in order to identify misleading information, showing that simply including expert-based guidance in prompt development was not sufficient in achieving accurate rating of the News Release. This is important as it means that the inner knowledge of the LLM exert a major effect on its output, and would suggest that LLMs be grown on large amounts of data for complex applications such as science and health, adding further to the issue with compute needs/cost.

But, we found that injecting insights from experts on why they judged a NewsRelease was able for both model to “rebalance” and enhance their performance to level not previously achieved by ML approaches (Gemini1-5 Pro MacroF1 0.96), giving us 1) a novel approach to prompt design, 2) a novel approach for LLM fine tuning which would reduce compute needs and of the potential 3) away forward for expert to be involved in AI progresses instead of sitting in margin. Our approach may help deal with the data wall

There are still several unresolved questions regarding the evolution of LLMs. For instance, their acceptability considering their rapid evolution, the relative opacity around the training data (including use of copyrighted material, and publicly created content) and our still basic understanding of their growth process. The emergence of on-device LLMs, while providing potentially enhanced privacy, will also make the evolving nature of LLMs an important aspect to study.

## Data Availability

We utilized the openly available Fake Health dataset which can be accessed at<sup>54</sup>. All the analysis scripts and datasets used for automatically labeling the GPT4 response are available in the public repository that can be accessed at.<sup>55</sup>

## Ethics declarations

The authors declare no competing interests.

## Author contributions:

AS and MK performed the experiments and co-wrote the manuscript. MRT performed the experiments and assisted with the manuscript. KR assisted with claim annotation. AD and OZ identified initial resources for the dataset as misleading, contributed to the project's conceptual development, and provided critical input into the manuscript. MR participated in design and



provided critical input into the manuscript. FVB conceptualized the project, analyzed the data, obtained funding, and co-wrote the manuscript.

## **Supplementary material**

### **Supplementary Method file**

### **Supplementary Table 1. List of criteria used by experts**

### **Supplementary Table 2. Summaries GPT4 for prompt 3**

### **Supplementary Table 3. Results metrics**

### **Supplementary Table 4. Results GPT4**

### **Supplementary Table 5. Results Gemini1-5PRO**

## **References**

1. *Pew Research Center. Internet Health Resources. Pew Research Center: Internet, Science & Tech* <https://www.pewresearch.org/internet/2003/07/16/internet-health-resources/> (2003)

2. Suarez-Lledo, V. & Alvarez-Galvez, J. Prevalence of Health Misinformation on Social Media: Systematic Review. *J. Med. Internet Res.* **23**, e17187 (2021).
3. Swire-  
Thompson, B. *Public Health and Online Misinformation: Challenges and Recommendations*. (2020).

4. Dhankar, A., Zaiane, O.R. & Bolduc, F. UofA-TruthatFactify2022: Transformer And Transfer Learning Based Multi-Modal Fact-Checking. (2022).
5. Ahmad, A.R. & Murad, H.R. The Impact of Social Media on Panic During the COVID-19 Pandemic in Iraqi Kurdistan: Online Questionnaire Study. *J. Med. Internet Res.* **22**, e19556 (2020).
6. Sarrouti, M., Ben Abacha, A., M'rabet, Y. & Demner-Fushman, D. Evidence-based Fact-Checking of Health-related Claims. in *Findings of the Association for Computational Linguistics: EMNLP 2021* 3499–3512 (2021).
7. Tan, N. et al. *Multi2Claim: Generating Scientific Claims from Multi-Choice Questions for Scientific Fact-Checking*. in *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics* 2652–2664 (2023).
8. Barve, Y. & Saini, J.R. Detecting and classifying online health misinformation with 'Content Similarity Measure (CSM)' algorithm: an automated fact-checking-based approach. *J. Supercomput.* **79**, 9127–9156 (2023).
9. Pritam Deka Queen's University Belfast, UK, Anna Jurek-Loughrey Queen's University Belfast, UK & Deepak Queen's University Belfast, UK. Unsupervised Keyword Combination Query Generation from Online Health Related Content for Evidence-Based Fact Checking. <https://dl.acm.org/doi/10.1145/3487664.3487701> doi:10.1145/3487664.3487701.
10. Samuel, H. & Zaiane, O. MedFact: Towards Improving Veracity of Medical Information in Social Media Using Applied Machine Learning. *Advances in Artificial Intelligence* 108–120 (2018).
11. Silva, A., Luo, L., Karunasekera, S. & Leckie, C. Embracing Domain Differences in Fake

News: Cross-domain Fake News Detection using Multi-modal Data. *AAAI* **35**, 557–565 (2021).

12. Verschuuren, P.J., Gao, J., van Eeden, A., Oikonomou, S. & Bandhakavi, A. Logically at Factify2: A Multi-Modal Fact Checking System Based on Evidence Retrieval techniques and Transformer Encoder Architecture. (2023).
13. Sarrouiti, M., Ben Abacha, A., M'rabet, Y. & Demner-Fushman, D. Evidence-based Fact-Checking of Health-related Claims. in *Findings of the Association for Computational Linguistics: EMNLP 2021* 3499–3512 (2021).
14. Xie, Q. *et al.* Faithful AI in Medicine: A Systematic Review with Large Language Models and Beyond. *medRxiv* doi:10.1101/2023.04.18.23288752.
15. The state of human-centered NLP technology for fact-checking. *Inf. Process. Manag.* **60**, 103219 (2023).
16. *FakeHealth: This Repository (FakeHealth) Is Collected to Address Challenges in Fake Health News Detection.* (Github).
17. Debunking health fake news with domain specific pre-trained model. *Global Transitions Proceedings* **2**, 267–272 (2021).
18. Di Sotto, S. & Viviani, M. Health Misinformation Detection in the Social Web: An Overview and a Data Science Approach. *Int. J. Environ. Res. Public Health* **19**, 2173 (2022).
19. Rishabh Upadhyay Department of Informatics, Systems, and Communication, University of Milano-Bicocca, Italy, Gabriella Pasi Department of Informatics, Systems, and Communication, University of Milano-Bicocca, Italy & Marco Viviani Department of Informatics, Systems, and Communication, University of Milano-Bicocca, Italy. Leveraging

- Socio-contextual Information in BERT for Fake Health News Detection in Social Media. <https://dl.acm.org/doi/10.1145/3599696.3612902>doi:10.1145/3599696.3612902.
20. Ziyu, Z. *et al.* Through the Lens of Core Competency: Survey on Evaluation of Large Language Models. in *Proceedings of the 22nd Chinese National Conference on Computational Linguistics (Volume 2: Frontier Forum)* 88–109 (2023).
  21. Chang, Y. *et al.* A Survey on Evaluation of Large Language Models. (2023).
  22. Clusmann, J. *et al.* The future landscape of large language models in medicine. *Communications Medicine* **3**, 1–8 (2023).
  23. Thirunavukarasu, A. J. *et al.* Large language models in medicine. *Nat. Med.* **29**, 1930–1940 (2023).
  24. Panda, S. & Kaur, N. Exploring the viability of ChatGPT as an alternative to traditional chatbot systems in library and information centers. *Library Hi Tech News* **40**, 22–25 (2023).
  25. Xu, R., Feng, Y. & Chen, H. ChatGPT vs. Google: A Comparative Study of Search Performance and User Experience. (2023).
  26. Marita Skjuve SINTEF Digital, Norway, Asbjørn Følstad SINTEF, N. & Petter Bae Brandtzaeg SINTEF, Norway and University of Oslo, Norway. The User Experience of ChatGPT: Findings from a Questionnaire Study of Early Users. <https://dl.acm.org/doi/10.1145/3571884.3597144>doi:10.1145/3571884.3597144.
  27. Zakka, C. *et al.* Almanac: Retrieval-Augmented Language Models for Clinical Medicine. *Research Square* doi:10.21203/rs.3.rs-2883198/v1.
  28. Hoes, E., Altay, S. & Bermeo, J. Leveraging ChatGPT for efficient fact-checking. *PsyArXiv*. April **3**, (2023).
  29. Leite, J. A., Razuvayevskaya, O., Bontcheva, K. & Scarton, C. Detecting Misinformation

with LLM-Predicted Credibility Signals and Weak Supervision. (2023).

30. Emsley, R. ChatGPT: these are not hallucinations—they're refabrications and falsifications. *Schizophrenia* **9**, 1–2 (2023).
31. Alkaissi, H. & McFarlane, S. I. Artificial Hallucinations in ChatGPT: Implications in Scientific Writing. *Cureus* **15**, e35179 (2023).
32. Azamfirei, R., Kudchadkar, S. R. & Fackler, J. Large language models and the peril of their hallucinations. *Crit. Care* **27**, 1–2 (2023).
33. Website.
34. Omiye, J. A., Lester, J. C., Spichak, S., Rotemberg, V. & Daneshjou, R. Large language models propagate race-based medicine. *npj Digital Medicine* **6**, 1–4 (2023).
35. Summary of ChatGPT-Related research and perspective toward the future of large language models. *Meta-Radiology* **1**, 100017 (2023).
36. DeAngelis, L. *et al.* ChatGPT and the rise of large language models: the new AI-driven infodemic threat in public health. *Front. Public Health* **11**, 1166120 (2023).
37. Giray, L. Prompt Engineering with ChatGPT: A Guide for Academic Writers. *Ann. Biomed. Eng.* **51**, 2629–2633 (2023).
38. Meskó, B. Prompt Engineering as an Important Emerging Skill for Medical Professionals: Tutorial. *J. Med. Internet Res.* **25**, e50638 (2023).
39. White, J. *et al.* A Prompt Pattern Catalog to Enhance Prompt Engineering with ChatGPT. (2023).
40. Product. <https://openai.com/product>.
41. Dai, E., Sun, Y. & Wang, S. Ginger Cannot Cure Cancer: Battling Fake Health News with a Comprehensive Data Repository. (2020).

42. Introduction.[https://python.langchain.com/v0.1/docs/get\\_started/introduction/](https://python.langchain.com/v0.1/docs/get_started/introduction/).

43. Kaur,M.*MedicalFactChecker*.(Github).

44. *Classification\_report.scikit-learn*

[https://scikit-learn.org/stable/modules/generated/sklearn.metrics.classification\\_report.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.classification_report.html).

45. Satapara,S.,Mehta,P.,Ganguly,D.&Modha,S.FightingFirewithFire:AdversarialPrompting to Generate a Misinformation Detection Dataset. (2024).

46. Ni,J.*etal*.AFaCTA: AssistingtheAnnotationofFactualClaimDetectionwithReliableLLM Annotators. (2024).

47. O'Connor,S.*etal*.Promptengineering whenusinggenerativeAIinnursingeducation. *Nurse Educ. Pract.***74**, 103825(2024).

48. Heston,T.F.*PromptEngineering:ForStudentsofMedicineandTheirTeachers*. (Independently Published, 2023).

49. Choi,E.C.&Ferrara,E.FACT-GPT:Fact-CheckingAugmentationviaClaimMatchingwith LLMs. (2024).

50. Xie,Q.*etal*.FaithfulAIinMedicine:ASystematicReviewwithLargeLanguageModelsand Beyond.*medRxiv*doi:10.1101/2023.04.18.23288752.

51. Yao, Z., Cao, Y., Yang, Z., Deshpande, V. & Yu, H. Extracting Biomedical FactualKnowledgeUsingPretrainedLanguageModelandElectronicHealthRecordContext. *AMIA Annu. Symp. Proc.***2022**, 1188 (2022).

52. Wang,J.*etal*.PromptEngineeringforHealthcare:MethodologiesandApplications.(2023).

53. GoldenGateClaude.<https://www.anthropic.com/news/golden-gate-claude>.

54. GitHub - EnyanDai/FakeHealth: This repository (FakeHealth) is collected to address

*challenges in Fake Health News detection.GitHub*

<https://github.com/EnyanDai/FakeHealth>.

55. *GitHub - ashwani227/gpt4AnalysisFakeHealth.GitHub*

<https://github.com/ashwani227/gpt4AnalysisFakeHealth>.



## Supplementary Files

## Multimedia Appendixes

Supplementary table 1.

URL: <http://asset.jmir.pub/assets/220e56921d1a84550c4b5d95fafa8e77.xlsx>

Supplementary Table 2.

URL: <http://asset.jmir.pub/assets/3ca4786b273ab8505cbfcfa0bae06bd6.xlsx>

Supplementary Table 3.

URL: <http://asset.jmir.pub/assets/c8214c19a3a276af8017e8cac4219be2.pdf>

Supplementary Table 4.

URL: <http://asset.jmir.pub/assets/c4399230caaada2e0f5dc7df6a0c4279.xlsx>

Supplementary table 5.

URL: <http://asset.jmir.pub/assets/2afced551fa6e7fa2b2cc55846e88c0e.xlsx>

Supplementary Method file.

URL: <http://asset.jmir.pub/assets/9081f738ecc5b25519af6d05f22e1626.docx>