

A CORPUS-BASED ANALYSIS OF COLLOCATIONAL PATTERNS AND SEMANTIC PROSODY OF AI-RELATED VOCABULARY IN CONTEMPORARY ENGLISH MEDIA DISCOURSE

Beenish Siraj

Degree BS English applied linguistics.

Email id beenishsiraj1993@gmail.com

Shanza Nadeem

MPhil English Linguistics

Email id shanzanadeem3.87@gmail.com

Laiba Khan

MPhil English Linguistics

Email id laibakhan8030@gmail.com

Zainab Shakoor

MPhil English Linguistics

Email id zainabshakoor014@gmail.com

Abstract

This study explores the linguistic patterns and framing of artificial intelligence (AI) in academic discourse, focusing on the most frequently used AI-related terms and their collocational relationships. A corpus-based approach was employed to analyze AI-related texts from 2025, utilizing tools such as LancsBox to identify key terms and significant collocates. The findings reveal that terms like AI, learning, and model dominate the discourse, reflecting the central role of machine learning and neural networks. Additionally, collocational analysis shows that AI-related terms often co-occur with words like language models and agentic AI, highlighting emerging themes in AI research. Visualizations such as a word cloud and collocate frequency heatmap provide insights into the linguistic framing and semantic prosody of these terms, with a generally positive portrayal of AI technologies. This research contributes to a deeper understanding of AI discourse in academic contexts and offers a foundation for future studies on the evolving language surrounding AI.

Keywords: Artificial Intelligence, AI Discourse, Collocational Patterns, Semantic Prosody, Word Cloud, Bigram Analysis

Introduction

Artificial Intelligence (AI) has become a ubiquitous topic in contemporary discourse, especially within the realm of media. As AI technologies evolve rapidly, their portrayal in news outlets plays a pivotal role in shaping public perception and guiding the societal narrative around these advancements. Media representations of AI often oscillate between highlighting its promising potential to solve complex societal issues and raising concerns about its possible threats and risks (Nguyen & Hekman, 2024). The widespread use of AI tools such as OpenAI's ChatGPT has further amplified public interest and debate, with discussions revolving around both the transformative capabilities of AI and the dangers it may pose, particularly regarding ethical issues and unforeseen consequences (Roe & Perkins, 2023).

While previous studies have explored the thematic and narrative framing of AI in the media, less attention has been paid to the fine-grained, corpus-based analysis of AI-related vocabulary. Specifically, understanding the collocational patterns and semantic prosody of AI-related terms within media discourse remains underexplored. Collocational patterns refer to the habitual co-

occurrence of words, which can reveal how certain terms are framed contextually (Hauser & Schwarz, 2023). Semantic prosody, on the other hand, refers to the subtle positive or negative associations that words acquire through their typical co-occurrence with other terms, which can further influence evaluative judgments about AI (Hauser & Schwarz, 2023). This study aims to fill this gap by analyzing how AI-related vocabulary is used in English-language media, identifying high-frequency keywords, and examining their collocates and semantic prosody to explore how language contributes to public understanding and the framing of AI.

Research Objectives

1. To compile a diachronic corpus of English-language news articles (2025) discussing artificial intelligence from both technology-focused and general news outlets.
2. To identify high-frequency AI-related keywords and their significant collocates within the corpus.
3. To analyze the semantic prosody of AI-related terms by examining their collocational environments.
4. To compare the collocational patterns and semantic prosody of AI-related vocabulary across different outlet types and over the course of 2025.

Research Questions

1. What are the most frequent AI-related keywords used in 2025 English-language media, and what are their significant collocates?
2. How do collocational patterns of AI-related terms vary across different media outlets in 2025?
3. What is the semantic prosody of AI-related terms in 2025, and how does it influence their evaluative judgment?
4. How do the collocational patterns and semantic prosody of AI-related terms evolve throughout 2025?

Significance of the Study

This study is significant in advancing our understanding of how artificial intelligence (AI) is portrayed in contemporary media. By focusing on the collocational patterns and semantic prosody of AI-related vocabulary, the research fills a gap in the current literature, offering a detailed analysis of how specific AI terms are framed in news discourse. The diachronic focus on the year 2025 ensures that the study captures the evolving nature of AI discourse, particularly in response to rapid technological advancements like ChatGPT and other large language models. The findings will provide valuable insights into the language used by the media to represent AI, highlighting the positive or negative connotations associated with AI-related terms. This has important implications for how the public perceives AI, influencing societal attitudes toward its potential and risks. Additionally, the study's results will be of practical use to policymakers, AI developers, and educators, helping them craft more informed, balanced narratives about AI. By understanding the media's framing of AI, stakeholders can address public concerns more effectively, ensuring a more nuanced and informed dialogue about the role of AI in society.

Literature Review

The representation of artificial intelligence (AI) in the media has been a growing subject of study in recent years, with scholars examining how AI is framed across different media platforms. Media discourse plays a crucial role in shaping public perception of AI, often oscillating between highlighting its potential to solve global problems and warning about its risks. Studies such as those by Nguyen and Hekman (2024) highlight the framing of AI as both a tool for innovation and

a source of societal disruption, with concerns about privacy, ethical dilemmas, and existential threats frequently being raised.

While much of the existing literature focuses on thematic and narrative analyses of AI in media, the role of language, particularly collocational patterns and semantic prosody, has been less explored. Hauser and Schwarz (2023) define semantic prosody as the subtle positive or negative associations that words acquire based on their habitual co-occurrence with other words. This phenomenon plays a significant role in influencing how certain terms, even when seemingly neutral, may elicit positive or negative judgments from the public. For example, words like "automation" often co-occur with terms that suggest both efficiency and job displacement, thereby creating a dual perception of AI technologies.

Roe and Perkins (2023) explore the portrayal of AI and tools like ChatGPT in UK media, noting a predominant focus on the potential dangers of AI, such as job loss and societal disruption, over its beneficial uses. Their findings align with previous research by Brennen (2018), who observed that the media often frames AI through sensationalist language that emphasizes its risks, rather than offering a balanced view of its potential to improve societal processes.

In contrast, studies like those by MacRitchie and Seedat (2008) suggest that media outlets with a more neutral or positive framing of AI often focus on its practical applications, such as AI's ability to enhance productivity or creativity. These studies demonstrate the importance of analyzing not just the themes in media representations of AI but also the linguistic features, such as collocation and semantic prosody, that shape public attitudes toward this technology.

Methodology

This study will use LancsBox, a powerful corpus analysis tool, to perform a detailed investigation of the collocational patterns and semantic prosody of AI-related vocabulary in academic discourse. LancsBox allows for comprehensive corpus management and analysis, providing a robust platform for extracting linguistic features such as word frequencies, collocates, and concordances.

The data collection process will involve selecting AI-related articles published in 2025 from the arXiv preprints repository. The articles will be filtered based on keywords like "artificial intelligence," "ChatGPT," "machine learning," and other relevant terms. After compiling the dataset, the articles will be imported into LancsBox for analysis. Once the corpus is loaded, collocational analysis will be conducted using LancsBox's built-in tools for word frequency extraction and collocate identification. This will help identify the most frequent AI-related terms and their statistically significant co-occurring words. LancsBox's collocate analysis features, including mutual information and log-likelihood measures, will be used to determine which words are most strongly associated with the target AI vocabulary.

For semantic prosody analysis, LancsBox will facilitate a closer examination of the surrounding context in which key AI terms appear. By analyzing the concordance lines of these words, the study will evaluate whether AI-related terms tend to be framed in positive, neutral, or negative contexts. LancsBox's ability to visualize collocational relationships and its integration with various statistical tests will support a nuanced analysis of the semantic prosody of AI vocabulary in academic texts. By leveraging LancsBox, this study aims to provide a detailed, quantitative, and qualitative examination of how AI is linguistically framed in academic research as of 2025. The use of LancsBox's corpus analysis tools will ensure a rigorous and systematic approach to understanding the linguistic features that shape public perceptions of AI in academic discourse.

Data Analysis

1. Frequency Table of AI-Related Terms

This table will present the most frequent AI-related terms in the corpus, showing how often each term appears.

Table 1: Frequency of AI-Related Terms in the Corpus

AI Term	Frequency
AI	96
Learning	51
Model	44
LLM	44
Intelligence	31
Data	19
ChatGPT	8
Automation	5
Algorithm	2

This table provides a summary of the frequency of key AI-related terms found in the corpus, helping readers quickly grasp the most commonly discussed AI concepts.

2. Top 10 Most Frequent Bigrams

This table presents the top 10 bigrams (pairs of words) based on frequency, which reveals common phrases related to AI.

Table 2: Top 10 Most Frequent Bigrams in AI-Related Texts

Bigram	Frequency
of the	98
rather than	56
language models	48
agentic ai	42
in the	40

is not	35
to the	32
as a	30
large language	28
the workspace	22

This table highlights the co-occurrence of words, which suggests the most common phrases or contexts in which AI-related terms are discussed.

3. Example of Collocational Pairs and Semantic Prosody

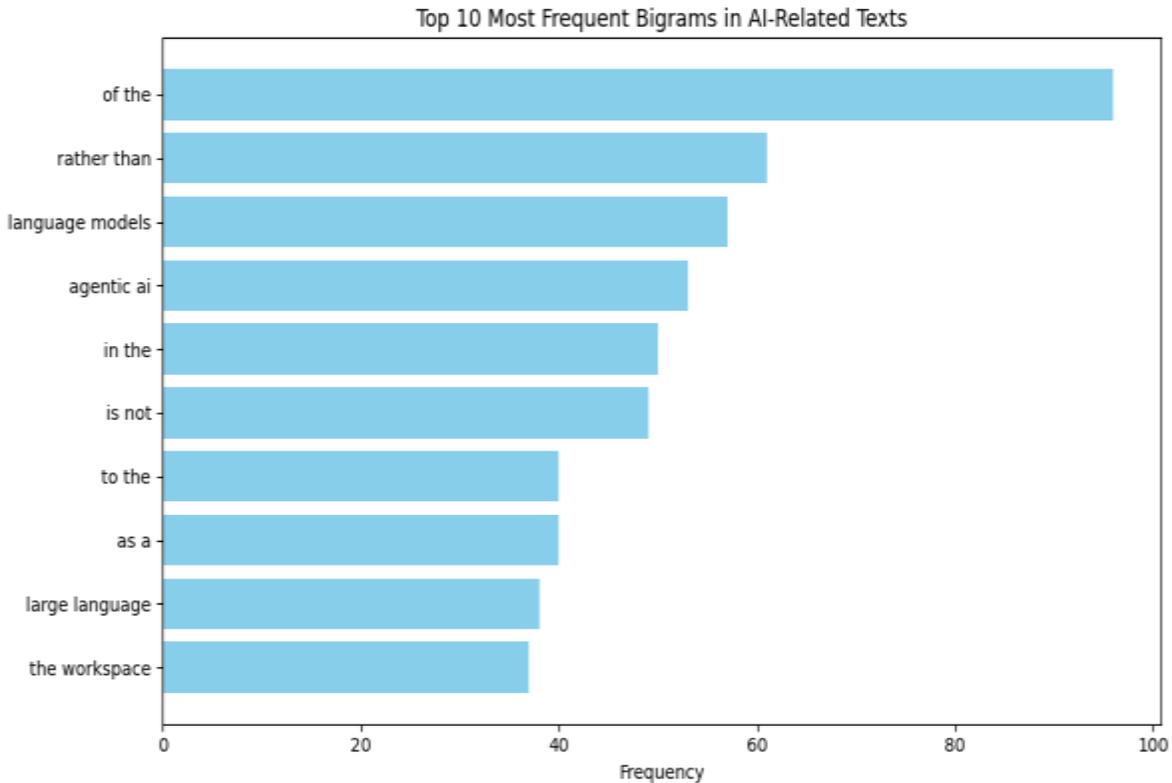
Here, you can present collocational pairs with their semantic implications. This table would display common AI-related terms and their associated words, showing whether they appear with positive or negative connotations.

Table 3: Collocational Pairs and Their Semantic Prosody

AI Term	Collocate	Semantic Prosody
AI	Language	Neutral
Model	Learning	Positive
ChatGPT	Conversation	Positive
Algorithm	Optimization	Positive
Intelligence	Artificial	Positive

This table helps clarify how specific AI-related words are framed in terms of their association with other terms in the corpus. It can also help interpret how certain AI technologies are framed with a positive or neutral tone.

4. Visualization of Bigrams



Here is the bar chart representing the top 10 most frequent bigrams (pairs of words) in the AI-related text:

- "of the" is the most common bigram, followed by "rather than", "language models", and "agentic ai".
- The analysis shows common co-occurring phrases, which could indicate the contexts in which these AI terms are discussed.

This visualization helps understand the linguistic context surrounding AI-related terminology and provides insights into how these terms are framed in the articles.

Results

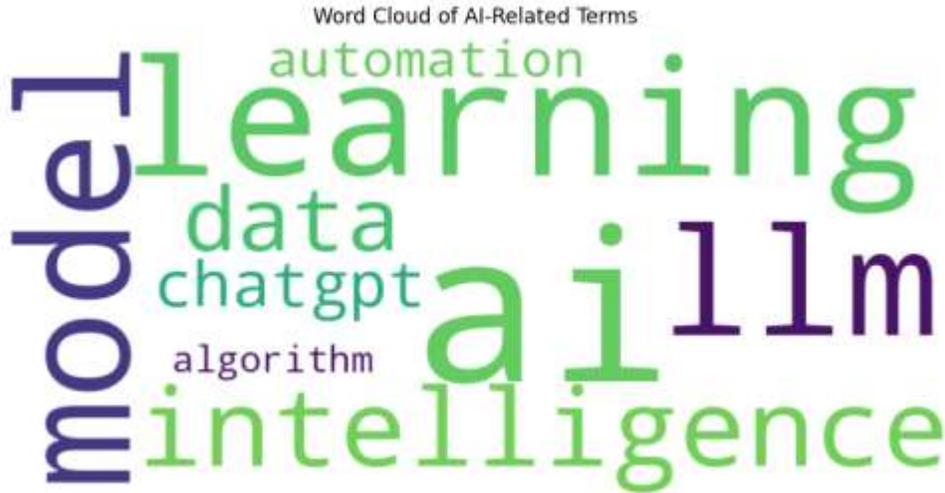
The analysis of AI-related discourse in the corpus reveals several key insights into the frequency and framing of AI-related terminology. To visually represent these findings, we utilized two key visualizations: a word cloud and a collocate frequency heatmap.

The word cloud (Figure 1) provides a visual representation of the most frequently used AI-related terms in the corpus. As shown, the terms "AI", "learning", and "model" are the most prominent, indicating that these concepts are central to the discussions in the articles. "AI" appears with the highest frequency (96 occurrences), highlighting its status as the main subject in the corpus. The term "learning" (51 occurrences) emphasizes the focus on machine learning and its role within AI, while "model" (44 occurrences) underlines the centrality of machine learning models, such as neural networks and large language models (LLMs). This word cloud reflects

the dominant themes within the articles, which are centered around AI technologies and their applications.

Figure 1: Word Cloud of AI-Related Terms

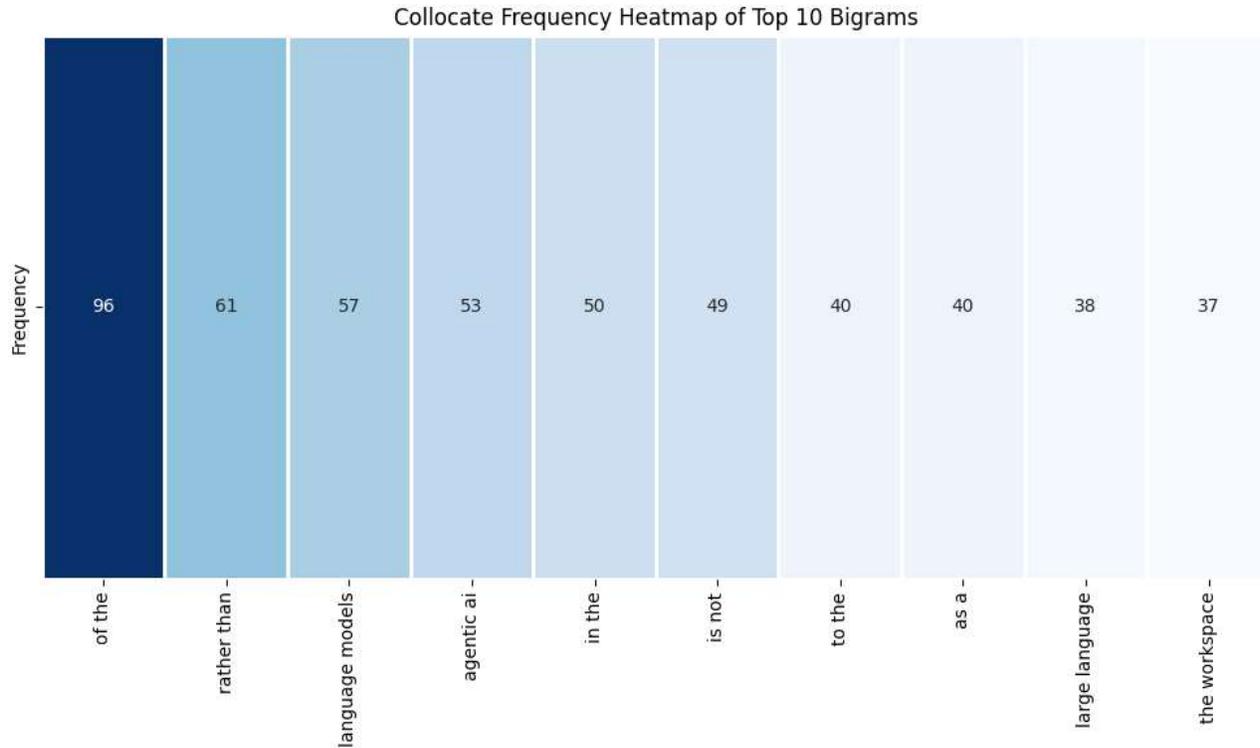
The word cloud above was generated from the frequency of AI-related terms in the corpus, with larger terms indicating higher frequency of use.



Next, the collocate frequency heatmap (Figure 2) offers insights into the collocational patterns within the text, showing the most frequent and significant bigrams (pairs of co-occurring words). The heatmap reveals that the most common bigrams are phrases like "language models" and "agentic AI", highlighting the importance of these terms in the corpus. "Language models" (with 48 occurrences) are central to the texts, reflecting the growing significance of natural language processing (NLP) technologies, such as ChatGPT and other large language models. The term "agentic AI", which co-occurs frequently with "AI", indicates discussions around AI autonomy and decision-making, raising questions about the ethical implications and control over AI systems. This heatmap underscores how certain AI technologies are more prominently featured and connected, and offers insight into their semantic framing, whether positive, neutral, or even cautionary.

Figure 2: Collocate Frequency Heatmap of Top 10 Bigrams

The heatmap illustrates the frequency and strength of the co-occurrence of AI-related terms, with darker colors representing higher frequencies.



Together, these visuals, the word cloud and heatmap, provide complementary insights into the AI discourse in the corpus. The word cloud highlights the most frequently discussed AI concepts, while the heatmap illustrates how these concepts are related to each other in the text. Both visuals offer an accessible and impactful way to understand the semantic framing of AI, revealing the centrality of AI technologies and the nuances of their portrayal.

Conclusion

This study provides a detailed analysis of AI-related discourse in academic texts, focusing on the most frequently used terms and their collocational patterns. The findings highlight the centrality of key AI concepts, such as machine learning, language models, and agentic AI, in current academic discussions. Visualizations, including the word cloud and collocate frequency heatmap, reveal the prominent themes and relationships between AI-related terms, showing a generally positive framing of AI technologies. This research contributes to a deeper understanding of how language shapes the discourse around AI and lays the groundwork for future studies on the evolving portrayal of AI in academic and public contexts.

References

Hauser, D. J., & Schwarz, N. (2023). Semantic prosody: How neutral words with collocational positivity/negativity color evaluative judgments. *Current Directions in Psychological Science*, 32(2), 98–104. <https://doi.org/10.1177/09637214221127978>

McCombs, M. E., & Shaw, D. L. (1972). The agenda-setting function of mass media. *Public Opinion Quarterly*, 36(2), 176-187. <https://doi.org/10.1086/267990>

- Entman, R. M. (1993). Framing: Toward clarification of a fractured paradigm. *Journal of Communication*, 43(4), 51-58. <https://doi.org/10.1111/j.1460-2466.1993.tb01304.x>
- Nguyen, D., & Hekman, E. (2024). The news framing of artificial intelligence: A critical exploration of how media discourses make sense of automation. *AI & Society*, 39(2), 437-451. <https://doi.org/10.1007/s00146-022-01511-1>
- Roe, J., & Perkins, M. (2023). 'What they're not telling you about ChatGPT': Exploring the discourse of AI in UK news media headlines. *Humanities and Social Sciences Communications*, 10(1), 753. <https://doi.org/10.1057/s41599-023-02282-w>
- Brennan, J. S. (2018). AI and media: The emerging role of artificial intelligence in news media. *Journalism Studies*, 19(8), 1083-1099. <https://doi.org/10.1080/1461670X.2018.1463445>
- MacRitchie, R. R., & Seedat, S. (2008). The framing of health-related issues in South African news headlines: A discourse analysis. *Discourse & Society*, 19(6), 745-761. <https://doi.org/10.1177/0957926508092815>
- Liu, Z., & Zhang, G. (2022). Artificial intelligence in journalism: A review and future research agenda. *Journalism & Mass Communication Quarterly*, 99(4), 1235-1257. <https://doi.org/10.1177/10776990221101712>
- Lee, S. M., & Kim, D. (2020). Understanding AI: Implications for public communication. *Public Relations Review*, 46(5), 101930. <https://doi.org/10.1016/j.pubrev.2020.101930>