**Academic Integrity in Crisis: A Systematic Analysis of Questionable Research Practices**

(Technical Paper)

**Trust, Science, and the Digital War on Health Information**

(STS Paper)


A Thesis Prospectus Submitted to the

Faculty of the School of Engineering and Applied Science

University of Virginia • Charlottesville, Virginia

In Partial Fulfillment of the Requirements of the Degree

Bachelor of Science, School of Engineering


**Riley Tomek**

Fall, 2024


Technical Project Team Members

David Downer

Sean Ferguson

Anna Fisher


On my honor as a University Student, I have neither given nor received unauthorized aid on this assignment as defined by the Honor Guidelines for Thesis-Related Assignments

Signature _____ Date _____

Riley Tomek

Approved _____ Date _____

William T. Scherer, Department of Systems and Information Engineering

Approved _____ Date _____

Matthew Bolton, Department of Systems and Information Engineering

Approved _____ Date _____

Richard Jacques, Department of Engineering and Society

**A. Technical Topic Introduction**

Societal knowledge is upheld by publications and journal proceedings which dictate lifestyles, formulate ways of thinking, and advance the greater world of academia. Prior to addressing any research question, a general literature review is performed to establish a lens into the topic of interest and develop a framework of thinking. Typically, a literature review acts as a dependable source of credible research findings. However, over the past decade, the rise of digital information and publicly accessible artificial intelligence (AI) tools has disrupted these practices, introducing both opportunities for scientific advancement and risks of exploitation (Else, 2023).

Scientific research serves as a rich environment to be taken advantage of and is riddled with predatory journals, weak institutional oversight, manipulation of authorship metrics, and one unique incentive familiar to all higher education: to publish or perish (Callaway, 2023). In consequence, one estimate indicates that approximately 1 in 7 papers produced have some measure of fraudulent activity (Pachankis, Hatzenbuehler, & Safren, n.d.). In resistance, non-profit organizations, such as Retraction Watch and Pub Peer, serve as so-called 'journalist watchdogs' with online flagging forums and blogs calling out prominent retracted papers and journals (Retraction Watch, 2020). However, the current system of researchers, journals, and institutions is not capable of effectively detecting all misinformation from infiltrating into the greater research community at the rate it is increasing currently.

Additionally, the novelty of large language models developed by technology giants such as OpenAI, Google, Meta, and other institutions reveals a significant gap in oversight, as these large language models are rapidly evolving, often outpacing the traditional review and validation mechanisms that govern scientific research (Page et al., 2023). Alongside an existing flawed

researcher/publisher system, there is a new barrage of poorly understood research practices that raise concerns. Paper mills, or illegal organizations selling fabricated papers and authorships, are increasingly being identified. In a study by the Committee on Publication Ethics, the percentage of journal publications affected by paper milling can range from 2%-40% (Committee on Publication Ethics, n.d.). Varying by subject area, this ongoing crisis prevailing scientific research is a source of profit for illegal organizations and scaling significantly.

The consequences of paper-milling are sweeping into internationally renowned publishers and trusted journals. For instance, in January 2021, the Royal Society of Chemistry, a reputable academic organization based in the UK, announced several retractions from its journals. Their journal RSC Advances retracted a record number of 68 articles due to the "systematic publication of falsified studies." (Mullin, 2021). Another multinational publishing company, Wiley, indicated that 19 journals of one of its subsidiaries had been shut down due to "systematic manipulation of the publishing process" (Hindawi, n.d.). Thus, a compromised knowledge base is shaping perceptions and influencing decisions in ways that are not easily quantifiable nor detectable. Upholding research integrity in the digital era is essential for fostering a publishing system of trustworthy science.

**B. Technical Project Description**

The technical project aims to address the research question: *How do components of the current academic publishing system contribute to the dissemination of unreliable information and distort the foundation of credible, evidence-based knowledge, and how might artificial intelligence both exacerbate and help detect these fraudulent practices*? The developing malicious practices are extensive; Retraction Watch lists 127 reasons alone for why a paper may

be retracted, including but not limited to, paper mill activity, ethical violations, plagiarism, manipulation, and other frequent activities of misconduct. (Retraction Watch Data, n.d.). To consolidate existing, known methods of questionable research practices and identify new activities coupled with the rise of artificial intelligence, a literature survey will be conducted to classify the dimensions of research misconduct. This includes, but is not limited to, the actions of individuals engaged in research misconduct, the institutional frameworks which fail to detect such behavior, predatory journal activities, and the systemic pressures which drive the reduction of academic integrity. Following a review of current literature, a taxonomy will be produced and validated of all fraudulent and questionable research practices. This will be in effort to distinguish emerging fraudulent practices from less systemic 'sloppy' science. Lastly, an initial exploration into concepts, theories, and methods involving signal detection theory, graph theory, time series methods, and game theory will be proposed to identify these varied dimensions of fraudulent activity.

Prevalent threats to academic integrity on an institutional level include predatory journal behavior, paper-milling, pay-to-publish schemes, and artificial intelligence misuse. Encompassing predatory publishers, we will explore case studies and identify common patterns for identified predatory journals. Some common behaviors exhibited by predatory journals include claiming unverified impact factors, publication of low-quality or unrelated work, aggressive email marketing, unresponsive authors, and lack of transparency around open access and article processing charges (Johnson & Clark, 2020). As we study metrics and indicators of predatory journals, a preliminary investigation and literature survey will dictate the scope and prevalence of the effect on scientific research. Data available from Retraction Watch, CrossRef,

Google Scholar, Scopus, and other entities will provide analytical context into the scale predatory journal behavior.

Within the greater realm of publishers, our team will investigate schemes existing such as pay-to-publish, where journals will publish any article given a sum of money without regard to quality of work. For example, in 2019, an academic publisher OCIMS Group made fraudulent claims about their credentials, relayed hidden publication fees, and designated themselves as an impactful journal with phony metrics, looping in researchers in the process. This deception resulted in hefty fines by the Federal Trade Commission (Federal Trade Commission, 2019). For individuals with existing pressure to publish or in regions with limited access to reputable journals, the practice of giving into these kinds of publishers is alarmingly common. This demand also fuels pay-to-write businesses to engage in the act of "ghost-writing," an ongoing trend of falsified authorship to satisfy the demand of research requirements (Hu & Wu, 2013). The practice of "pay-to-publish" or "ghost-writing" not only dilutes the quality of available scientific literature, but it also creates a burden on researchers who may unintentionally fall prey to these journals and suffer reputational damage as a result. The infiltration of this institutional activity is polarizing to research integrity and poses a risk to the broader research community. Compiling a scope and quantity of each threat will assist in generalizing the current descriptive scenario.

Artificial intelligence plays a complex role in contemporary scientific research, with both positive and negative implications. On one hand, AI tools can assist in automating the detection of plagiarism, promoting the efficient production of high-quality scientific research. On the other hand, cases of misconduct involving AI are becoming increasingly prominent. These include the generation of falsified data and images, automated content creation, citation "hallucinations," and

even unethical peer review practices (Lee & Zhang, 2024). Such instances indicate the emergence of a larger issue within the research community.

Ironically, natural language processing (NLP) technologies can both facilitate the production of fraudulent scientific literature and aid in its detection. For example, Clear Skies is a company which utilizes machine learning, aims to investigate organized research fraud, and detects anomalies in author history to examine long-term patterns (Clear Skies, n.d.) Assistive tools and methods being developed by companies such as Clear Skies possess the potential to protect research and journal integrity by assisting in investigations and identifying bad actors before information is cited, credited, and circulated (Byrne & Christopher, 2020). Our team may examine and consider these tools and methods as well in our analysis.

As for the team's proposal, the detection of such fraudulent anomalies in the research community may be explored via graph theory, mapping connections across citations, authors, and networks of actors (Christopher, 2021). Additionally, as misconduct indicators are validated, the framework for flagging predatory institutions will be presented via signal detection theory, time series analysis, and machine learning methods, effectively improving the identification of trends that questionable author and institutional entities employ. Game theory is also another method which can model the academic pressure simulations in the researcher/publishing system, and applications of such event simulations will be considered (Lazebnik et al., 2023). In practice, formulation of standardized practices, methods, and frameworks for the greater research community will assist in predicting and identifying perpetrators of fraudulent activity.

**STS Topic Introduction**

The following discussion aims to identify and investigate the widespread prevalence of misinformation, particularly concerning health, nutrition, hormones, and the reproductive system in women. The current landscape is dominated by non-expert sources, harmful dieting fads, and invalid content that often fail to consider the unique physiological needs of women at various life stages. The consequence of this widespread reach in online content and research in medicine leads to confusion and dangerous practices, contributing to mental and physical health challenges such as hormonal imbalances, stress, and fatigue.

Leading health risks in the United States, such as cardiovascular disease, have continually manifested uniquely in women and have not been investigated thoroughly. The presence of risk factors for heart disease, autoimmune disorders, diabetes, and reproductive issues are misunderstood and dictate the outcomes of diagnosis and personal care experience. For example, primary care physicians are more concerned with diagnosing health issues related to weight issues and breast cancer in women as opposed to heart disease, which stands as the leading cause of death for women (Bairey Merz et al., 2017). Enabling healthcare providers, patients, and the public with reliable information in an ever-growing digital age is critical when such disparities exist between concerns versus respective treatments. Specifically, there is a lack of evidence that women possess adequate health literacy regarding the prenatal period, a critical time frame that determines reproductive outcomes. These outcomes are significantly exacerbated by age, ethnicity, and socioeconomic background (Meldgaard 2022). Misled diagnoses from professional resources force women to seek treatment elsewhere, setting up opportunities for falsehoods and health fads to attempt to solve the problems not addressed by primary care physicians.

In addressing the compromised knowledge base, identifying the platforms and health topics which source the most health-related information to the public is critical to recognizing where this problem lies. Platforms contributing to misinformation include, but are not limited to, TikTok, Instagram, and prevalent social media sites, where studies have shown less than half of health-related video content is factual (Dimitroyannis et al., 2024). Health information is readily accessible on these platforms, but sources are not consistently verified. Additionally, around 36% of Americans have poor health literacy, contributing to a feedback loop of creator's producing questionable content not backed by science (Shieh 2009). Medical research involving women stands at a unique vantage point where the existing knowledge base is flawed, and social media presents the opportunity for the dissemination of misinformation.

Knowledge about the menstrual cycle, prenatal period, and value of nutrition during all phases of life is critical to providing life-changing treatments. However, there have existed consistent barriers to providing quality and reliable research in women's health. Up until the early 90's, women were not required to be considered in clinical trials unless a distinct reasoning for exclusion was provided (Mastroianni 1999). For example, the study of a disease that affects both men and women may select only male participants to control for the anomalies detected during hormonal fluctuations or pregnancy (Shieh 2009). While this reasoning provides the basis for inconsistent knowledge from the greater medical community, the age of big data only compounds this effect. Novelty large language models (LLM's) are also being recorded as supplier's of significant health-related research, but have been proven to lack safeguards and may hallucinate citations in sourcing health information (Menz et al., 2024). The addition of rapidly produced digital content and mass of information presents a unique challenge to this demographic.

**Conclusion**

As the knowledge base builds upon itself, information that is inconsistent with science may hurt the public. The technical deliverable will assist in preventing infiltration of fraudulent content being processed by society via the implementation of a variety of data analytics and graphing techniques. The science, technology, and society topic discussion aims to address and distinguish misleading information relating to a disparity in health literacy, the leading sources and methods for evaluating misinformation, and implications relating to this concept in the digital age. Both deliverables are loosely coupled and touch on concepts of a compromised knowledge base from questionable research practices, the widespread adoption of AI, and systemic bureaucracies.

**References**

1) Bairey Merz, C, Andersen, H, Sprague, E. et al. Knowledge, Attitudes, and Beliefs Regarding Cardiovascular Disease in Women: The Women's Heart Alliance. *JACC*. 2017 Jul, 70 (2) 123–132.https://doi.org/10.1016/j.jacc.2017.05.024

2) Byrne, J. A., & Christopher, J. (2020). Digital magic, or the dark arts of the 21st century-how can journals and peer reviewers detect manuscripts and publications from paper mills?. *FEBS letters*, *594*(4), 583–589. https://doi.org/10.1002/1873-3468.13747

3) Callaway, E. (2023, September 25). AI tools that generate research papers are becoming ubiquitous—but do they work? Nature. https://www.nature.com/articles/d41586-023-02553-1

4) Clear Skies. (n.d.). Homepage. https://clear-skies.co.uk/

5) Christopher J. (2021). The raw truth about paper mills. *FEBS letters*, *595*(13), 1751–1757. https://doi.org/10.1002/1873-3468.14143

6) Committee on Publication Ethics. (n.d.). Retractions in the scientific literature. COPE. https://publicationethics.org/node/55256

7) Dimitroyannis, R., Fenton, D., Cho, S., Nordgren, R., Pinto, J. M., & Roxbury, C. R. (2024). A Social Media Quality Review of Popular Sinusitis Videos on TikTok. *Otolaryngology--head and neck surgery : official journal of American Academy of Otolaryngology-Head and Neck Surgery*, *170*(5), 1456–1466. https://doi.org/10.1002/ohn.688

8) Else, H. (2023, January 27). The battle against paper mills that churn out fake science. Nature. https://www.nature.com/articles/d41586-023-00191-1

9) Federal Trade Commission. (2019, April 3). Court rules in FTC's favor against predatory academic publisher OMICS Group, imposes $50.1 million judgment. Federal Trade Commission. https://www.ftc.gov/news-events/news/press-releases/2019/04/court-rules-ftcs-favor-against-predatory-academic-publisher-omics-group-imposes-501-million-judgment

10) Hindawi. (n.d.). Evolving our portfolio in response to integrity challenges. Hindawi. https://www.hindawi.com/post/evolving-our-portfolio-response-integrity-challenges

11) Hu, Z., & Wu, Y. (2013). An empirical analysis on number and monetary value of ghostwritten papers in China. *Current Science, 105*(9), 1230–1234. http://www.jstor.org/stable/24098932

12) Institute of Medicine (US) Committee on the Ethical and Legal Issues Relating to the Inclusion of Women in Clinical Studies; Mastroianni, A. C., Faden, R., & Federman, D. (Eds.). (1999). Women and health research: Ethical and legal issues of including women in clinical studies: Volume 2: Workshop and commissioned papers. National Academies Press. https://www.ncbi.nlm.nih.gov/books/NBK236575/

13) Johnson, D. S., & Clark, T. A. (2013). Retractions in scientific publishing: Root causes and implications. PubMed Central. https://pmc.ncbi.nlm.nih.gov/articles/PMC3702092/#ref7

14) Lee, C., & Zhang, H. (2024). Research integrity in the era of artificial intelligence: Challenges and opportunities. Medicine (Baltimore), 103(27). https://journals.lww.com/md-journal/fulltext/2024/07050/research_integrity_in_the_era_of_artificial.41.aspx

15) Lazebnik, T., Beck, S. & Shami, L. Academic co-authorship is a risky game. *Scientometrics* **128**, 6495–6507 (2023). https://doi.org/10.1007/s11192-023-04843-x

16) Meldgaard, M., Gamborg, M., & Maindal, H. T. (2022). Health literacy levels among women in the prenatal period: A systematic review. Sexual & Reproductive Healthcare, 34, 100796. https://doi.org/10.1016/j.srhc.2022.100796

17) Menz B D, Kuderer N M, Bacchi S, Modi N D, Chin-Yee B, Hu T et al. Current safeguards, risk mitigation, and transparency measures of large language models against the generation of health disinformation: repeated cross sectional analysis *BMJ* 2024; 384 :e078538 doi:10.1136/bmj-2023-078538

18) Mullin, R. (2021, January 27). Paper mill hits RSC journals. Chemical & Engineering News. https://cen.acs.org/policy/publishing/Paper-mill-hits-RSC-journals/99/web/2021/01

19) Pachankis, J. E., Hatzenbuehler, M. L., & Safren, S. A. (n.d.). [Dataset]. OSF. https://osf.io/5rf2m/

20) Page, M. J., McKenzie, J. E., Bossuyt, P. M., Boutron, I., Hoffmann, T. C., Mulrow, C. D.,... & Moher, D. (2023). PRISMA 2020 explanation and elaboration: Updated guidance and exemplars for reporting systematic reviews. *BMJ*, 372(n160). https://pmc.ncbi.nlm.nih.gov/articles/PMC10166793/

21) Patel, R. K., & Chawla, N. (2023). AI-driven research ethics. PubMed. https://pubmed.ncbi.nlm.nih.gov/38431902/

22) Retraction Watch. (2020, October 28). Following Retraction Watch and PubPeer posts, journal upgrades correction to a retraction. Retraction Watch. https://retractionwatch.com/2020/10/28/following-retraction-watch-and-pubpeer-posts-journal-upgrades-correction-to-a-retraction/

23) Retraction Watch Data. (n.d.). [Database]. Crossref. https://gitlab.com/crossref/retraction-watch-data

24) Shieh, C., & Halstead, J. A. (2009). Understanding the impact of health literacy on women's health. Journal of Obstetric, Gynecologic, & Neonatal Nursing, 38(6), 601–612. https://doi.org/10.1111/j.1552-6909.2009.01059.x