

# AI in the Newsroom: Analyzing the Increase in ChatGPT-Favored Words in News Articles

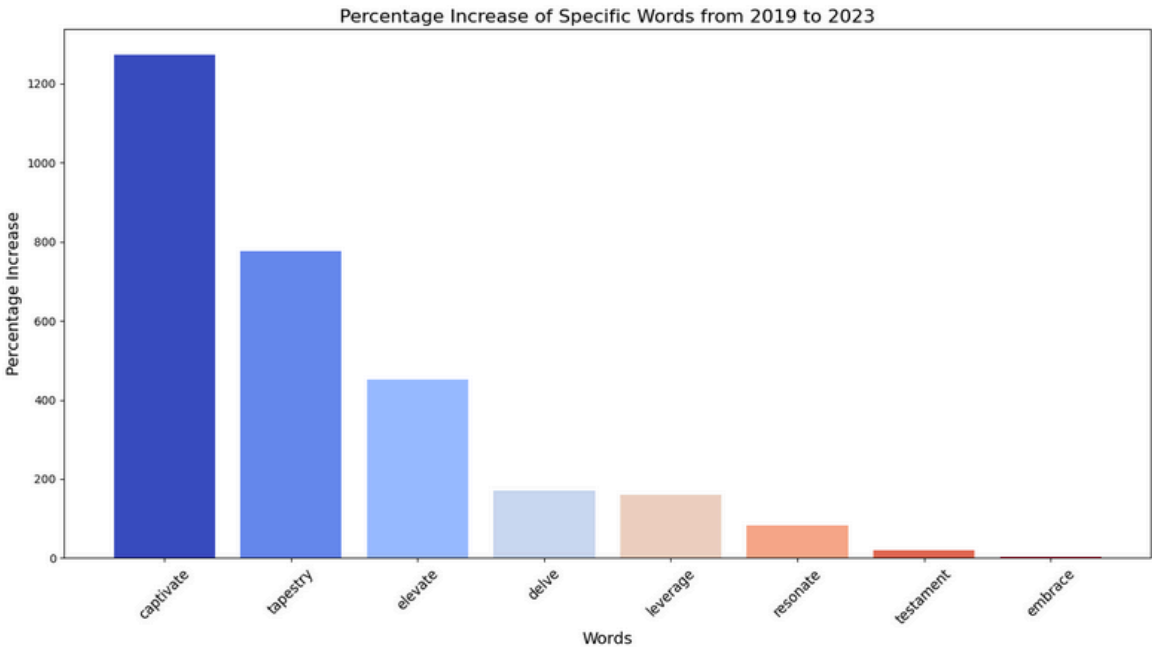
**Proliferation of AI-generated news**  
Media plays a vital role in informing the public and upholding accountability in democracy. However, recent trends indicate declining trust in the media, and there are increasing concerns that artificial intelligence might exacerbate this trend. In addition to direct harms from misinformation, the proliferation of AI-generated inauthentic content could result in people not knowing what to believe in. This could further lead to lower trust in established institutions.

**Increase in the occurrence of words favored by ChatGPT in news articles**  
Our demonstration shows a significant increase in the occurrence of words favored by ChatGPT in news articles from reputable sources in 2023 compared to previous years. The findings provide tentative evidence of LLMs potentially being used to produce news articles.

**Capability advances imply a continuation in the trend, unless precautions are taken**

As LLMs become more and more capable, we expect the trend suggested

Figure 1: Overall change of the frequency of the words in news articles favored by ChatGPT over time (2019-2023)



by our findings to continue. This could result in lower trust in the media, presenting a risk to democracy. However, it is also possible that LLM watermarking becomes technically viable in the coming years, and regulation is put in place in order to mitigate the risks.

**Mitigation strategies**  
We propose further research to LLM watermarking, which would help to identify whether a particular article is written by an LLM. Recent research suggests that there is still much work to be done before such technology is ready to be deployed. Furthermore, we encourage policymakers to invent ways of ensuring that LLM watermarking, once technically feasible, is used in the most efficient way possible to ensure that citizens' trust in institutions is protected. Examples of such efforts would be mandating news agencies to tell their users about the usage of LLMs, and giving users themselves easy-to-use tools for detecting whether a given source is written by AI.

# Appendix

## AI in the Newsroom: Analyzing the Increase in ChatGPT-Favored Words in News Articles

Aayush Kucheria<sup>1\*</sup>, Okko Katajamäki<sup>1\*</sup>, Santeri Koivula<sup>1\*</sup>, Andrea La Mantia<sup>2\*</sup>, Norman Piotrowski<sup>3\*\*</sup>

\*Co-first authorship  
\*\*Second authorship

### Abstract

Media plays a vital role in informing the public and upholding accountability in democracy. However, recent trends indicate declining trust in media, and there are increasing concerns that artificial intelligence might exacerbate this through the proliferation of inauthentic content. We investigate the usage of large language models (LLMs) in news articles, analyzing the frequency of words commonly associated with ChatGPT-generated content from a dataset of 75,000 articles. Our findings reveal a significant increase in the occurrence of words favored by ChatGPT after the release of the model, while control words saw minimal changes. This suggests a rise in AI-generated content in journalism.

<b>Abstract.....</b>	<b>1</b>
<b>1. Introduction.....</b>	<b>2</b>
<b>2. Methodology.....</b>	<b>2</b>
<b>3. Results.....</b>	<b>3</b>
<b>4. Discussion.....</b>	<b>6</b>
<b>5. Acknowledgements.....</b>	<b>7</b>
<b>6. References.....</b>	<b>7</b>

<sup>1</sup> Aalto University  
<sup>2</sup> University of Helsinki  
<sup>3</sup> Liverpool John Moores University

# 1. Introduction

Trust plays a critical part in democracy, forming the foundation for meaningful exchange between citizens and their government. Here democratic trust involves two aspects: citizens trusting institutions, such as the media and the government, and political leaders having confidence in the accuracy of information about public viewpoints.

Media is often referred to as the fourth estate due to its crucial role in informing the public (Gentzkow, 2006). Its primary purpose in terms of democracy is to provide accurate information that allows citizens to make educated decisions. Acting as a watchdog for democracy, journalism holds those in power accountable by reporting about issues such as corruption, shortcomings of policies, and other scandals (Allan, 2022). However, trust in media has in recent years gone down especially in the United States (Edelman, 2018). There is growing fear among academics that AI could exacerbate this phenomenon. In addition to direct harms from misinformation, the proliferation of AI-generated inauthentic content, such as deep fakes, AI-generated images and text, could result in people not knowing what to believe in (Kreps & Kriner, 2023). This could further lead to lower trust in established institutions.

Some reputable news sites have already publicly made some experimental articles written by LLM (Sweeney, 2023). In addition to verified uses of LLMs, we note that it is also possible that there is non-disclosed LLM usage at the organizational level and shadow usage by individual journalists.

There are several existing efforts demonstrating the increase of LLM use in academic articles and peer reviews (Liang et al., 2024b; Shapira, 2024; Stokel, 2024). In this paper, we examine whether news articles have seen a surge in AI-generated content in recent years, using similar techniques to (Shapira, 2024).

## 2. Methodology

To investigate the usage of LLMs in news articles and public comments, we adopted a similar approach to Shapira (2024) and Liang et al. (2024b), based on word frequency comparisons between different years. The words used were eight of the ten words overrepresented in the outputs of ChatGPT, the most popular LLM, mentioned in the blog post by AI Phrase Finder (2024). The words were: "captivate", "tapestry", "leverage", "embrace", "resonate", "testament", "delve", and "elevate". Additionally, we used eight control words to showcase how much the occurrence of the words that are not overrepresented in the outputs of ChatGPT have changed over the years.

We curated a dataset consisting of news articles from various reputable sources. The dataset comprises 75,000 articles published in the years 2019-2023, and it was combined from four different datasets (Saksham, 2023;

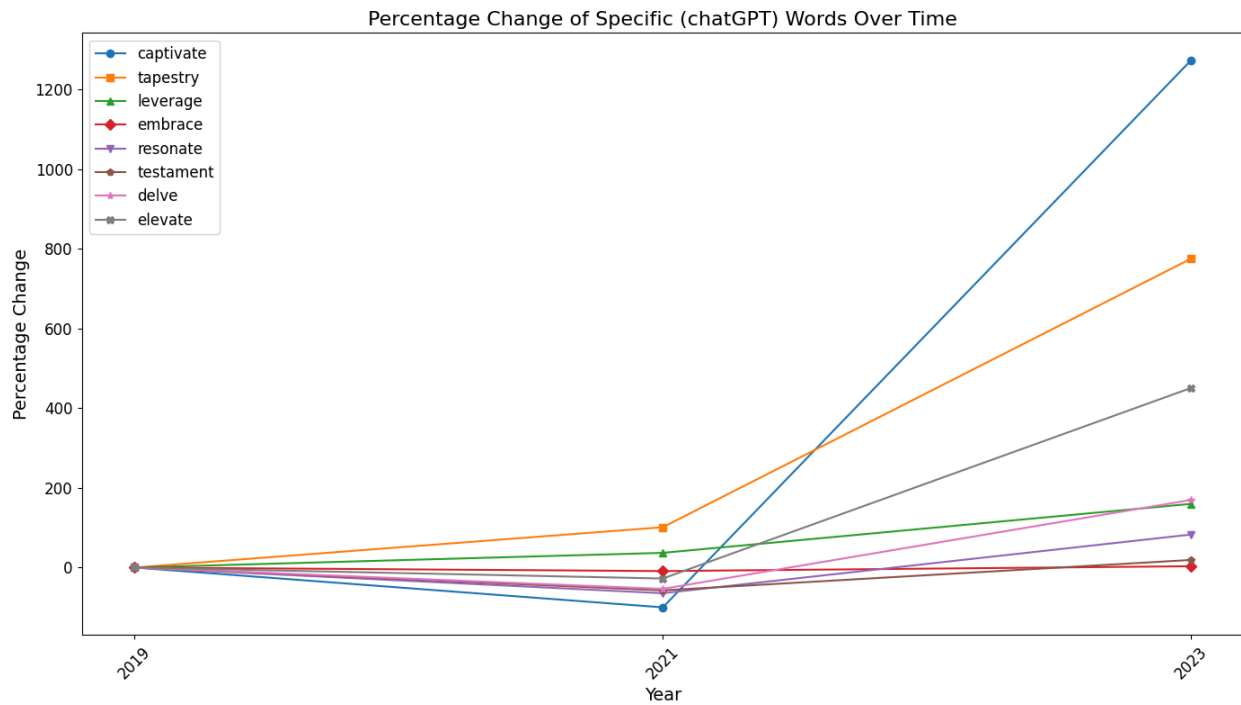
Petukhova & Fachada, 2022; Hugging Face, 2023; Hugging Face, 2021). The datasets contain news articles from several popular news sources, such as ABC News, Al-Jazeera, and Forbes.

After the data was collected and preprocessed, we performed a word frequency analysis to determine the evolution of the occurrence of the eight target words in the datasets over several years. The relative frequency of each word was compared to the baseline frequency of 2019. The code for this section can be found on [Github](#).

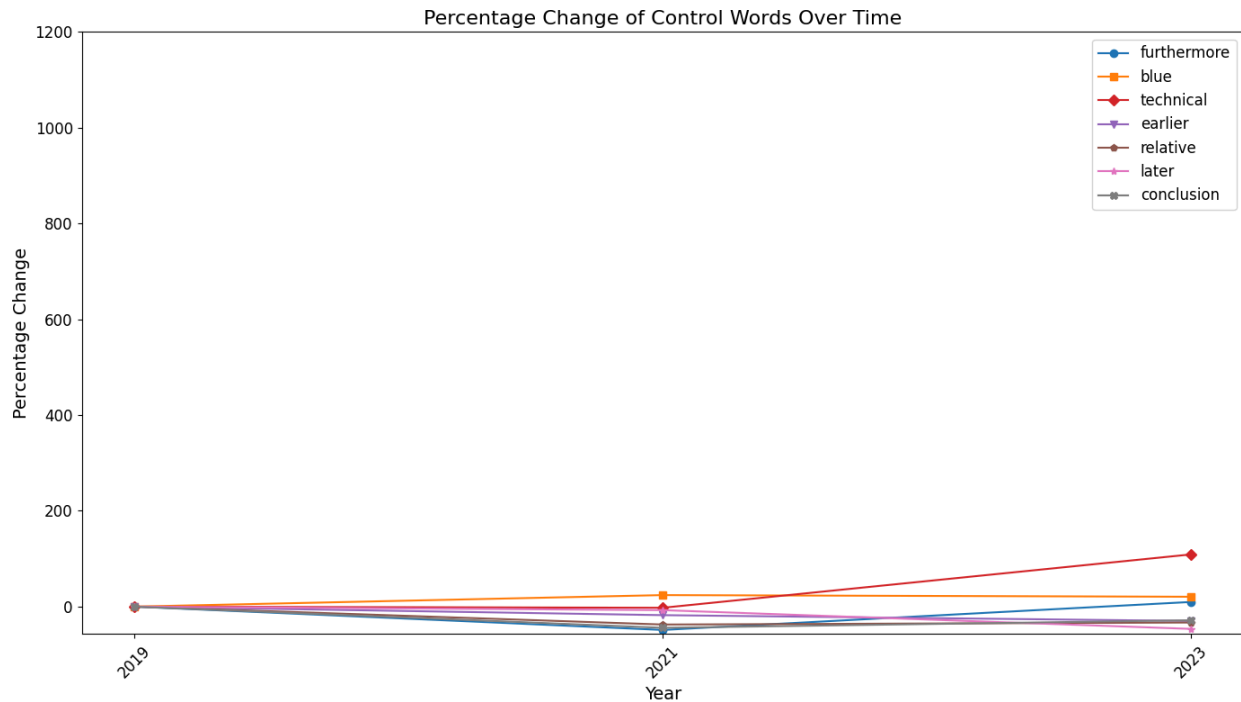
### 3. Results

Our results show a significant increase in the frequency in many of the words favored by ChatGPT in 2023 compared to the frequencies in 2019. That time window coincides with the release of ChatGPT in 2022. Figure 3 shows that three of the words, “captivate”, “tapestry” and “elevate” saw around 1200%, 800% and 400% increase, respectively. Two other words had almost a 200% increase, and the other words had a smaller increase. Meanwhile the changes are much smaller in the control words. As seen in Figure 4, Only the word “technical” had a significant increase of around 100%, while the others had very little changes. Overall the frequency of words favoured by ChatGPT increased significantly compared to the change in control words. The findings suggest that LLMs might have been used in the creation of news articles.

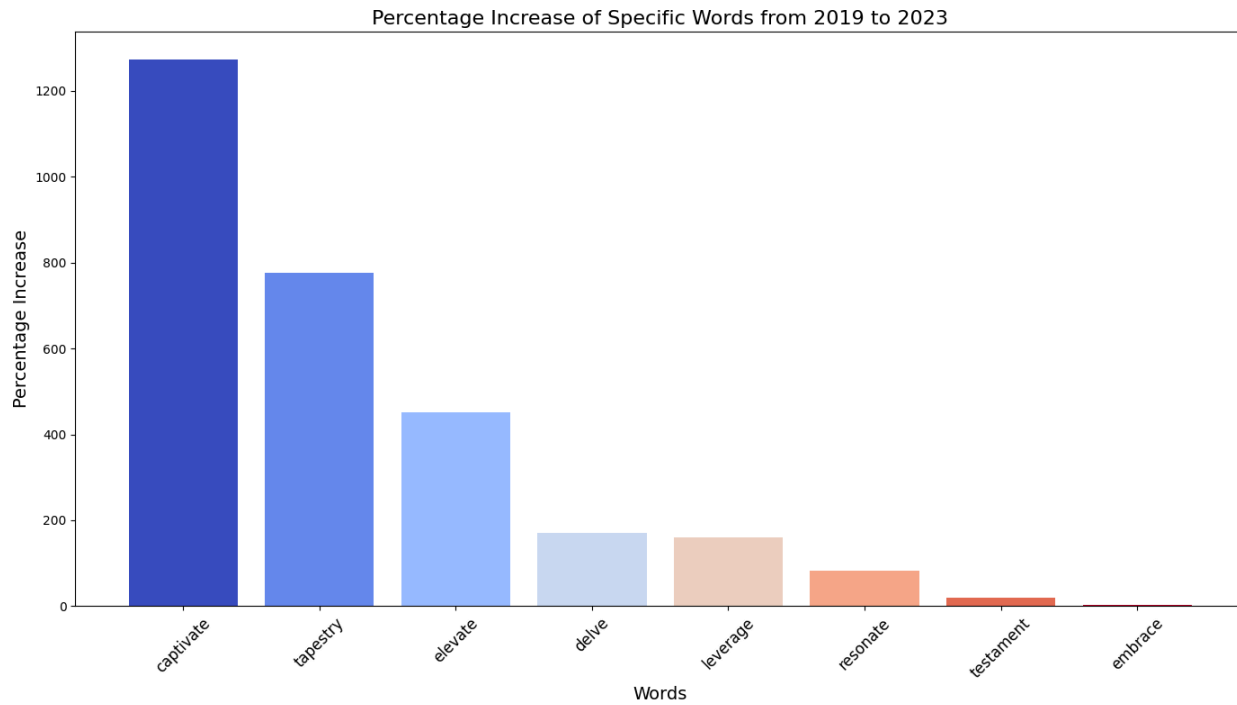
**Figure 1. Change of the frequency of words favored by ChatGPT over time**



**Figure 2. Change of the frequency of the control words over time**

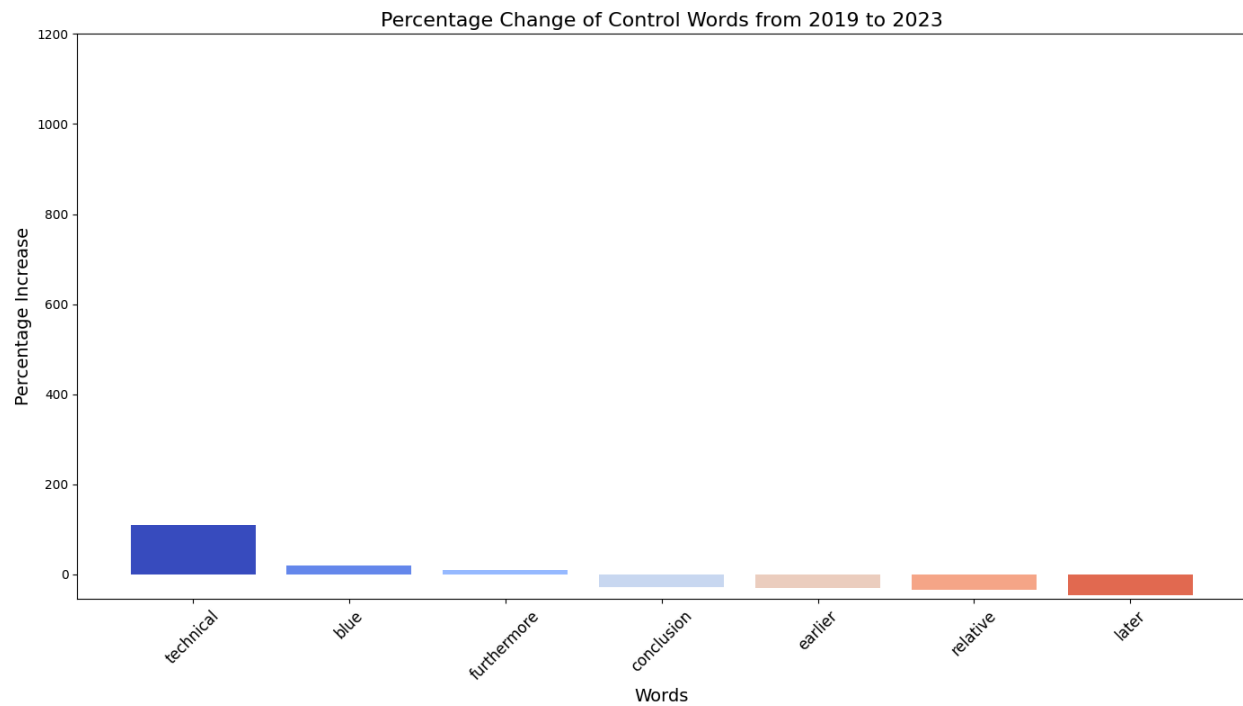


**Figure 3. Overall change of the frequency of the words favored by ChatGPT over time**



**Figure 4.**

**Overall change of the frequency of the control words over time**



## 4. Discussion

Our findings provide tentative evidence of LLMs potentially being used to produce news articles. However, there are several limitations to our approach. First, the words selected to measure the prevalence of LLM-generated text were selected on a relatively ad-hoc basis, based on a blog post (AI Phrase Finder, 2024). Additionally, the approach does not give a quantitative measure of how much LLM-generated text has increased; it merely provides evidence that *some* increase has happened.

A more meticulous analysis would give a quantitative measure of the fraction of the text that is likely to be produced or modified by an LLM, as in (Liang et al., 2024a). Additionally, further research could involve a more methodological approach to choosing the words associated with AI.

Our approach combined different datasets which can cause unnecessary variance in the frequency of different words. There could also be more analysis to show that the changes in the words are statistically significant.

While the increase in AI-generated content in media could undermine citizens' trust in institutions, there is a possibility that AI could also undermine trust in the other direction. In particular, AI could be used to misrepresent public opinions in order to advance a particular viewpoint (Kreps & Kriner, 2023). This could undermine political leaders' trust in receiving reliable information from the public. In practice, LLMs could for example be used to automate the composition of public comments for new regulation. Therefore, an especially promising direction would be to analyze these public comments, openly available on websites such as [regulations.gov](https://www.regulations.gov), to detect whether there is evidence of LLM usage increasing. Furthermore, there are many other types of text that could be explored. Further research could focus on other datasets that have relevance to democracy, such as tweets by politicians and ordinary citizens, or political communication of different parties.

In addition to directly advancing knowledge of LLM use in different text types, we encourage further research on topics relevant to risks of AI undermining democratic trust. For example, LLM watermarking could help to identify whether a particular article is written by an LLM. Recent research suggests that there is still much work to be done before such technology is ready to be deployed (Jovanovic et al. 2024).

Finally, we encourage policymakers to invent ways of ensuring that LLM watermarking, once technically feasible, is used in the most efficient way possible to ensure that citizens' trust in institutions is protected. Examples of such efforts would be mandating news agencies to tell their users about the usage of LLMs, and giving users themselves easy-to-use tools for detecting whether a given source is written by AI.

## 5. Acknowledgements

The paper was written as a submission to the AI x Democracy Hackathon organized by Apart Research. We would like to thank Apart Research for organizing the hackathon, Aalto Effective Altruism for hosting a local version of it, and Aalto Design Factory for offering a comfortable co-working space and a foosball table. Finally, we would like to thank ChatGPT and Claude for assistance in ideating, coding, and writing.

## 6. References

- AI Phrase Finder. (2024, March 9). *The 10 Most Common ChatGPT Words*. AI Phrase Finder. Retrieved May 5, 2024, from <https://aiphrasefinder.com/common-chatgpt-words/>
- Allan, S. (Ed.). (2022). *The Routledge Companion to News and Journalism*. Routledge.
- Edelman. (2018, January 22). *2018 Edelman Trust Barometer Reveals Record-Breaking Drop in Trust in the U.S.* Edelman. Retrieved May 5, 2024, from <https://www.edelman.com/news-awards/2018-edelman-trust-barometer-reveals-record-breaking-drop-trust-in-the-us>
- Gentzkow, M., Glaeser, E. L., & Goldin, C. (2006, 03). *The Rise of the Fourth Estate. How Newspapers Became Informative and Why It Mattered*.
- Hugging Face. (2021). *News\_Seq\_2021 Dataset*. Hugging Face. Retrieved May 6, 2024, from [https://huggingface.co/datasets/RealTimeData/News\\_Seq\\_2021](https://huggingface.co/datasets/RealTimeData/News_Seq_2021)
- Hugging Face. (2023). *News\_August\_2023 Dataset*. Hugging Face. Retrieved May 6, 2024, from [https://huggingface.co/datasets/RealTimeData/News\\_August\\_2023?row=0](https://huggingface.co/datasets/RealTimeData/News_August_2023?row=0)
- Jovanovic, N., Staab, R., & Vechev, M. (2024). *Watermark Stealing in Large Language Models*. <https://watermark-stealing.org/>
- Kreps, S., & Kriner, D. (2023). *How AI Threatens Democracy*. *Journal of Democracy*, 34(4), 122-31. <https://www.journalofdemocracy.org/articles/how-ai-threatens-democracy/>



- Liang, W., & et al. (2024, March 11). Monitoring AI-Modified Content at Scale: A Case Study on the Impact of ChatGPT on AI Conference Peer Reviews. *arXiv*. <https://arxiv.org/abs/2403.07183>
- Liang, W., & et al. (2024, April 1). Mapping the Increasing Use of LLMs in Scientific Papers. *arXiv*. <https://arxiv.org/abs/2404.01268>
- Petukhova, A., & Fachada, N. (2022, December 3). *MN-DS: A Multilabeled News Dataset for News Articles Hierarchical Classification*. Zenodo. Retrieved May 6, 2024, from <https://zenodo.org/records/7394851>
- Saksham, K. (2023). *Global News Dataset*. Kaggle. Retrieved May 6, 2024, from <https://www.kaggle.com/datasets/everydaycodings/global-news-dataset>
- Shapira, P. (2024, March 31). *Delving into “delve”*. Philip Shapira. Retrieved May 5, 2024, from <https://pshapira.net/2024/03/31/delving-into-delve/>
- Stokel, C. (2024, May 1). *AI Chatbots Have Thoroughly Infiltrated Scientific Publishing*. Scientific American. Retrieved May 5, 2024, from <https://www.scientificamerican.com/article/chatbots-have-thoroughly-infiltrated-scientific-publishing/>
- Sweney, M. (2023, March 7). Mirror and Express owner publishes first articles written using AI. *The Guardian*. <https://www.theguardian.com/business/2023/mar/07/mirror-and-express-owner-publishes-first-articles-written-using-ai>