

Brief Report

Not peer-reviewed version

---

# Advancements in Corpus Analysis Tools: A Comprehensive Guide for Linguistic Research and Novice Users

---

[Maryyam Jabeen](#)<sup>\*</sup> and Rabia Rashid

Posted Date: 10 February 2025

doi: 10.20944/preprints202502.0639.v1

Keywords: corpus; technical report; artificial intelligence; linguistics



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a Creative Commons CC BY 4.0 license, which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

*Brief Report*

# Advancements in Corpus Analysis Tools: A Comprehensive Guide for Linguistic Research and Novice Users

Maryyam Jabeen \* and Rabia Rashid

Department of English, University of Sahiwal, Pakistan; rabiarasheed589@gmail.com

\* Correspondence: jabeenm666@gmail.com

**Abstract:** Corpus analysis has become an essential methodology in linguistic research, providing critical insights into language use across large datasets. The development of various corpus analysis tools has greatly facilitated the exploration of textual data, allowing researchers to conduct sophisticated analyses of linguistic patterns and structures. This report offers a comprehensive overview of the most commonly used corpus analysis tools available today, with a focus on their functionality, applications, and suitability for beginner researchers. Additionally, this report discusses how these tools can provide crucial support to novice linguists in their corpus-based investigations, relieving them of the burden of complex tasks and offering user-friendly features. By introducing these tools and outlining their core capabilities, this report serves as a guide for those new to corpus analysis, providing reassurance and guidance, while also providing a foundation for further research and advanced linguistic exploration.

**Keywords:** corpus; technical report; artificial intelligence; linguistics

## I. Introduction

Corpus linguistics, the study of language as expressed in textual corpora, has experienced a significant surge, largely due to the advancement of computational tools. These tools have revolutionized linguistic research by enabling the systematic investigation of language patterns, often revealing insights that are not easily detectable through manual analysis (Jockers & Witten, 2010; Imran & Ain, 2019). The automated processing of vast amounts of text, facilitated by these tools, has uncovered everything from simple word frequencies to complex syntactic structures, further cementing their importance in linguistic research (Imran & Almusharraf, 2023, 2024).

Over the past decade, the field of corpus-based research has witnessed a proliferation of tools, catering to the needs of researchers at all levels (Imran et al., 2024a, 2024b). These tools, ranging from basic concordance software to advanced platforms for corpus creation and annotation, are now more accessible than ever. It is crucial to emphasize that these tools are within reach for beginner researchers, providing a clear guide on how to utilize them effectively (Perkins & Roe, 2024; Maqbool et al., 2024; Jabeen, 2023). This report offers a detailed examination of the most prominent corpus analysis tools, along with guidance for novices on how to integrate them into their research, thereby fostering a sense of inclusion and community among novice researchers.

Tools for Corpus Analysis:

A wide range of tools is available for corpus analysis, each offering unique features tailored to different research needs. The following sections describe key tools currently used in corpus linguistics, highlighting their specific capabilities and applications.

### 1. AntConc

AntConc, developed by Laurence Anthony, is a free and highly accessible tool designed for conducting basic corpus analyses. It is particularly suitable for beginners due to its intuitive interface and straightforward functionalities. AntConc's core features include:

*Concordance Analysis:* This tool allows users to examine the context in which specific words or phrases appear within a text. The concordance function generates Key Word in Context (KWIC) displays, enabling researchers to analyze how words are used in different linguistic environments.

*Word Frequency Analysis:* AntConc can generate word frequency lists, showing how often each word appears in a corpus. This is useful for identifying high-frequency terms and studying language usage patterns.

*Collocation Analysis:* The tool also allows for the analysis of collocations, or words that frequently occur together. This is particularly valuable for exploring lexical patterns and understanding the relationships between words.

*Keyword Analysis:* AntConc helps identify keywords by comparing the frequency of terms across different corpora. This allows researchers to pinpoint words that are significantly more common in one corpus compared to another.

Beginners benefit from AntConc's simplicity and extensive documentation. Numerous online tutorials and videos guide users through its various features, making it an ideal entry point for novice linguists (Anthony, 2020).

### 2. Sketch Engine

Sketch Engine is a powerful web-based tool for corpus analysis and corpus creation, widely used in professional linguistic research. Although it is a subscription-based platform, its range of features makes it worth the investment for advanced users. Key features of Sketch Engine include:

*Word Sketches:* These are one-page summaries of a word's grammatical and collocational behavior, providing a quick overview of how a word functions in different syntactic contexts. This feature is especially useful for lexicography and linguistic analysis.

*Corpus Building:* Sketch Engine allows users to create their own corpora by uploading text files or by using its integrated web-crawling feature to gather data from the internet.

*Keyword and Frequency Analysis:* Similar to AntConc, Sketch Engine can generate frequency lists and perform keyword analysis, but it also provides more sophisticated statistical measures.

*Thesaurus Creation:* Another powerful feature of Sketch Engine is its ability to automatically generate thesauruses for a given corpus, highlighting synonyms and related words based on actual language use in the corpus. While Sketch Engine offers more advanced features than AntConc, its user-friendly interface ensures that beginners can still navigate the software with relative ease after a brief learning period. This makes it a versatile tool for both novice and experienced researchers (Kilgariff et al., 2014).

### 3. WordSmith Tools

WordSmith Tools, developed by Mike Scott, is another prominent software package used for corpus analysis. It provides a range of text analysis tools that are highly useful for intermediate and advanced researchers, including:

*Concordance:* Similar to AntConc, WordSmith provides KWIC displays that allow researchers to explore how specific words are used in context.

*Frequency Lists:* The software can generate detailed word frequency lists, with options for further customization based on the user's needs.

*Keyness Analysis:* WordSmith helps identify key words in a text, defined as words that occur with unusual frequency compared to other corpora. This is useful for comparative studies across different datasets.

*Clusters and N-Grams:* This tool is also capable of analyzing word clusters and n-grams, which are groups of words that frequently co-occur. This functionality is valuable for identifying common phrases and examining language patterns. WordSmith's comprehensive set of tools makes it a preferred choice for advanced linguistic studies. However, its interface may present a steep learning curve for beginners. For novice researchers, it is advisable to start with simpler tools like AntConc before progressing to WordSmith Tools (Scott, 2010).

#### 4. Constituent Likelihood Automatic Word-tagging System (CLAWS)

CLAWS is a well-known tool for part-of-speech (PoS) tagging, developed at Lancaster University. It is designed to automatically tag words with their corresponding parts of speech, offering an accuracy rate of over 96%. CLAWS is particularly noted for its role in tagging the British National Corpus (BNC), one of the largest and most well-known corpora in English linguistics.

*PoS Tagging:* CLAWS assigns syntactic categories (such as nouns, verbs, adjectives) to each word in a text, which is essential for many types of linguistic analysis.

*Custom Tagsets:* The tool allows users to work with a variety of tagsets, making it adaptable to different research requirements.

*Batch Processing:* CLAWS is capable of processing large corpora quickly, making it an excellent choice for large-scale studies.

While CLAWS is a specialized tool, its high level of accuracy and efficiency make it a valuable resource for researchers working with tagged corpora (Garside, 1987).

#### 5. UAM CorpusTool

The UAM CorpusTool, developed by Mick O'Donnell, is designed for multilayer corpus annotation. It allows researchers to annotate texts for a variety of linguistic features, including syntax, semantics, and discourse structure. This tool is particularly useful for conducting fine-grained analyses that go beyond basic text statistics (Perkins & Roe, 2024).

*Multilayer Annotation:* Users can apply different layers of annotation to a single text, enabling the simultaneous analysis of syntax, discourse, and semantics.

*Visualization Features:* The UAM CorpusTool includes visualization options that help researchers interpret the results of their annotations. Graphs and charts offer a clear representation of complex linguistic data.

*Query Functions:* The tool allows users to perform complex searches based on the annotated features of the corpus, making it easier to identify patterns and relationships in the data. The UAM CorpusTool is ideal for more advanced researchers who need detailed annotations for their corpora. Beginners may find the tool challenging at first, but its extensive documentation and tutorials can help users get started (Perkins & Roe, 2024).

#### 6. BNCweb

BNCweb is a web-based interface for querying the British National Corpus (BNC), providing a user-friendly way to explore this extensive collection of English texts. Key features include:

*Corpus Access:* BNCweb allows users to search and retrieve data from the British National Corpus, one of the largest corpora available for the study of contemporary British English.

*Collocation Analysis:* The platform supports the analysis of word collocations within the corpus, enabling users to examine how words co-occur in various contexts.

*Frequency and Distribution Analysis:* BNCweb provides tools for analyzing word frequencies and their distribution across different genres and text types within the corpus. BNCweb is particularly useful for beginners, as it offers easy access to a pre-built, well-annotated corpus. Novice users can quickly get started with linguistic analysis without the need to build or compile their own corpus (Burnard & Aston, 1998).

## 7. WordNet

WordNet is a lexical database that groups English words into sets of synonyms (synsets), providing information about their meanings, relationships, and usage. While not strictly a corpus analysis tool, WordNet is frequently used in conjunction with corpus tools to provide semantic analysis.

*Semantic Relationships:* WordNet organizes words based on their meanings, offering semantic information such as synonyms, antonyms, and hyponyms.

*Integration with Other Tools:* WordNet can be integrated with tools like Sketch Engine to enhance semantic analysis of corpora by providing additional lexical resources. For beginners interested in the semantics of words, WordNet is an invaluable resource that can complement their corpus analysis work (Miller, 1995).

## II. Supporting Beginners in Corpus Analysis

For beginners, navigating the world of corpus analysis tools can seem daunting. However, with the right guidance and tools, novice researchers can quickly develop their skills in this area. The following sections highlight how these tools can assist beginners in their linguistic research:

### 1. Accessible Tools for Entry-Level Research

AntConc, BNCweb, and Sketch Engine are highly accessible tools for beginners. AntConc, in particular, is well-suited for those starting out in corpus analysis due to its simple interface and extensive online tutorials. BNCweb offers access to a pre-compiled corpus, allowing beginners to focus on analysis rather than corpus creation. Meanwhile, Sketch Engine, though more advanced, provides a comprehensive set of features that can gradually introduce novice users to more complex forms of analysis.

### 2. Step-by-Step Learning

For beginners, hands-on practice with user-friendly tools is critical. AntConc, for example, offers basic word frequency and concordance analysis, which are ideal for developing an understanding of corpus analysis principles. Sketch Engine provides tutorials and documentation that allow users to gradually explore its more sophisticated functions, such as thesaurus creation and collocation analysis.

### 3. Visualization and Interpretation

Visualization tools are especially important for beginners who are new to interpreting corpus data. Tools like Sketch Engine and UAM CorpusTool offer visualization options that display linguistic patterns in graphs and charts, making it easier for beginners to grasp complex relationships between words and phrases.

### 4. Extensive Documentation and Support

The majority of the tools described in this report come with extensive documentation, user guides, and tutorials. This is particularly valuable for beginners, as these resources provide step-by-step instructions on how to utilize each tool effectively. AntConc and Sketch Engine, for instance, offer video tutorials and forums where users can seek advice and share experiences.

## III. Conclusion

The diverse range of tools available for corpus analysis has transformed the way linguistic research is conducted. Tools like AntConc, Sketch Engine, WordSmith Tools, and UAM CorpusTool provide researchers with the ability to analyze vast amounts of text data quickly and efficiently. For



beginners, these tools offer accessible entry points into the field of corpus linguistics, allowing them to develop their skills in textual analysis while gradually expanding their methodological toolkit.

As corpus linguistics continues to evolve, mastering these tools will be essential for both novice and experienced researchers. The combination of user-friendly interfaces, comprehensive documentation, and powerful functionalities ensures that these tools will remain integral to linguistic research for years to come.

## References

- Anthony, L. (2020). AntConc (Version 3.4. 3)[Computer Software]. Waseda University, 2014. <https://doi.org/10.5281/zenodo.3818746>
- Burnard, L. (1998). *The BNC Handbook: Exploring the British National Corpus with SARA*. Edinburgh University Press. <https://doi.org/10.3366/edinburgh/9780748609901.001.0001>
- Garside, R. (1988). *The computational analysis of English: A corpus-based approach* (Vol. 57). G. Sampson, & G. Leech (Eds.). Longman. [https://doi.org/10.1007/978-1-349-19422-4\\_3](https://doi.org/10.1007/978-1-349-19422-4_3)
- Imran, M., & Ain, Q. (2019). Effects of non-native instructors' L1, beliefs and priorities on pronunciation pedagogy at secondary level in district Rajanpur. *Pakistan. Journal of Language and Cultural Education*, 7(2), 108-121.
- Imran, M., & Almusharraf, N. (2023). A review on the development of education in the post-COVID-19 era: Teaching in the post-COVID-19 era: World education dilemmas, teaching innovations and solutions in the age of crisis, edited by Ismail Fayed and Jill Cummings. *Social Identities*, 29(3), 323–325. <https://doi.org/10.1080/13504630.2023.2227569>
- Imran, M., & Almusharraf, N. (2024). Digital Learning Demand and Applicability of Quality 4.0 for Future Education: A Systematic Review. *International Journal of Engineering Pedagogy*, 14(4).
- Jabeen, M. (2023). Exploring ChatGPT's Role in Creative Writing: A Short Review. *International Review of Literary Studies*, 5(2), 32-34.
- Jockers, M. L., & Witten, D. M. (2010). A comparative study of machine learning methods for authorship attribution. *Literary and Linguistic Computing*, 25(2), 215-223.
- Kilgarrriff, A., Baisa, V., Bušta, J., Jakubíček, M., Kovář, V., Michelfeit, J., ... & Suchomel, V. (2014). The Sketch Engine: ten years on. *Lexicography*, 1(1), 7-36. <https://doi.org/10.1007/s40607-014-0009-9>
- Maqbool, M. A., Asif, M., Imran, M., Bibi, S., & Almusharraf, N. (2024). Emerging e-learning trends: a study of faculty perceptions and impact of collaborative techniques using fuzzy interface system. *Social Sciences & Humanities Open*, 10, 101035.
- Miller, G. A. (1995). WordNet: a lexical database for English. *Communications of the ACM*, 38(11), 39-41. <https://doi.org/10.1145/219717.219748>
- Perkins, M., & Roe, J. (2024). Decoding academic integrity policies: A corpus linguistics investigation of AI and other technological threats. *Higher Education Policy*, 37(3), 633-653.
- Scott, M. (2010). WordSmith Tools (Version 6.0) [Computer software]. Lexical Analysis Software. <https://doi.org/10.1093/lc/fqq001>
- Imran, M., Almusharraf, N., Ahmed, S., & Mansoor, M. I. (2024a). Personalization of E-Learning: Future Trends, Opportunities, and Challenges. *International Journal of Interactive Mobile Technologies*, 18(10).
- Imran, M., Almusharraf, N., Abdellatif, M. S., & Ghaffar, A. (2024b). Teachers' perspectives on effective English language teaching practices at the elementary level: A phenomenological study. *Heliyon*, 10(8).

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.