# Is ChatGPT Transforming Academics' Writing Style?

**Mingmeng Geng** [1]   **Roberto Trotta** [1 2]

## Abstract

Based on one million arXiv papers submitted from May 2018 to January 2024, we assess the textual density of ChatGPT's writing style in their abstracts by means of a statistical analysis of word frequency changes. Our model is calibrated and validated on a mixture of real abstracts and ChatGPT-modified abstracts (simulated data) after a careful noise analysis. We find that ChatGPT is having an increasing impact on arXiv abstracts, especially in the field of computer science, where the fraction of ChatGPT-revised abstracts is estimated to be approximately 35%, if we take the output of one of the simplest prompts, "revise the following sentences", as a baseline. We conclude with an analysis of both positive and negative aspects of the penetration of ChatGPT into academics' writing style.

## 1. Introduction

Since its official release on November 30, 2022, ChatGPT (Chat Generative Pre-trained Transformer) has impacted many aspects of our lives, and academic writing has not been immune. While ChatGPT does increase productivity and may help scientific discovery (Noy & Zhang, 2023; AI4Science & Quantum, 2023), we must be wary of its potential risks and the possibility of negative impacts. A large number of papers have already explored the advantages and disadvantages of large language models (LLMs) (Kasneci et al., 2023); here, we focus on ChatGPT, which is being very widely used (reaching an unprecedented 100 million active users three months after its release) and is considered to be one of the two major milestones of language models along with GPT-4. (Zhao et al., 2023; von Garrel & Mayer, 2023)

While there is already a corpus of current research on using ChatGPT in academia (Casal & Kessler, 2023; Lingard et al., 2023; Fergus et al., 2023; Lund et al., 2023), to our knowledge only a handful of works have attempted to quantify its impact on the whole academic community. As this article was being finalized, two preprints appeared that addressed related questions: one focuses on AI conferences peer reviews (Liang et al., 2024a) and, even more recently, analyzes scientific papers(Liang et al., 2024b). They claim that the usage of LLMs is evident in AI conference reviews and scientific writings, especially in computer science papers. Within the broad field of academic writing and publishing, we chose the abstracts of articles as the focus of this work, as they have a relatively uniform format across disciplines, are supposed to condense an entire research article and thus are often highly polished, and can be considered short articles of pure text, not involving pictures nor tables.

ChatGPT is of course able to generate abstracts directly given a suitable prompt (Luo et al., 2023), and studies have shown that identifying such abstracts is not easy even if they remain unedited by humans (Gao et al., 2023; Cheng et al., 2023) – watermarking being a possible strategy to enable such identification (Kirchenbauer et al., 2023). Determining whether a given few sentences were generated by ChatGPT is difficult, but determining that millions of sentences were influenced by ChatGPT is statistically feasible, as we demonstrate here. We analyzed the fingerprints of ChatGPT on scientific abstracts as a function of time in order to tease out a statistical signature, rather than attempting to detect whether a given abstract was generated or polished with ChatGPT.

In fact, that the abstract of a paper shows what we call the "ChatGPT style" does not necessarily mean that the authors directly utilized ChatGPT to generate or modify it. It is also possible that the authors used ChatGPT in another context and that, as a result, their writing habits were influenced by the ChatGPT style – not a remote possibility. It is worth considering in this context that reading and writing in English is more difficult for non-native English academics (Amano et al., 2023). Before ChatGPT was released, the pros and cons of other tools were discussed, such as Google Translate (Mundt & Groves, 2016) and Grammarly (Fitria, 2021), but ChatGPT has a much wider range of application scenarios – not to mention, a much higher flexibility.

[1]Theoretical and Scientific Data Science Group, Scuola Internazionale Superiore di Studi Avanzati (SISSA), Trieste, Italy. [2]Department of Physics, Imperial College London, London, UK. Correspondence to: Mingmeng Geng <mgeng@sissa.it>, Roberto Trotta <rtrotta@sissa.it>.

We have seen similar AI-induced seismic shifts in the past: after AlphaGo (Silver et al., 2017) shocked the world, professional Go players have begun training with AI, and the sport of Go has been profoundly changed as a result (Kang et al., 2022). AlphaFold brings new opportunities for life science research (Varadi & Velankar, 2023) and ChatGPT has also been used for data extraction in materials science. (Polak & Morgan, 2023). A similar story may be happening with academic writing, especially for researchers whose first language is not English (Hwang et al., 2023). This paper is a first effort at establishing whether this is the case.

## 2. Data

**arXiv dataset:** The metadata of arXiv papers are provided by Kaggle (arXiv.org submitters, 2024). Because the abstracts in this dataset are updated when authors submit changes, we used the first version in 2024 (version 161) as well as the last version before the ChatGPT era (version 105). Our observations and analysis are based on one million arXiv articles submitted from May 2018 to January 2024.

**English word frequency:** Google Ngram dataset is chosen for comparison and reference (Michel et al., 2011). Specifically, we used the freely available mirrors on Kaggle (`http://kaggle.com/datasets/wheelercode/english-word-frequency-list`) covering word frequencies from the 1800s to 2019 as established from Google Books.

## 3. Observations and analysis

### 3.1. Changes in word frequency

We approach the problem by analyzing how the frequency of words changes after ChatGPT has been deployed, so we define the change factor in the frequency of word $i$, $R_i$, as follows:

$$R_i = \frac{\max_t(f_i(t)) - \min_t(f_i(t))}{\max_t(f_i(t))} \quad (1)$$

where $f_i(t)$ is the count of word $i$ during the time period $t$. We divided the 1 million abstracts into 100 uneven time-periods, each encompassing 10,000 abstracts.

Figure 1 illustrates that most of the words with the largest change rate in the time period considered (generally, an increase) in the abstracts are related to hot research topics of the last few years, such as "Covid-19", "LLMs", "AI". But this is not the case for all words with the largest growth in terms of frequency. Indeed, the frequency of some non-specialized words also starts to skyrocket in early 2023, as presented in Figure 2. How could the frequencies of words like "**significant**" grow **significantly** together?
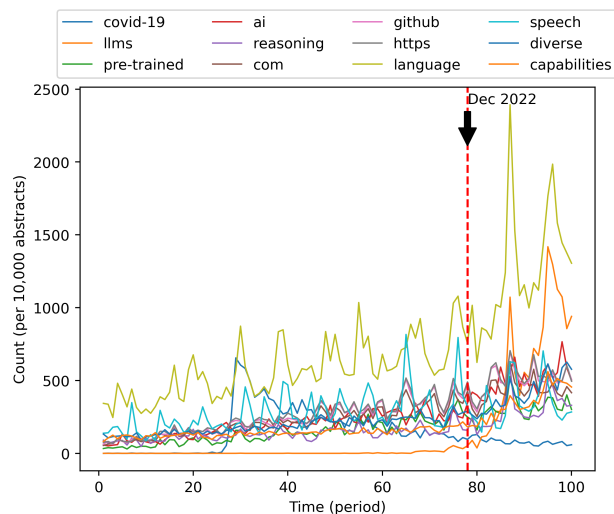


*Figure 1.* The 12 words with the highest change rate $R_i$ and satisfying $\max_t(f_i(t)) > 500$. The vertical red dashed line demarcates the first time period after ChatGPT's release.

Another striking example is the frequency change of the words "are" and "is", as depicted in Figure 3. The counts in 10,000 abstracts of these two words were quite stable before 2023. However, the frequency of these two terms has dropped by more than 10% in 2023.

Because the average length of abstracts tends to grow over time, we also considered normalizing frequencies to abstract length. The corresponding figures are displayed in the Appendix, and show similar trends.

These examples, anecdotal as they are, may represent the tip of the iceberg of a wider and growing phenomenon: the rapid increase in the usage of ChatGPT. The rise and fall in frequency of specific technical nouns may well be related to the changing popularity of certain research topics, but that a research trend is responsible for the change in usage of adjectives appears implausible – even less so for words like "is" and "are".

### 3.2. ChatGPT simulations

We wanted to be more specific about the impact of ChatGPT on articles from different disciplines, so we examined arXiv abstracts from different categories separately. The one million arXiv articles were divided into 20 periods in this part in order to increase the number of articles per period and reduce estimation error, which is not the same as the previous part. The identifier numbers of the first and last arXiv articles corresponding to each period are given in the Appendix.

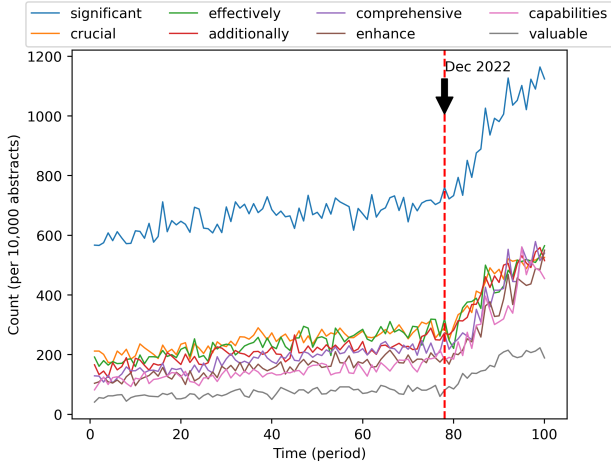Previous studies have shown that ChatGPT has its own

*Figure 2.* Examples of words with rapidly growing frequency in arXiv abstracts.



*Figure 3.* The words "are" and "is" are decreasing in frequency in arXiv abstracts.

linguistic style ([AlAfnan & MohdZuki, 2023](#)), and that likely includes the frequency of some words. Although there is no direct way to investigate ChatGPT's word preference, we can ask ChatGPT to polish or rewrite real, pre-2023 abstracts, and use the resulting simulation data to calculate the estimated frequency change rate $\hat{r}_{ij}$ of word $i$ in category $j$:

$$\hat{r}_{ij} = \frac{\tilde{q}_{ij}^d - q_{ij}^d}{q_{ij}^d} = \frac{\tilde{q}_{ij}^d}{q_{ij}^d} - 1 \tag{2}$$

where $q_{ij}^d$ represents the word frequency of real abstracts in the dataset and $\tilde{q}_{ij}^d$ means the frequency after ChatGPT processing. We have no way of knowing the real usage scenarios of ChatGPT, so some simple prompts were used, for example,

*"Revise the following sentences:"*

GPT-3.5 was utilized in our simulations for 10,000 abstracts in period 14 (April 2022 to July 2022), although it may have different word preferences than the more recent GPT-4. Many words have different frequencies before and after ChatGPT processing, such as the words "is", "are", and "significant" that we mentioned earlier. For simplicity, the results of the 4 categories with the highest number of articles are shown in Table 1 and the rest parts in this paper, namely *cs* (computer science), *math* (mathematics), *astro* (astrophysics), and *cond-mat* (condensed matter). This corroborates the hypothesis, formulated earlier, that the drop in the frequency of these two words observed in real abstracts in 2023 may have been caused by ChatGPT.

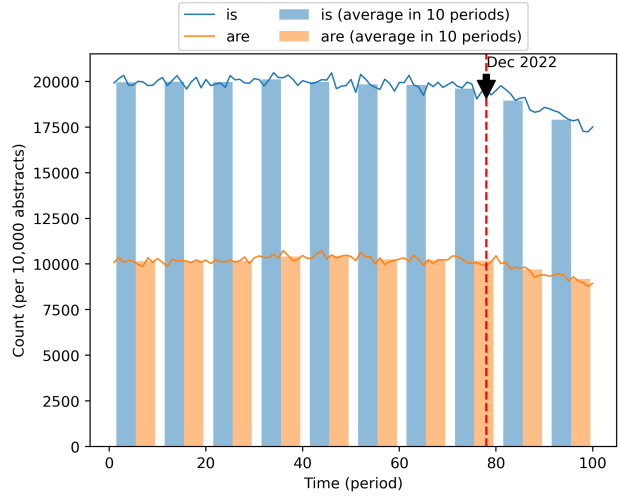In the meantime, we also defined the word frequency change

*Table 1.* Word frequency (per abstract) before and after ChatGPT processing in simulation data of period 14.

| WORD | CATEGORY | BEFORE | AFTER | CHANGE |
|---|---|---|---|---|
| IS | CS | 2.01 | 1.73 | -14% |
| IS | MATH | 1.78 | 1.61 | -9% |
| IS | ASTRO | 2.13 | 1.90 | -11% |
| IS | COND-MAT | 2.00 | 1.68 | -16% |
| ARE | CS | 1.00 | 0.83 | -17% |
| ARE | MATH | 0.74 | 0.71 | -5% |
| ARE | ASTRO | 1.39 | 1.25 | -1% |
| ARE | COND-MAT | 0.92 | 0.80 | -13% |
| SIGNIFICANT | CS | 0.09 | 0.18 | 99% |
| SIGNIFICANT | MATH | 0.01 | 0.03 | 308% |
| SIGNIFICANT | ASTRO | 0.17 | 0.26 | 53% |
| SIGNIFICANT | COND-MAT | 0.07 | 0.18 | 171% |

in all abstracts from year $t-1$ to year $t$, $R_{ij,t}$:

$$R_{ij,t} = \frac{F_{ij,t} - F_{ij,t-1}}{F_{ij,t-1}}, \tag{3}$$

where $F_{ij,t}$ represent frequency of word $i$ per arXiv abstract in category $j$ in year $t$.

Only words with a frequency larger than 0.1 times per abstract before ChatGPT processing are plotted in Figure 4 and Figure 5. The correlation coefficient between the word frequency change in arXiv abstracts and our estimated ChatGPT-induced word frequency change is very small in all four categories of abstracts, as shown in Figure 4.

However, Figure 5 presents a totally different pattern, where $\hat{r}_{ij}$ and $R_{ij,2023}$ are strongly correlated, especially in com-
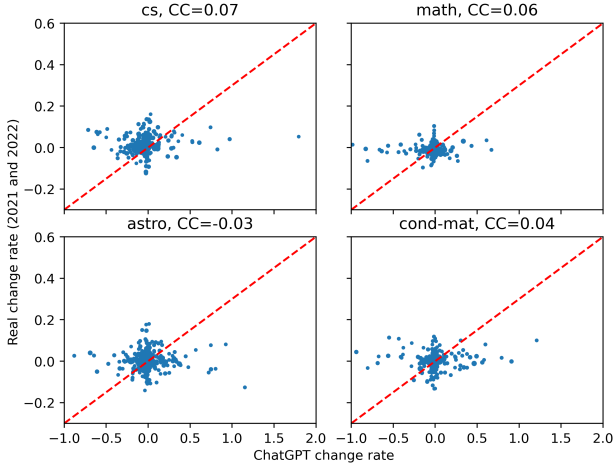
*Figure 4.* Comparison of the predicted frequency change rate due to ChatGPT $\hat{r}_{ij}$ (x-axis) and the actual word frequency change for all abstracts from 2021 to 2022 $R_{ij,2022}$ (y-axis). CC indicates the correlation coefficient.
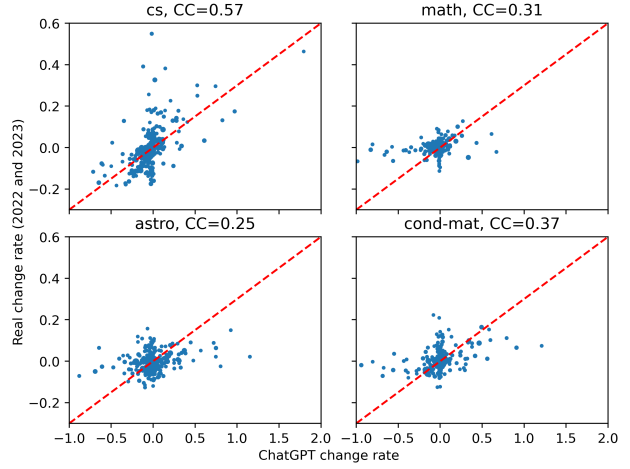


*Figure 5.* Comparison of the predicted frequency change rate due to ChatGPT $\hat{r}_{ij}$ (x-axis) and the actual word frequency change for all abstracts from 2022 to 2023 $R_{ij,2023}$ (y-axis). The correlation coefficients (CC) are now significantly positive.

puter science abstracts. Although many words seem insensitive to ChatGPT, we can still see a positive correlation for some words in this figure, even among the other categories.

Taken together, our consideration point to ChatGPT as one of the important reasons, possibly even the main reason, for the recent word frequency change in abstracts. Our next step is to start by modeling ChatGPT impact, as well as estimating the impact based on real data and simulations.

## 4. Models and methods

### 4.1. ChatGPT impact

Imagine different scenarios of using ChatGPT in scientific writing: a researcher might simply use it to correct grammatical errors, another employs it for translating native sentences into English, and yet another one wants it to polish their draft in English very purposefully. In theory, each of these use cases contributes the same proportion of ChatGPT usage. But, as is well known, different prompts will lead to different outputs, which means different word frequency changes. Therefore, we use the more neutral term "ChatGPT impact" instead of "proportion" in our estimation part.

We start with a simple model, ignoring noise and variability for this subsection. Suppose that the frequency of word $i$ for abstracts in subject category $j$ changes from $f_{ij}^*$ to $\tilde{f}_{ij}^*$ after being processed by ChatGPT, when it's used as a means to polish and improve the abstract (if not to fully generate it).

The corresponding word change rate is defined as

$$\bar{r}_{ij} = \frac{\tilde{f}_{ij}^* - f_{ij}^*}{f_{ij}^*} = \frac{\tilde{f}_{ij}^*}{f_{ij}^*} - 1\,. \tag{4}$$

Suppose that $\bar{f}_{ij}(t)$ is the word frequency for word $i$ in category $j$ at time period $t$, this can be written as:

$$\begin{aligned}
\bar{f}_{ij}(t) &= (1 - \eta_j(t))f_{ij}^*(t) + \eta_j(t)f_{ij}^*(t)(\bar{r}_{ij} + 1) \\
&= f_{ij}^*(t) + \eta_j(t)f_{ij}^*(t)\bar{r}_{ij}
\end{aligned} \tag{5}$$

where $\eta_j(t)$ denotes the proportion of abstracts in category $j$ affected by ChatGPT, and $f_{ij}^*(t)$ represents the original evolution in word frequency without ChatGPT.

This model is highly idealized: we have to additionally consider the effects of noise (such as randomness inside ChatGPT), uncertainty in word usage evolution without ChatGPT, and the epistemic uncertainty in how users actually prompt ChatGPT.

### 4.2. Noise model

We now consider the noise terms, which might be modelled in many different ways.

For example, we denote the word frequency for word $i$ in category $j$ by $f_{ij}^d$, which represents the word frequency observed in the data:

$$f_{ij}^d = f_{ij}^* + \delta_{ij}(f_{ij}^*) \tag{6}$$

where $\delta_{ij}(\cdot)$ represents noise and word usage variability which are not directly related to the internal parameters of ChatGPT.

4

Thus, we can define

$$f_{ij}^{\delta,\eta}(t) = \eta_j(t)f_{ij}^*(t) + \delta_{ij}(\eta_j(t)f_{ij}^*(t)),\qquad(7)$$

and

$$f_{ij}^{\delta,1-\eta}(t) = (1 - \eta_j(t))f_{ij}^*(t)\\ + \delta_{ij}((1 - \eta_j(t))f_{ij}^*(t)).\qquad(8)$$

In this case, the equation corresponding to Eq. (5) is

$$f_{ij}^d(t) = f_{ij}^{\delta,1-\eta}(t) + \mathrm{C}_{ij}(f_{ij}^{\delta,\eta}(t))\qquad(9)$$

where the function $\mathrm{C}_{ij}(\cdot)$ means the frequency after Chat-GPT process.

We assume that the noise for word $i$ due to ChatGPT processing can be represented as $\epsilon_{ij}(\cdot)$ and $\epsilon_{ij}^s(\cdot)$, then Eq. (2) and Eq. (4) are related by

$$\frac{\tilde{f}_{ij}^* - \epsilon_{ij}(f_{ij}^*) - f_{ij}^*}{f_{ij}^*} = \frac{\tilde{q}_{ij}^d - \epsilon_{ij}^s(q_{ij}^d) - q_{ij}^d}{q_{ij}^d}.\qquad(10)$$

Therefore,

$$\mathrm{C}_{ij}(f_{ij}^{\delta,\eta}(t)) = f_{ij}^{\delta,\eta}(t)(\hat{r}_{ij} + 1 + \epsilon_{ij}^\eta(q,f,t))\qquad(11)$$

where

$$\epsilon_{ij}^\eta(q,f,t) = \frac{\epsilon_{ij}(f_{ij}^{\delta,\eta}(t))}{f_{ij}^{\delta,\eta}(t)} - \frac{\epsilon_{ij}^s(q_{ij}^d)}{q_{ij}^d}.\qquad(12)$$

Then, Eq. (9) – representing the difference in word frequency before and after ChatGPT processing – can be rewritten as

$$f_{ij}^d(t) - f_{ij}^*(t) = \eta_j(t)x_{ij}(t) + g_{ij}(t) + \xi_{ij}(t)\qquad(13)$$

where

$$x_{ij}(t) = f_{ij}^*(t)\hat{r}_{ij}\qquad(14)$$
$$g_{ij}(t) = \eta_j(t)f_{ij}^*(t)\epsilon_{ij}^\eta(q,f,t)\qquad(15)$$
$$\xi_{ij}(t) = (\hat{r}_{ij} + 1 + \epsilon_{ij}^\eta(q,f,t))\delta_{ij}(\eta_j(t)f_{ij}^*(t))\\ + \delta_{ij}'((1 - \eta_j(t))f_{ij}^*(t)).\qquad(16)$$

where $\delta_{ij}'(\cdot)$ follows the same distribution as $\delta_{ij}(\cdot)$. It should be noted that $g_{ij}(t)$ includes only ChatGPT-related noise $\epsilon_{ij}(\cdot)$ and $\epsilon_{ij}^s(\cdot)$, however $\xi_{ij}(t)$ contains $\delta_{ij}(\cdot)$ and $\delta_{ij}'(\cdot)$ that are unrelated to ChatGPT.

## 4.3. Impact estimation and bias analysis

In many data analysis applications, more data point (in our case, using a larger number of words) means better estimates. But in our case, the effect of noise is different for each data point (word), and choosing wisely which words to include can improve our estimates.

For simplicity, we define

$$h_{ij}(t) = f_{ij}^d(t) - f_{ij}^*(t).\qquad(17)$$

For abstracts in category $j$, we use the words in the subset $I_j$ (whose determination is discussed below), of numerosity $n_j$. In order to estimate $\eta_j(t)$, we can use the quadratic loss function

$$L_{j,t}(\eta_j) = \frac{1}{n_j}\sum_{i\in I_j}(h_{ij}(t) - \eta_j(t)x_{ij}(t))^2\\ = \frac{1}{n_j}\sum_{i\in I_j}(g_{ij}(t) + \xi_{ij}(t))^2.\qquad(18)$$

If we ignored the dependency of $g_{ij}(t)$ and $\xi_{ij}(t)$ on $\eta_j(t)$, the estimate of ChatGPT impact would simply be given by Ordinary Least Squares (OLS) as

$$\hat{\eta}_j(t) = \frac{\sum_{i\in I_j} h_{ij}(t)x_{ij}(t)}{\sum_{i\in I_j} x_{ij}^2(t)}.\qquad(19)$$

However, since $g_{ij}(t)$ also depends on $\eta_j(t)$ and $\xi_{ij}$ contains $\eta_j(t)$ as described in Eq. (15) and Eq. (16), we need to make additional assumptions to progress further.

**Case 1:** if the effect of $\eta_j(t)$ on $\xi_{ij}(t)$ can be ignored compared to other terms, e.g., the following simple scenario,

$$\mathrm{Var}[\delta_{ij}(\eta_j(t)f_{ij}^*(t))] \ll \eta_j(t)f_{ij}^*(t)\mathrm{Var}[\epsilon_{ij}^\eta(q,f,t)]\quad(20)$$

One can also derive the approximation below:

$$f_{ij}^{\delta,\eta}(t) \approx \eta_j(t)f_{ij}^*(t) + \delta_{ij}(*)\qquad(21)$$

where $\delta_{ij}(*)$ is a random variable with zero mean and variance much smaller than $\eta_j(t)f_{ij}^*(t)$, and its derivative with respect to $\eta_j(t)$ is negligible compared to $f_{ij}^*(t)$.

Therefore, the loss function under this assumption is:

$$L_{j,t,g}(\eta_j)\\ = \frac{1}{n_j}\sum_{i\in I_j}(h_{ij}(t) - \eta_j(t)x_{ij}(t) - g_{ij}(t))^2\\ = \frac{1}{n_j}\sum_{i\in I_j}\xi_{ij}^2(t).\qquad(22)$$

Thus,

$$\frac{\partial L_{j,t,g}(\eta_j)}{\partial\eta_j}\\ = \frac{2}{n_j}\sum_{i\in I_j}\left(\eta_j(t)x_{ij}^2(t) - h_{ij}(t)x_{ij}(t)\right)\\ - \frac{2}{n_j}\sum_{i\in I_j}\frac{\partial g_{ij}(t)}{\partial\eta_j(t)}(h_{ij}(t) - \eta_j(t)x_{ij}(t) - g_{ij}(t))\\ + \frac{2}{n_j}\sum_{i\in I_j}x_{ij}(t)g_{ij}(t)$$

If we require a minimum by setting $\frac{\partial L_{j,t,g}(\eta_j)}{\partial \eta_j} = 0$, we obtain a new estimate $\hat{\eta}_j^g(t)$, which is equal to the OLS $\hat{\eta}_j(t)$ in Eq. (19) corrected for bias and noise,

$$
\begin{aligned}
(\hat{\eta}_j^g(t) &- \hat{\eta}_j(t)) \sum_{i \in I_j} x_{ij}^2(t) \\
&= \sum_{i \in I_j} \frac{\partial g_{ij}(t)}{\partial \eta_j(t)} (h_{ij}(t) - \eta_j(t) x_{ij}(t)) \\
&- \sum_{i \in I_j} x_{ij}(t) g_{ij}(t) - \sum_{i \in I_j} g_{ij}(t) \frac{\partial g_{ij}(t)}{\partial \eta_j(t)} .
\end{aligned}
\tag{23}
$$

But without knowing the distribution of $\epsilon_{ij}(\cdot)$ and $\epsilon_{ij}^s(\cdot)$, we have no way of estimating the value of this bias, so we assume that $\epsilon_{ij}(f_{ij}) \sim \mathcal{N}(0, f_{ij}\sigma_{ij,\epsilon}^2)$ and $\epsilon_{ij}^s(f_{ij}) \sim \mathcal{N}(0, f_{ij}\sigma_{ij,\epsilon}^2)$, e.g., $\epsilon_{ij}(1) \sim \mathcal{N}(0, \sigma_{ij,\epsilon}^2)$, then we can obtain an expression for $\epsilon_{ij}^\eta(q, f, t)$:

$$
\epsilon_{ij}^\eta(q, f, t) = \frac{\epsilon_{ij}(1)}{\sqrt{\eta_j(t) f_{ij}^*(t) + \delta_{ij}(*)}} - \frac{\epsilon_{ij}^s(1)}{\sqrt{q_{ij}^d}}
\tag{24}
$$

$$
g_{ij}(t) = \frac{\eta_j(t) f_{ij}^*(t) \epsilon_{ij}(1)}{\sqrt{\eta_j(t) f_{ij}^*(t) + \delta_{ij}(*)}}
\tag{25}
$$

$$
- \frac{\eta_j(t) f_{ij}^*(t) \epsilon_{ij}^s(1)}{\sqrt{q_{ij}^d}} .
\tag{26}
$$

Therefore, all terms on the right-hand side of Eq. (23) are zero-mean noise, except for the last one:

$$
\begin{aligned}
g_{ij}(t) &\frac{\partial g_{ij}(t)}{\partial \eta_j(t)} \\
= g_{ij}(t) &\frac{f_{ij}^*(t)(\eta_j(t) f_{ij}^*(t) + 2\delta_{ij}(*)) \epsilon_{ij}(1)}{2(\eta_j(t) f_{ij}^*(t) + \delta_{ij}(*))^{\frac{3}{2}}} \\
- g_{ij}(t) &\frac{f_{ij}^* \epsilon_{ij}^s(1)}{\sqrt{q_{ij}^d}} .
\end{aligned}
\tag{27}
$$

Removing the items with zero means, we get

$$
\begin{aligned}
\mathrm{E}\left[ g_{ij}(t) \frac{\partial g_{ij}(t)}{\partial \eta_j(t)} \right] & \\
= \frac{\eta_j(t)(f_{ij}^*(t))^2 (\eta_j(t) f_{ij}^*(t) + 2\delta_{ij}(*)) \sigma_{ij,\epsilon}^2}{2(\eta_j(t) f_{ij}^*(t) + \delta_{ij}(*))^2} & \\
+ \frac{\eta_j(t)(f_{ij}^*(t))^2 \sigma_{ij,\epsilon}^2}{q_{ij}^d} . &
\end{aligned}
\tag{28}
$$

And the bias part is expressed as

$$
\hat{\eta}_j(t) - \hat{\eta}_j^g(t) = \frac{\sum_{i \in I_j} \mathrm{E}\left[ g_{ij}(t) \frac{\partial g_{ij}(t)}{\partial \eta_j(t)} \right]}{\sum_{i \in I_j} (f_{ij}^*(t)\hat{r}_{ij})^2} .
\tag{29}
$$

Some insights can be gained from the results above. As by definition $\eta_j(t) \geq 0$, the estimate $\hat{\eta}_j(t)$ given by Eq. (19) tends to be biased high in our model. The value of $\hat{r}_{ij}$ plays a role in the minimization of bias, as it only appears in the denominator in Eq. (29). Similarly, if the value of $\hat{r}_{ij}$ is similar for different words, then larger values of $q_{ij}^d$ and $f_{ij}^*$ will reduce the bias, as seen from Eq. (28) – therefore, we should consider including preferentially in our analysis words with larger values of $q_{ij}^d$, $f_{ij}^*$ and $\hat{r}_{ij}$. Considering that the value of $\eta_j(t)$ affects the bias as well, which is not simply linear, we are led to consider adaptive or iterative criteria for word choice, which will in general depend on the true (and unknown) value of $\eta_j(t)$.

**Case 2:** Gaussian distribution for $\delta_{ij}(f_{ij})$, e.g., $\delta_{ij}(f_{ij}) \sim \mathcal{N}(0, f_{ij}\sigma_{ij}^2)$, which is justified empirically in the Appendix, Figure 14 and Figure 16. As a result,

$$
\begin{aligned}
\xi_{ij}(t) &= (\hat{r}_{ij} + \epsilon_{ij}^\eta(q, f, t)) \delta_{ij}(\eta_j(t) f_{ij}^*(t)) \\
&\quad + \delta_{ij}'(f_{ij}^*(t)) \\
&= \sqrt{\eta_j(t) f_{ij}^*(t)} (\hat{r}_{ij} + \epsilon_{ij}^\eta(q, f, t)) \delta_{ij}(1) \\
&\quad + \sqrt{f_{ij}^*(t)} \delta_{ij}'(1)
\end{aligned}
\tag{30}
$$

Therefore, we can define $g_{ij}^c(t)$ and $\xi_{ij}^c(t)$:

$$
\begin{aligned}
g_{ij}^c(t) &= \eta_j(t) f_{ij}^*(t) \epsilon_{ij}^\eta(q, f, t) \\
&\quad + \sqrt{\eta_j(t) f_{ij}^*(t)} (\hat{r}_{ij} + \epsilon_{ij}^\eta(q, f, t)) \delta_{ij}(1)
\end{aligned}
\tag{31}
$$

$$
\xi_{ij}^c(t) = \sqrt{f_{ij}^*(t)} \delta_{ij}'(1)
\tag{32}
$$

As $\xi_{ij}^c(t)$ doesn't depend on $\eta_j(t)$, the loss function under this assumption is:

$$
\begin{aligned}
L_{j,t,g}^c(\eta_j) & \\
= \frac{1}{n_j} \sum_{i \in I_j} (h_{ij}(t) &- \eta_j(t) x_{ij}(t) - g_{ij}^c(t))^2 \\
= \frac{1}{n_j} \sum_{i \in I_j} (\xi_{ij}^c(t))^2 . &
\end{aligned}
\tag{33}
$$

And we will get a complex expression for the bias part like Eq. (23), which gives us similar conclusions. (Some calculations are in the Appendix.)

Finding criteria for selecting the words that are included in the frequency change analysis greatly reduces the computational complexity compared to trying different word combinations. Our analysis of noise models gives some insights into these criteria, such as $q_{ij}^d$ and $\hat{r}_{ij}$.

### 4.4. Approximations in real data

Unfortunately, we cannot know the true value of $f_{ij}^*(t)$ in the ChatGPT era, but we can replace it with the estimation

$\hat{f}_{ij}^*(t)$ based on the word frequency before ChatGPT was introduced. As our objective is to identify the words that ChatGPT "likes" (or "dislikes") to use compared to academic researchers on average, we assume that the frequencies of these words should remain stable without ChatGPT, i.e., we take the average of the pre-ChatGPT periods before $T_0$ as following:

$$f_{ij}^*(t) = \frac{1}{\#\{t \le T_0\}} \sum_{t \le T_0} f_{ij}^d(t), \text{if } t > T_0. \quad (34)$$

Considering that the noise in real data is likely highly complex, we did not estimate the variance of $\epsilon_{ij}(\cdot)$. Instead, we used ChatGPT to process additional abstracts (on top of those used to estimate $r_{ij}$), and used the resulting frequencies as calibration for the bias and noise.

### 4.5. Calibration and test

In order to verify the theoretical and practical validity of our approach, we used calibrations and tests, with ChatGPT-processed abstracts mixed with real abstracts.

As previous analyses have demonstrated, with the goal of reducing bias in estimation, different selected words are likely to correspond to the different (unknown) ground truth value of $\eta_j(t)$. Therefore, we construct $N$ different sets of abstract data for calibration, $D_n$, with its correspond mixed ratio of ChatGPT-processed abstracts, $\eta_n$, as

$$(D_n, \eta_n), n \in \{1, 2, \ldots, N\} \quad (35)$$

and similarly for test data $T_{n'}$ and $\eta'_{n'}$,

$$(T'_n, \eta'_n), n' \in \{1, 2, \ldots, N'\}. \quad (36)$$

And for one pair of $(D_n, \eta_n)$ and a specific word choice requirement $q_k$ (for example, $r_{ij}^* > 0.1$ and $\frac{\hat{r}_{ij}+1}{\hat{r}_{ij}^2} < \frac{0.1+1}{0.1^2}$), the efficiency can be defined as

$$e(D_n, \eta_n, q_k) = |\eta_n - \hat{\eta}_n(D_n, q_k)| \quad (37)$$

where $\hat{\eta}_n(D_n, q_k)$ is the estimate of $\eta_n$ using Eq. (19) and the words set $I_j$ can be derived from $q_k$, denoted $I_j(q_k)$

For a given set of $q_k$ (examples can be found in the Appendix), we are looking for the best one minimizing $e(D_n, \eta_n, q_k)$, denoted $q(D_n, \eta_n)$, which is the calibration part.

For the test data $T_{n'}$, the estimate of $\eta_{n'}$ is calculated from Eq. (19) with different $I_j$, based on different $q(D_n, \eta_n)$ obtained in the calibration procedure.

Because of the goal of the calibration, word choice may well actually introduce a new bias to neutralize the original bias, so that the estimate is not necessarily higher in the test results than the ground truth.

## 5. Results

### 5.1. Calibration and test results

To calibrate the choice of set $I_j$, we use different mixing ratios, in proportion to the value of $\eta_j(t)$. In addition, we only consider the 10,000 words with the highest frequency in the Google Ngram dataset.

We continue our simulations based on GPT-3.5. As the training data for GPT-3.5 is up to September 2021, abstracts submitted later than this time are considered: 20,000 abstracts in period 13 to estimate $r_{ij}$, 10,000 abstracts in period 12 for calibration, and 10,000 abstracts in period 14 for testing.

We used the first 10 periods before ChatGPT was introduced, to estimate $f_{ij}^*(t)$, as they weren't influenced by ChatGPT, which means $T_0 = 10$ and $\#\{t \le T_0\} = 10$ in Eq. (34).

We take $\{\eta_n\} = \{0, 0.05, 1, \ldots, 0.45, 0.5\}$ and $m = 1$, which means $N = \#\{(D_n, \eta_n)\} = 11$. Then the 11 $I_j$ (with possible repetitions), obtained from mixed data with 11 corresponding $\eta_n$ of period 12, were used for $\eta'_n$ estimation in the test data (period 14). Other parameters can be found in the Appendix.
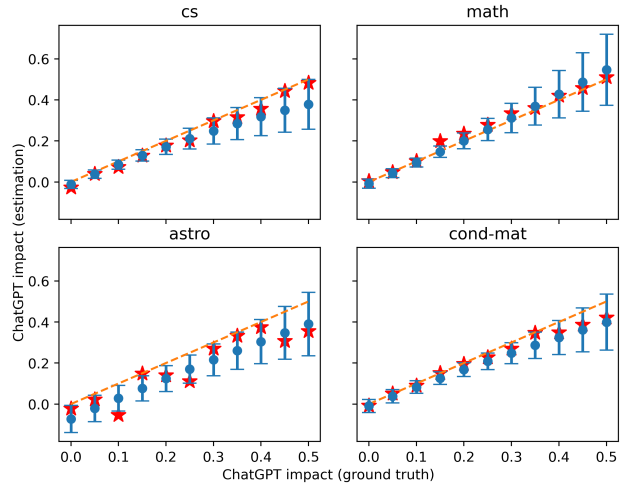


*Figure 6.* Test results for simulated admixtures of abstracts in period 14. The error bars represent the standard deviation of the estimation results, and the red star is the estimated value of $\eta'_n$ from test data based on optimal $I_j$ with the same mixed ratio $\eta_n$ as in the calibration data. The orange dashed lines correspond to perfect estimation.

The results using the same prompt for generating calibration and test data are shown in Figure 6, with injected mixed ratio (i.e., ChatGPT impact) $\eta'_n$ from 0 to 0.5. It is clear that when the calibration and test sets are mixed in the same ratio, word combinations that achieve better estimates on the calibration set generally work better on the test set, as

well.

Unlike in Figure 6 where we normalized the word frequency by the total number of abstracts, we normalized it by the total number of words for one period in Figure 7. The trends remain similar, albeit different in detail.
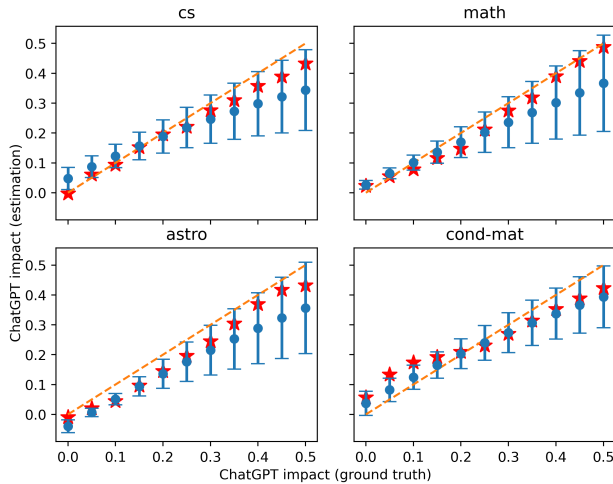


*Figure 7.* Similar to Figure 6, but normalized to the total number of words, rather than the total number of abstracts.

Because one may use a wide variety of prompts in practical applications, we also evaluated the robustness of our approach by adopting a different prompt for generating the test data than the one we used for calibration. The corresponding results in Figure 8 use the following prompt:

*"Please rewrite the following paragraph from an academic paper:"*

In this example, we add the word "please" and make it clear that this comes from an "academic paper", replacing "revise" with "rewrite".

Although the quantitative results of our tests were not as good as before, the errors were still small at lower mixed ratios, which also illustrates the robustness of our method. This is understandable because in data generated with different prompts, not all of our previous assumptions hold, and the estimate of $\hat{r}_{ij}$ on $r_{ij}$ in our model may be biased. We can also note that most of our estimates in Figure 8 are on the high side relative to the ground truth, most likely because we use a more precise prompt for the test data here, making the frequency change rate of the relevant words higher.

### 5.2. Estimation from real data

The estimates of ChatGPT impact on the real data are shown in Figure 9 and Figure 10. Based on our calibration results,



*Figure 8.* Similar to Figure 7, but with a different prompt for test data than used in calibration data.

we chose 11 words set $I_j$ for different injected values of $\eta_n$. According to the results of the first estimation about $\eta_j(t)$, we found the three values of $\eta_n$ that were closest to the mean of the first estimation and used their optimal word set $I_j$ in the calibration procedure for a second estimation, leading to the triangle points shown in the figures.

Despite mild differences in the estimates under the two different normalizations, the conclusions are essentially the same. Our estimates on $\eta_j(t)$ hover around 0 until 2023, which gives reassurance of the reliability of our methodology. More and more abstracts are being influenced by ChatGPT, especially in the *cs* category, starting from December 2022, after the release of ChatGPT.

Our estimate indicates that the density of ChatGPT style texts of the most recent time period in this category is around 35%, when we use the results of one simple prompt, "revise the following sentences", as a baseline. By contrast, we detected a much smaller uptick in ChatGPT impact in *math*, while *astro* and *cond-mat* both reach values between 10% and 20%, approximately.

It is important to note that our ChatGPT impact here is a relative value that corresponds to the change in word frequency from the use of simple prompts. More precise prompts, both in reality and in simulation, could potentially lead to an impact value greater than 1.

## 6. Conclusions

Is ChatGPT transforming academics' writing style? An important question before these discussions is the evaluation of the actual penetration of the usage of ChatGPT in aca-
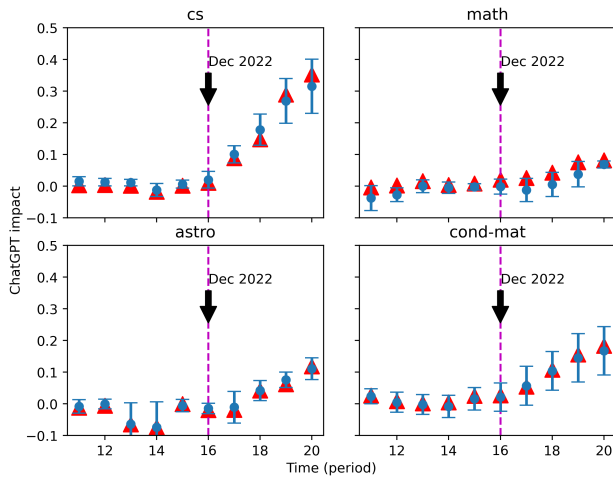
*Figure 9.* Estimates of $\eta_j(t)$ (i.e., ChatGPT impact) from real data. Word frequencies were normalized on the number of abstracts in each period before the estimation was performed. The error bars represent the standard deviation of the estimation results, using 11 different word sets $I_j$ obtained in the calibration procedure with 11 different $\eta_n$. The points of the triangle represent the average of the 3 estimates, corresponding to the 3 word selection requirements $q$ based on the 3 $\eta_n$ closest to the mean of the previous 11 estimates.

demic writing – without a quantitative estimate, the debate is founded on anecdotal evidence.

We have demonstrated here that a simple statistical analysis of word frequency is sufficient to detect and analyze the impact of ChatGPT in arXiv abstracts, which is easily extendable to other subjects and to the complete text of articles. We found convincing evidence of a change in word frequency after ChatGPT's release, consistent with predictions obtained from simulating ChatGPT's impact from possible users' prompts. The most enthusiastic community (among the four we investigated) in terms of ChatGPT adoption appears to be that of computer scientists, a result that is perhaps unsurprising. Mathematicians, by contrast, are the least keen.

Our estimates are founded on a population level and based on the output of simple prompts. Using more precise prompts, it is entirely possible to achieve abstracts that are more ChatGPT-like texts than our simulations. In addition, in the real world people might use large language models other than ChatGPT to revise articles, which may have similar but not identical word preferences to ChatGPT.

Another central take-away from this article is that we can monitor the impact of ChaGPT on academic writing by using simple and transparent statistical methods (e.g., word frequencies) rather than black-box GPT detectors. At the same time, we also believe that better estimates can be

*Figure 10.* The same estimates as in Figure 9, but word frequencies were normalized by the number of words.

made by more rigorous analysis, such as considering more complex noise terms.

## 7. Discussion

The debate around the usage of generative models such as ChatGPT in academic writing is multi-faceted: from fears of lowering rigour due to "hallucinations" to uncertainty about the actual sources of AI-produced text. It is however indisputable that tools such as ChatGPT also have positive impacts: they help non-English native writers to improve the quality and flow of their text, as well as to translate into English from their mother tongue or vice versa. In this sense, generative AI is a great leveller, and as such it is a welcome addition to the academic's toolbox. What we need to be wary of is its use in fully generative mode, without expert human supervision – something that we have not addressed in this paper.

We are aware that our methods can be further improved. For example, our results follow from analyzing a set of words selected based on the value of $q_{ij}^d$ and $\hat{r}_{ij}$. It is actually possible to fine-tune this criterium for a more accurate word selection, which would theoretically give better results, but would be more computationally expensive. Similarly, trying a larger range of prompts should theoretically result in better estimates. We are more interested in the density of ChatGPT style texts and its relative value (comparisons between categories and over time) than in establishing how many people are using ChatGPT – this can be estimated with the help of questionnaires, and it is not possible to get an accurate estimate only based on simulated data.

As our results have shown, ChatGPT is having an increasing

impact on academic publications. This trend is hard to avoid, and we need to adapt gradually. With the increasing influx of young researchers, especially non-native English speakers, tools based on large language models, represented by ChatGPT, are transforming academic writing, at least for some disciplines. Even if you refuse to use them, you are likely to be influenced indirectly.

# References

AI4Science, M. R. and Quantum, M. A. The impact of large language models on scientific discovery: a preliminary study using gpt-4. *arXiv preprint arXiv:2311.07361*, 2023.

AlAfnan, M. A. and MohdZuki, S. F. Do artificial intelligence chatbots have a writing style? an investigation into the stylistic features of chatgpt-4. *Journal of Artificial intelligence and technology*, 3(3):85–94, 2023.

Amano, T., Ramírez-Castañeda, V., Berdejo-Espinola, V., Borokini, I., Chowdhury, S., Golivets, M., González-Trujillo, J. D., Montaño-Centellas, F., Paudel, K., White, R. L., et al. The manifold costs of being a non-native english speaker in science. *PLoS Biology*, 21(7):e3002184, 2023.

arXiv.org submitters. arxiv dataset, 2024. URL https://www.kaggle.com/dsv/7352739.

Casal, J. E. and Kessler, M. Can linguists distinguish between chatgpt/ai and human writing?: A study of research ethics and academic publishing. *Research Methods in Applied Linguistics*, 2(3):100068, 2023.

Cheng, S.-L., Tsai, S.-J., Bai, Y.-M., Ko, C.-H., Hsu, C.-W., Yang, F.-C., Tsai, C.-K., Tu, Y.-K., Yang, S.-N., Tseng, P.-T., et al. Comparisons of quality, correctness, and similarity between chatgpt-generated and human-written abstracts for basic research: Cross-sectional study. *Journal of Medical Internet Research*, 25:e51229, 2023.

Fergus, S., Botha, M., and Ostovar, M. Evaluating academic answers generated using chatgpt. *Journal of Chemical Education*, 100(4):1672–1675, 2023.

Fitria, T. N. Grammarly as ai-powered english writing assistant: Students' alternative for writing english. *Metathesis: Journal of English Language, Literature, and Teaching*, 5 (1):65–78, 2021.

Gao, C. A., Howard, F. M., Markov, N. S., Dyer, E. C., Ramesh, S., Luo, Y., and Pearson, A. T. Comparing scientific abstracts generated by chatgpt to real abstracts with detectors and blinded human reviewers. *NPJ Digital Medicine*, 6(1):75, 2023.

Hwang, S. I., Lim, J. S., Lee, R. W., Matsui, Y., Iguchi, T., Hiraki, T., and Ahn, H. Is chatgpt a "fire of prometheus" for non-native english-speaking researchers in academic writing? *Korean Journal of Radiology*, 24(10):952, 2023.

Kang, J., Yoon, J. S., and Lee, B. How ai-based training affected the performance of professional go players. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pp. 1–12, 2022.

Kasneci, E., Seßler, K., Küchemann, S., Bannert, M., Dementieva, D., Fischer, F., Gasser, U., Groh, G., Günnemann, S., Hüllermeier, E., et al. Chatgpt for good? on opportunities and challenges of large language models for education. *Learning and individual differences*, 103: 102274, 2023.

Kirchenbauer, J., Geiping, J., Wen, Y., Katz, J., Miers, I., and Goldstein, T. A watermark for large language models. In *International Conference on Machine Learning*, pp. 17061–17084. PMLR, 2023.

Liang, W., Izzo, Z., Zhang, Y., Lepp, H., Cao, H., Zhao, X., Chen, L., Ye, H., Liu, S., Huang, Z., et al. Monitoring ai-modified content at scale: A case study on the impact of chatgpt on ai conference peer reviews. *arXiv preprint arXiv:2403.07183*, 2024a.

Liang, W., Zhang, Y., Wu, Z., Lepp, H., Ji, W., Zhao, X., Cao, H., Liu, S., He, S., Huang, Z., et al. Mapping the increasing use of llms in scientific papers. *arXiv preprint arXiv:2404.01268*, 2024b.

Lingard, L., Chandritilake, M., de Heer, M., Klasen, J., Maulina, F., Olmos-Vega, F., and St-Onge, C. Will chatgpt's free language editing service level the playing field in science communication?: Insights from a collaborative project with non-native english scholars. *Perspectives on Medical Education*, 12(1):565, 2023.

Lund, B. D., Wang, T., Mannuru, N. R., Nie, B., Shimray, S., and Wang, Z. Chatgpt and a new academic reality: Artificial intelligence-written research papers and the ethics of the large language models in scholarly publishing. *Journal of the Association for Information Science and Technology*, 74(5):570–581, 2023.

Luo, Z., Xie, Q., and Ananiadou, S. Chatgpt as a factual inconsistency evaluator for abstractive text summarization. *arXiv preprint arXiv:2303.15621*, 2023.

Michel, J.-B., Shen, Y. K., Aiden, A. P., Veres, A., Gray, M. K., Team, G. B., Pickett, J. P., Hoiberg, D., Clancy, D., Norvig, P., et al. Quantitative analysis of culture using millions of digitized books. *science*, 331(6014):176–182, 2011.

Mundt, K. and Groves, M. A double-edged sword: the merits and the policy implications of google translate in higher education. *European Journal of Higher Education*, 6(4):387–401, 2016.

Noy, S. and Zhang, W. Experimental evidence on the productivity effects of generative artificial intelligence. *Science*, 381(6654):187–192, 2023.

Polak, M. P. and Morgan, D. Extracting accurate materials data from research papers with conversational language models and prompt engineering–example of chatgpt. *arXiv preprint arXiv:2303.05352*, 2023.

Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.

Varadi, M. and Velankar, S. The impact of alphafold protein structure database on the fields of life sciences. *Proteomics*, 23(17):2200128, 2023.

von Garrel, J. and Mayer, J. Artificial intelligence in studies—use of chatgpt and ai-based tools among students in germany. *Humanities and Social Sciences Communications*, 10(1):1–9, 2023.

Zhao, W. X., Zhou, K., Li, J., Tang, T., Wang, X., Hou, Y., Min, Y., Zhang, B., Zhang, J., Dong, Z., et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.

# A. Bonus!



*Figure 11.* Revise the following sentences:

# B. Period divisions

*Table 2.* First and last arXiv paper identifier of 20 periods.

| PERIOD | FIRST PAPER | LAST PAPER |
|---|---|---|
| 1 | 1805.08929 | 1810.00786 |
| 2 | 1810.00787 | 1902.00889 |
| 3 | 1902.00890 | 1905.13537 |
| 4 | 1905.13538 | 1909.11935 |
| 5 | 1909.11936 | 2001.06560 |
| 6 | 2001.06561 | 2005.02178 |
| 7 | 2005.02179 | 2008.04251 |
| 8 | 2008.04252 | 2011.09225 |
| 9 | 2011.09226 | 2103.01828 |
| 10 | 2103.01829 | 2106.04209 |
| 11 | 2106.04210 | 2109.09152 |
| 12 | 2109.09153 | 2112.12197 |
| 13 | 2112.12198 | 2204.01835 |
| 14 | 2204.01836 | 2207.06075 |
| 15 | 2207.06076 | 2210.10618 |
| 16 | 2210.10619 | 2301.10909 |
| 17 | 2301.10910 | 2304.13927 |
| 18 | 2304.13928 | 2307.10978 |
| 19 | 2307.10979 | 2310.09716 |
| 20 | 2310.09717 | 2401.02417 |



*Figure 12.* Please rewrite the following paragraph from an academic paper:

# C. arXiv categories

Formally, arXiv has 8 categories in total: physics, mathematics, computer science, quantitative biology, quantitative finance, statistics, electrical engineering and systems science, economics. The first 3 categories contribute the vast majority of arXiv articles, around 91% among the 1 million articles. Hence, we divided the physics papers into sub-categories: astrophysics, condensed matter, high energy physics, etc. The four categories (computer science, mathematics, astrophysics, condensed matter) we selected account for 70% of the total number of articles. To avoid repetition, we also only count the first category of the article for those that have multiple categories (cross-postings).

# D. Parameters

## D.1. ChatGPT simulations

- model: gpt-3.5-turbo-1106

- temperature: 0.7

- seed: 1106

- top_p: 0.2

## D.2. Calibration

- $\frac{1}{q_{ij}^d}$: 10, 20, 30, 40, 50, 60, 70, 80, 100, 150, 200, 500

- $\hat{r}_{ij}$: 0.1, 0.15, 0.2, 0,3, 0.4, 0.5, 0.6, 0.7, 0.8 (corresponding value of $\frac{\hat{r}_{ij}+1}{\hat{r}_{ij}^2}$)

For example, when we take $\frac{1}{q_{ij}^{d_j}} < 10$ and $\frac{\hat{r}_{ij}+1}{\hat{r}_{ij}^2} < \frac{0.1+1}{0.1^2}$ for abstracts in computer science, the words that satisfy the conditions are: 'the', 'is', 'for', 'by', 'be', 'this', 'are', 'i', 'at', 'which', 'an', 'have', 'but', 'we', 'all', 'they', 'one', 'has', 'their', 'other', 'there', 'more', 'new', 'any', 'these', 'time', 'than', 'some', 'only', 'two', 'into', 'them', 'our', 'under', 'first', 'most', 'then', 'over', 'work', 'where', 'many', 'through', 'well', 'how', 'even', 'while', 'however', 'high', 'given', 'present', 'large', 'research', 'different', 'set', 'study', 'important', 'several', 'e', 'further', 'including', 'often', 'provide', 'due', 'using', 'better', 'various', 'problem', 'show', 'problems', 'design', 'proposed', 'g', 'across', 'approach', 'existing', 'compared', 'task', 'learn', 'improve', 'achieve', 'novel', 'domain', 'demonstrate', 'introduce', 'propose', 'prediction'.

And when $\frac{1}{q_{ij}^{d_j}} < 50$ and $\frac{\hat{r}_{ij}+1}{\hat{r}_{ij}^2} < \frac{0.8+1}{0.8^2}$, the words are: 'i', 'would', 'so', 'some', 'what', 'out', 'work', 'very', 'because', 'much', 'good', 'way', 'great', 'here', 'since', 'might', 'last', 'end', 'means', 'having', 'thus', 'above', 'give', 'e', 'further', 'far', 'find', 'although', 'show', 'n', 'help', 'together', 'particular', 'whose', 'issue', 'according', 'addition', 'usually', 'art', 'especially', 'respect', 'works', 'shows', 'g', 'makes', 'hard', 'significant', 'run', 'address', 'particularly', 'idea', 'consider', 'includes', 'built', 'adopted', 'obtain', 'establish', 'useful', 'leading', 'performed', 'create', 'named', 'conducted', 'resulting', 'hence', 'findings', 'towards', 'prove', 'build', 'perform', 'moreover', 'describe', 'besides', 'demonstrated', 'via', 'presents', 'mainly', 'fail', 'namely', 'allowing', 'demonstrate', 'advances', 'suffer', 'overcome', 'introduce', 'accurately', 'identifying', 'enhance', 'crucial', 'etc', 'utilize', 'demonstrates', 'additionally', 'focuses', 'motivated', 'characterize'.

## E. Noise analysis

### E.1. Variance in real data

Abstracts in the *cs* category among the first 500,000 articles were divided into groups in chronological order, with the same number in each group. We counted the number of occurrences of each word within each group, and calculated the variance between the different groups. This was repeated as a function of the number of abstracts included in each group, and the results are shown in Figure 13.

Then we also analyzed the variance-to-mean ratio (defined as the variance of the sum of a word's counts divided by the mean of the sum) and the coefficient of variation (defined as the standard deviation of the sum divided by the mean of the
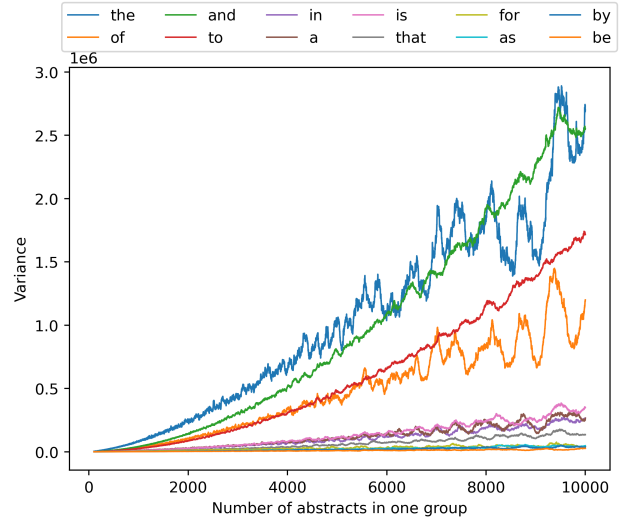


*Figure 13.* Variance of the word counts between groups of abstracts as a function of the number of abstracts included in each group.

sum) for the 12 most frequent words, as shown in Figure 14 and Figure 15, and the variance-mean ratio of further words as in Figure 16.
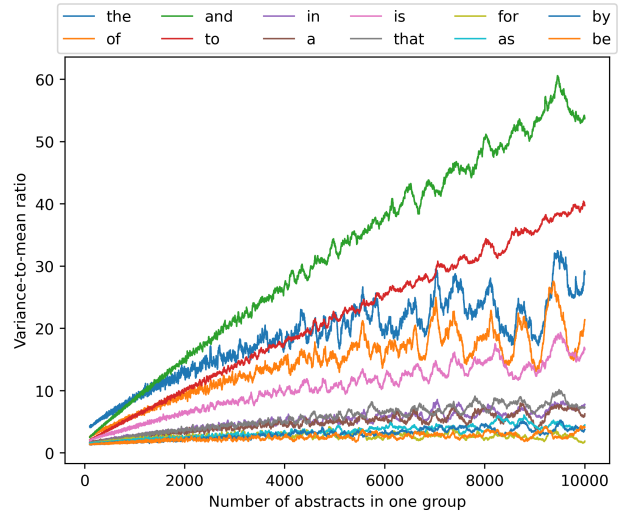


*Figure 14.* Variance-to-mean ratio.

We observe that, at least for a subset of the words considered here, the variance-to-mean ratios are essentially on the same scale (although there are words that do not follow this pattern). Therefore, a simple Gaussian distribution

$$\delta_{ij}(f_{ij}) \sim \mathcal{N}(0, f_{ij}\sigma_{ij}^2). \tag{38}$$
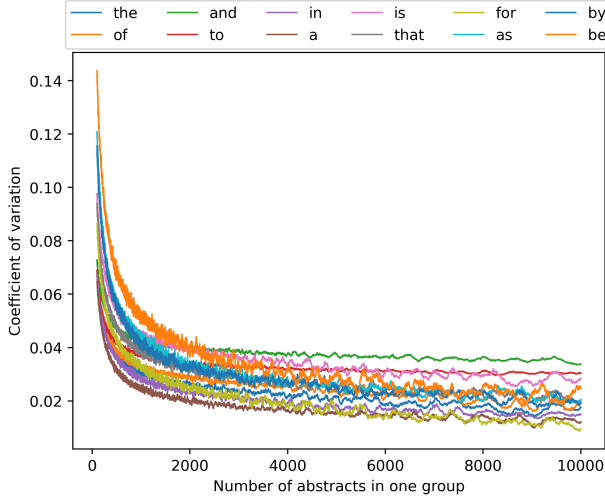
which corresponds to case 2, seems to be a reasonable ap-

*Figure 15.* Coefficient of variation.



*Figure 16.* Variance-to-mean ratio for some further words.

proximation.

### E.2. Calculation details

As in case 1, we set $\frac{\partial L_{j,t,g}^c(\eta_j)}{\partial \eta_j} = 0$ to obtain the new estimate $\hat{\eta}_j^g(t)$ corrected for bias and noise,

$$
\begin{aligned}
(\hat{\eta}_j^g(t) - \hat{\eta}_j(t)) &\sum_{i \in I_j} x_{ij}^2(t) \\
= \sum_{i \in I_j} &\frac{\partial g_{ij}^c(t)}{\partial \eta_j(t)} \left( h_{ij}(t) - \eta_j(t) x_{ij}(t) \right) \\
&- \sum_{i \in I_j} x_{ij}(t) g_{ij}^c(t) - \sum_{i \in I_j} g_{ij}^c(t) \frac{\partial g_{ij}^c(t)}{\partial \eta_j(t)}
\end{aligned}
\tag{39}
$$

where

$$
\begin{aligned}
\frac{\partial g_{ij}^c(t)}{\partial \eta_j(t)} =& f_{ij}^*(t) \epsilon_{ij}^\eta(q, f, t) + \eta_j(t) f_{ij}^* \frac{\partial \epsilon_{ij}^\eta(q, f, t)}{\partial \eta_j(t)} \\
&+ \frac{\sqrt{f_{ij}^*(t)}}{2\sqrt{\eta_j(t)}} (\hat{r}_{ij} + \epsilon_{ij}^\eta(q, f, t)) \delta_{ij}(1) \\
&+ \sqrt{\eta_j(t) f_{ij}^*(t)} \frac{\partial \epsilon_{ij}^\eta(q, f, t)}{\partial \eta_j(t)} \delta_{ij}(1) .
\end{aligned}
\tag{40}
$$

The bias part is also expressed as

$$
\hat{\eta}_j(t) - \hat{\eta}_j^g(t) = \frac{\sum_{i \in I_j} \mathrm{E}\left[ g_{ij}^c(t) \frac{\partial g_{ij}^c(t)}{\partial \eta_j(t)} \right]}{\sum_{i \in I_j} (f_{ij}^*(t) \hat{r}_{ij})^2} .
\tag{41}
$$

Also with the same assumptions for $\epsilon_{ij}(\cdot)$ and $\epsilon_{ij}^s(\cdot)$, $\epsilon_{ij}(f_{ij}) \sim \mathcal{N}(0, f_{ij}\sigma_{ij,\epsilon}^2)$ and $\epsilon_{ij}^s(f_{ij}) \sim \mathcal{N}(0, f_{ij}\sigma_{ij,\epsilon}^2)$.

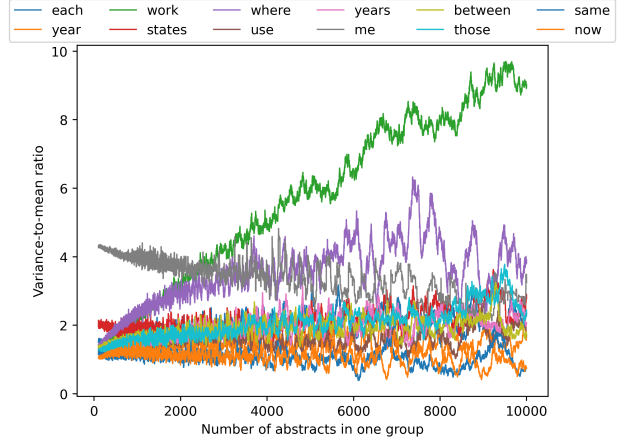then we can obtain an expression for $\epsilon_{ij}^\eta(q, f, t)$,

$$
\begin{aligned}
\epsilon_{ij}^\eta(q, f, t) =& \frac{\epsilon_{ij}(1)}{\sqrt{\eta_j(t) f_{ij}^*(t) + \sqrt{\eta_j(t) f_{ij}^*(t)} \delta_{ij}(1)}} \\
&- \frac{\epsilon_{ij}^s(1)}{\sqrt{q_{ij}^d}}
\end{aligned}
\tag{42}
$$

and its derivative,

$$
\begin{aligned}
&\frac{\partial \epsilon_{ij}^\eta(q, f, t)}{\partial \eta_j(t)} \\
&= \frac{-\left( 2f_{ij}^*(t)\sqrt{\eta_j(t)} + \sqrt{f_{ij}^*(t)}\delta_{ij}(1) \right) \epsilon_{ij}(1)}{4\sqrt{\eta_j(t)} \left( \eta_j(t) f_{ij}^*(t) + \sqrt{\eta_j(t) f_{ij}^*(t)}\delta_{ij}(1) \right)^{\frac{3}{2}}} .
\end{aligned}
\tag{43}
$$

Combining the above equations, we can get similar conclusions as in case 1.

## F. Other observations

We also define a change factor in the frequency of word $i$, $R_i'$, as follows:

$$
R_i' = \frac{\max_t(f_i'(t)) - \min_t(f_i'(t))}{\max_t(f_i'(t))}
\tag{44}
$$

where $f_i'(t)$ is the count of word $i$ in period $t$, normalized to the same value of $\sum_i f_i(t)$ for all periods $t$.

The total number of words in all abstracts of the first period is used as a base to normalize the frequency of words in the other periods, and the corresponding results are shown in Figure 17, Figure 18, and Figure 19.
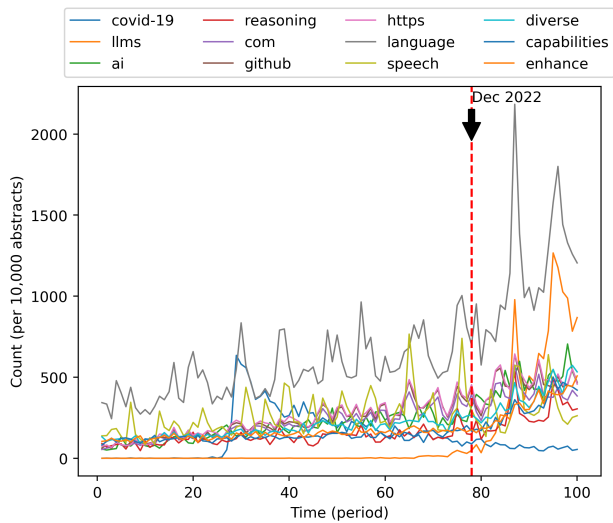
14

*Figure 17.* The 12 words with the highest change rate $R_i{}'$ and satisfying $\max_t(f_i{}'(t)) > 500$.
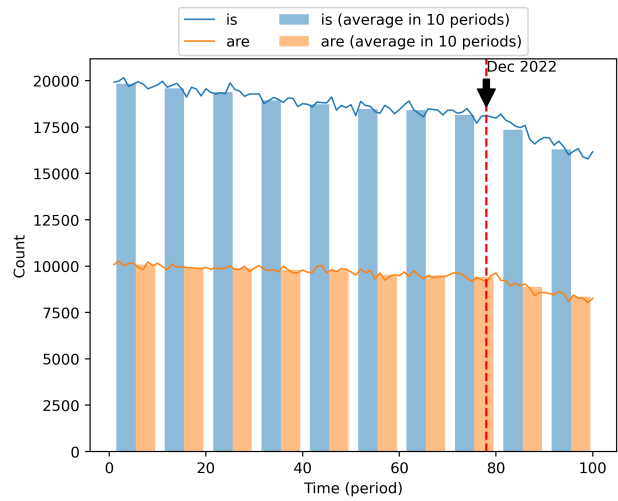


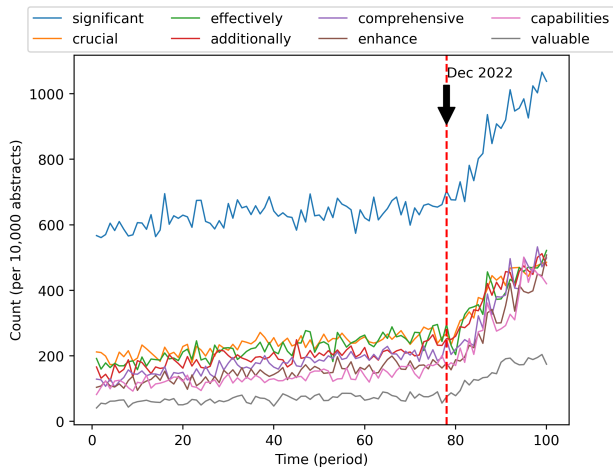*Figure 19.* The words "are" and "is" are decreasing in frequency in arXiv abstracts.



*Figure 18.* Examples of words with rapidly growing frequency in arXiv abstracts.