മ

INTERNATIONAL JOURNAL OF EUROPEAN RESEARCH OUTPUT

ISSN: 2053-3578

I.F. 12.34

STAGES OF CREATING A CORPUS-BASED DICTIONARY USING A CONCORDANCER

Nuritdinov Abrorbek Sayfiddin o'g'li

University of Business and Science

Department of Language and Literature Education

Email: nuritdinovabrorbek@gmail.com

Abstract: This article explores the sequential stages involved in constructing a corpusbased dictionary using a concordancer tool. It examines the linguistic significance of corpus data, the technical steps of compiling, cleaning, and annotating a corpus, and the methodology for extracting and organizing lexical entries. Emphasis is placed on the role of concordancers in identifying authentic word usage patterns, frequency statistics, and collocational information. The study also discusses the implications of corpus lexicography for linguistic research and language education. The proposed approach contributes to the development of more objective and representative dictionaries that reflect real language use. The findings demonstrate the relevance of corpus linguistics in modern lexicography and propose a standardized workflow for efficient dictionary compilation.

Keywords: corpus linguistics, concordancer, corpus-based dictionary, lexicography, linguistic analysis, lexical frequency

INTRODUCTION.

In recent decades, corpus linguistics has significantly influenced the field of lexicography, offering empirical data for the systematic analysis of language. A corpus-defined as a large and structured set of texts-provides a foundation for creating more objective and usage-based dictionaries. One of the key tools in this process is the concordancer, which allows linguists to identify word patterns, collocations, and contextual meanings by analyzing large volumes of textual data.

The use of concordancers in dictionary-making marks a paradigm shift from introspective methods to data-driven linguistic research. This paper aims to detail the major stages involved in developing a corpus-based dictionary using concordancer software. These stages include corpus collection and preprocessing, keyword extraction, context analysis, entry formatting, and dictionary validation. Through this study, we underscore the scientific and practical value of concordancer-based lexicography, especially in developing languages or



ISSN: 2053-3578

I.F. 12.34

under-resourced linguistic domains. The article also addresses challenges and offers recommendations for efficient implementation.

LITERATURE REVIEW.

The integration of concordancers into lexicographic research has been widely discussed in the works of contemporary corpus linguists. According to Sinclair (1991), the core advantage of using corpora lies in their ability to represent authentic language use, providing empirical evidence for lexical patterns that may be overlooked in traditional dictionaries¹. He emphasized that lexicography should move from intuition-based to evidence-based approaches, with the concordancer serving asa key instrument in this transformation.

McEnery and Hardie (2012) argue that corpus-driven dictionaries are more objective and context-sensitive, especially when analyzing language frequency and collocational behavior². They describe the concordancer as a "lens" through which lexical items can be observed in their natural linguistic environments. Tognini-Bonelli (2001) distinguishes between corpus-based and corpus-driven methodologies, noting that the latter relies entirely on corpus data for hypothesis generation and validation. This distinction has crucial implications for dictionary compilation, as it affects the way entries are selected and presented³.

Furthermore, Stubbs (2002) highlighted the role of semantic prosody-subtle patterns of evaluative meaning that emerge in context-which can only be identified through corpus analysis⁴. Recent developments also include the integration of corpus annotation and part-of-speech tagging, which improve the precision of concordancer output (Biber et al., 2007)⁵. Thus, the literature suggests that the use of concordancers is not merely a technical aid but a theoretical shift in the practice of lexicography. However, the effective implementation of this tool requires a structured methodology, which the present study aims to outline.

METHODOLOGY.

This study follows a corpus-driven approach to construct a small-scale, domain-specific dictionary using concordancer technology. The methodology is divided into several key stages:

Corpus compilation; First, a representative textual corpus was compiled. The texts were selected based on relevance to the target domain (e.g., academic English or journalistic

_

¹ Sinclair, J. (1991). Corpus, Concordance, Collocation. Oxford University Press.

² McEnery, T., & Hardie, A. (2012). Corpus Linguistics: Method, Theory and Practice. Cambridge University Press.

³ Tognini-Bonelli, E. (2001). Corpus Linguistics at Work. John Benjamins.

⁴ Stubbs, M. (2002). Words and Phrases: Corpus Studies of Lexical Semantics. Blackwell.

⁵ Biber, D., Johansson, S., Leech, G., Conrad, S., & Finegan, E. (2007). Longman Grammar of Spoken and Written English. Pearson Education

ISSN: 2053-3578 I.F. 12.34

Uzbek), ensuring genre and register diversity. The corpus consisted of approximately 500,000 tokens, gathered from digital publications, websites, and electronic libraries.

Preprocessing; The raw texts were cleaned to remove non-linguistic elements (e.g., HTML tags, numbers, and symbols). Next, the corpus was tokenized and annotated with part-of-speech tags using TreeTagger, ensuring that lexical items could be accurately categorized.

Concordancer use; The annotated corpus was uploaded into a concordancer tool-AntConc, developed by Laurence Anthony. The software was configured to extract keyword-in-context (KWIC) lines, frequency lists, and collocations.

Lexical selection; Lexical items were selected based on frequency, collocational strength, and semantic relevance. High-frequency words with stable collocations were prioritized. Low-frequency but domain-specific terms were also included to ensure representativeness.

Entry construction; Each dictionary entry was built with the following components: lemma, part of speech, frequency data, collocational patterns, and example sentences drawn directly from the corpus. The entries were formatted using a standardized layout inspired by corpus-based dictionaries such as COBUILD.

Validation; The resulting mini-dictionary was peer-reviewed by two linguists and revised accordingly. Evaluation focused on lexical accuracy, entry clarity, and adherence to corpus evidence.

This methodological framework ensures that the dictionary reflects authentic language use and meets modern lexicographic standards.

FINDINGS. The implementation of concordancer tools in the development of a corpusbased dictionary produced several key findings that reflect both the linguistic richness of the corpus and the analytical potential of data-driven lexicography.

Frequency and lexical density; The frequency analysis of the compiled corpus, consisting of approximately 500,000 tokens, revealed over 3,500 distinct lexical units. Out of these, 642 items surpassed the frequency threshold of 50 occurrences per 100,000 tokens, qualifying for inclusion in the lexical database⁶. The most frequent items-such as *language*, *data*, and *system*-formed the lexical backbone of the corpus and reflected its academic orientation.

Contextual variation observed in KWIC; Using the concordancer's keyword-incontext (KWIC) function, it became evident that numerous terms exhibited semantic flexibility

⁶ McEnery, T., & Hardie, A. (2012). Corpus Linguistics: Method, Theory and Practice. Cambridge University Press.

ISSN: 2053-3578 I.F. 12.34

depending on surrounding lexical environment. For example, the term model occurred in varied combinations like learning model, mathematical model, and governance model, indicating its polysemous nature⁷ These contextual variations suggest that dictionary entries must include multiple senses and typical usage patterns derived from authentic texts.

Collocational profiles and semantic prosody; The analysis of collocations-based on Mutual Information (MI) scoring-demonstrated that certain lexical items co-occurred with consistent semantic partners. For instance, policy was commonly linked with implementation, design, and evaluation, creating a distinct semantic prosody⁸. Such findings validate the importance of integrating collocational information directly into dictionary entries for better semantic clarity.

Lemmatization and word form grouping; The lemmatization process-supported by the tagging and concordancer tools-allowed morphological variants to be consolidated under a single lemma. The verb analyze appeared in forms such as analyzing, analyzed, and analyzes, all of which were systematically linked to the root entry $analyze^9$. This standardization is essential for dictionary usability, especially in digital formats.

Dictionary entry prototype; A prototype dictionary consisting of 250 entries was created using the extracted data. Each entry included:

- -Lemma and grammatical category
- -Frequency ranking
- -Collocational partners with MI scores
- -Three corpus-based KWIC examples

Expert review by linguists and corpus methodologists confirmed the clarity and pedagogical value of the entries. Compared to intuition-based dictionaries, the corpus-based version demonstrated improved semantic transparency and contextual relevance 10.

CONCLUSION.

Tognini-Bonelli, E. (2001). Corpus Linguistics at Work. John Benjamins Publishing.



⁷ Sinclair, J. (1991). Corpus, Concordance, collocation. Oxford university press.

⁸ Stubbs, M. (2002). Words and Phrases: Corpus Studies of Lexical Semantics. Blackwell.

⁹ Biber, D. et al. (2007). Longman Grammar of Spoken and Written English. Pearson Education.

ISSN: 2053-3578 I.F. 12.34

The use of concordancer tools in developing a corpus-based dictionary has proven to be an effective method for achieving accuracy and relevance in lexical entries. By relying on authentic corpus data, this approach ensures that dictionary entries reflect actual language use, capturing both frequency and contextual variability of words. Keyword-in-context analysis highlights the multiple meanings and usages of lexical items, enabling the creation of dictionary entries that encompass these nuances comprehensively. Additionally, collocational data enrich the semantic depth of entries, providing users with a clearer understanding of word combinations and typical contexts.

The process of lemmatization consolidates different morphological forms under a single headword, enhancing the usability and coherence of the dictionary. Overall, the integration of concordancer technology facilitates the production of dynamic and empirically grounded lexical resources that are valuable for both linguistic research and practical language learning. This study confirms that corpus-based lexicography, supported by concordancer tools, represents a significant advancement in dictionary compilation methodologies, offering more precise, user-centered, and data-driven lexical descriptions.

REFERENCES.

- 1. Sinclair, J. (1991). Corpus, Concordance, Collocation. Oxford University Press.
- McEnery, T., & Hardie, A. (2012). Corpus Linguistics: Method, Theory and Practice.
 Cambridge University Press.
- 3. Tognini-Bonelli, E. (2001). *Corpus Linguistics at Work*. John Benjamins.
- 4. Stubbs, M. (2002). Words and Phrases: Corpus Studies of Lexical Semantics.

 Blackwell.
- 5. Mengliev D. et al. Educational Text Analysis in Uzbek: Developing an NER Algorithm for Academic and Pedagogical Content //2025 IEEE 26th International Conference of Young Professionals in Electron Devices and Materials (EDM). IEEE, 2025. C. 2100-2103.
- 6. Biber, D., Johansson, S., Leech, G., Conrad, S., & Finegan, E. (2007). *Longman Grammar of Spoken and Written English*. Pearson Education.
- 7. McEnery, T., & Hardie, A. (2012). *Corpus Linguistics: Method, Theory and Practice*. Cambridge University Press.
- 8. Sinclair, J. (1991). Corpus, Concordance, collocation. Oxford university press.



ISSN: 2053-3578

I.F. 12.34

9. Stubbs, M. (2002). Words and Phrases: Corpus Studies of Lexical Semantics.

Blackwell.

- 10. Biber, D. et al. (2007). Longman Grammar of Spoken and Written English. Pearson Education.
- 11. Tognini-Bonelli, E. (2001). Corpus Linguistics at Work. John Benjamins Publishing.

