

CORPUS LINGUISTICS: A HISTORICAL OVERVIEW

Sukhrob Avezov Sobirovich

Associate Professor of the Department of Russian

Language and Literature Bukhara State University

senigama1990@mail.ru

Abstract

Corpus linguistics is a field of study that investigates language through large collections of texts (corpora). This article provides a historical overview of corpus linguistics, tracing its development from early textual analyses in pre-digital eras to the emergence of computerized corpora in the mid-20th century and the subsequent growth of the field. Key milestones such as the first computerized corpora in the 1960s, the expansion of corpus resources in the 1980s-1990s (the Brown Corpus, British National Corpus), and the integration of corpus-based methods into linguistics and language technology are discussed.

Keywords: Corpus linguistics, history of corpus linguistics, Brown Corpus, British National Corpus, text corpora, computational linguistics, lexicography, language analysis.

Introduction

Corpus linguistics is an approach to linguistic inquiry that relies on the empirical analysis of large collections of texts (corpora) to study language use. The approach contrasts with purely introspective methods by emphasizing attested usage of language. Although the term corpus linguistics gained currency in the late 20th century, the practice of analyzing language based on textual corpora has deep historical roots. This paper offers a historical overview of corpus linguistics, charting its trajectory from early attempts at compiling and analyzing text collections to the establishment of a distinct field. We review major milestones – from early manually compiled corpora to the first computerized corpora in the 1960s, the rapid expansion of corpora in the late 20th century, and current trends – to illustrate how technological advancements and theoretical shifts shaped the field's development. The goal is to understand how corpus linguistics evolved into a robust discipline and how it transformed linguistic research.

Methods

This study employs a historical-descriptive methodology. We conducted a literature review of seminal works and historical analyses in corpus linguistics to gather information on the field's development. Key publications and corpus documentation were examined to identify pivotal events, such as the creation of notable corpora and the introduction of corpus-based methods. The approach is qualitative, synthesizing information from published histories, research articles, and corpus project reports. By drawing on documented milestones and the reflections of prominent linguists, we reconstruct the timeline of corpus linguistics and highlight the contributions that have been most influential in its evolution. This method allows for a comprehensive overview of the field's history based on existing scholarly sources, without new experimental data.

Results and Discussion

Early pre-computer traditions: the concept of analyzing language by examining collections of texts predates modern linguistics. Long before computers, scholars compiled corpora of important texts for linguistic and literary analysis. In ancient times, for example, Indian linguist Pāṇini (4th century BCE) based his Sanskrit grammar in part on a corpus of Vedic texts, and early Arabic grammarians meticulously studied the language of the Quran. In medieval Europe, monks and scholars prepared concordances of the Bible – essentially manual corpora of religious texts – to analyze word usage in scripture. These efforts show that the corpus-based approach to language description has deep roots in philology and lexicography. Similarly, by the 18th and 19th centuries, lexicographers like Samuel Johnson and later the Oxford English Dictionary editors gathered vast collections of written examples (citations) to document word meanings. These pre-computer endeavors established the principle that insights about language could be gained by observing large samples of authentic text, even though they lacked today's computational tools.

Birth of modern corpus linguistics (1950s–1960s): the mid-20th century saw the first attempts to systematically compile machine-readable corpora. Early forerunners included specialized projects such as the *Trésor de la Langue Française* (TLF) in France, which began converting a large French text collection into electronic form in the late 1950s. Around the same time, in 1959, a team at RAND Corporation in the United States created a small electronic corpus for

experiments in machine translation. However, the true landmark was the creation of the *Brown Corpus* in the 1960s. Henry Kučera and W. Nelson Francis at Brown University compiled this corpus – one million words of American English from 1961, sampled from 500 texts across diverse genres – and published its analysis in *Computational Analysis of Present-Day American English* (1967). The Brown Corpus was the first computerized general corpus designed for linguistic research. It set new standards with its balanced design and availability to other researchers, effectively launching modern corpus linguistics. Notably, British linguist Randolph Quirk simultaneously began the *Survey of English Usage* (SEU) in London in 1959–60, aiming to collect a broad corpus of British English. Quirk’s project, while initially paper-based, was conceived with future computer analysis in mind and is considered the first modern corpus intended to represent an entire language variety. Early reception of these efforts was mixed. At the time, generative grammar dominated linguistics and some leading theorists were skeptical of corpus-based approaches. W. Nelson Francis later recalled that a prominent generativist derided the Brown Corpus project as “a useless and foolhardy enterprise” on the grounds that “the only legitimate source of grammatical knowledge was the intuitions of the native speaker, which could not be obtained from a corpus”. Despite such criticisms, the success of the Brown Corpus demonstrated the value of empirical data, and it opened a new era by making empirical language evidence widely available for analysis.

Expansion in the 1970s–1980s: following the Brown Corpus, corpus linguistics gradually gained momentum. The design of Brown inspired a series of *Brown-family corpora* – comparable 1-million-word corpora for other varieties of English. Notably, the *LOB Corpus* (Lancaster-Oslo/Bergen Corpus of British English) was compiled in the 1970s as a British English counterpart to Brown. Similar corpora were created for other English varieties (Kolhapur Corpus for Indian English, Wellington Corpus for New Zealand English) and for different time periods (Frown and FLOB corpora in the early 1990s, updating Brown and LOB to later decades). In tandem, spoken language corpora emerged. One early example was the *London–Lund Corpus* of spoken British English, derived from Quirk’s SEU recordings in the 1970s. The first computerized spoken corpus, however, was compiled in 1971 in Canada: the Montreal French Corpus of one million words of transcribed Quebec French speech. These projects demonstrated that corpus methods could be extended beyond written texts to spoken language.

During the 1980s, the use of corpora spread to lexicography and grammar writing. Publishers began to exploit corpora to produce better dictionaries. For instance, the *American Heritage Dictionary* (1969) was an early dictionary to use a corpus (providing a citation base from the Brown Corpus). In the UK, John Sinclair led the innovative *Collins COBUILD* project, which built the *Bank of English* corpus and produced a revolutionary learners' dictionary (1987) based entirely on corpus evidence. Corpus-based research also informed grammar reference works: Quirk et al.'s *Comprehensive Grammar of the English Language* (1985) was one of the first large grammars to incorporate findings from corpus data (using the SEU corpus). By the late 1980s, the term "corpus linguistics" was coming into use as scholars organized conferences and publications around this approach. What had begun as a few pioneering projects had grown into a small but distinct community of researchers dedicated to empirical language study.

Large-scale corpora and acceptance (1990s): the 1990s witnessed an explosion in both the size and number of corpora, along with mainstream acceptance of corpus linguistics. Technological advances in computing and storage enabled the construction of much larger corpora. The *British National Corpus (BNC)*, completed in 1994, contained 100 million words of British English (spoken and written), far exceeding earlier corpora. It was created by a consortium of universities, publishers, and the British Library, reflecting broad institutional support for corpus-based research. In the United States, a project to compile an American National Corpus was initiated, and although it stalled, other American English resources grew – notably the *Corpus of Contemporary American English (COCA)*, which by 1990s end had hundreds of millions of words available via a web interface. Crucially, the 1990s also saw corpus linguistics permeate *computational linguistics* and natural language processing (NLP). The rise of statistical methods in NLP – for tasks like speech recognition and machine translation – was fueled by the availability of large corpora. For example, IBM's famous experiments in statistical machine translation leveraged large *parallel corpora* (bilingual text corpora) such as Canadian parliamentary proceedings and European Union multilingual texts. This data-driven paradigm shift in NLP demonstrated that having vast corpora of real language could dramatically improve algorithmic performance. Additionally, corpora in many languages blossomed during this period, aided by international collaborations. The *International Corpus of English (ICE)* project, launched in the 1990s, built

comparable corpora for numerous varieties of English worldwide. Other languages established their first national corpora, and multilingual corpora became valuable resources for contrastive linguistics and translation studies. Another significant development was the enrichment of corpora with linguistic annotations. The *Penn Treebank* project (Marcus et al., 1993), for instance, provided a corpus of American English text that was *parsed* (syntactically annotated), enabling detailed grammatical studies and training data for parsing algorithms. By the end of the 1990s, corpus linguistics had moved from the periphery of linguistics to the mainstream. The once sharp divide between corpus-based empiricism and theoretical linguistics (epitomized by generative grammar) had softened, as many linguists recognized the value of combining approaches. Dedicated conferences (ICAME) and journals (International Journal of Corpus Linguistics, founded 1995) solidified the field's academic presence.

New millennium: bigger and richer corpora (2000s): in the 2000s, corpus linguistics benefited from the digital revolution and the advent of the internet as a data source. Corpora grew exponentially in size. One watershed moment came in 2006 when Google released a trillion-word corpus of English derived from web pages (the Google Web 1T corpus), providing frequency data for n-grams (word sequences up to length 5) across a teraword of text. This unprecedented resource, though comprised of raw web text with noise and errors, illustrated the scale at which corpus data had expanded – roughly one million times the size of the Brown Corpus. Such extremely large corpora enabled new research in lexical statistics and prompted discussions about the “unreasonable effectiveness of data”, suggesting that sheer quantity of linguistic data can reveal even rare phenomena. Alongside size, corpora also became more richly annotated and specialized. For example, the open-source *Wikipedia Corpus* and various *monitor corpora* (which are continually updated with new texts, such as news articles) allowed researchers to track linguistic changes in real time. The 2000s also saw corpora being built for *specialized domains* and languages that were previously under-resourced. Corpus linguists compiled extensive databases for languages in Asia, Africa, and other regions, as well as for varieties of English (Corpus of Global Web-Based English covering 20 countries). Spoken corpora expanded with audio and video recordings aligned to transcripts. Even sign languages began to have video corpora compiled for linguistic analysis. In computational applications, corpora of user-generated content (like social media corpora)

emerged, feeding into research on language variation and change. The synergy between corpus linguistics and computational linguistics grew stronger: corpora became the training ground for data-hungry models in machine learning and later, in the 2010s, for neural language models. By the end of the 2000s, corpus linguistics had firmly established itself as an essential approach in both descriptive and applied linguistics.

Contemporary developments and impact: in the 2010s and beyond, corpus linguistics continues to flourish and diversify. The field now underpins a wide range of linguistic sub-disciplines and interdisciplinary studies. In *sociolinguistics* and *discourse analysis*, researchers use corpora to examine language usage across social contexts and media. In *historical linguistics*, large diachronic corpora (the Historical Thesaurus of English, or various corpora of earlier stages of languages) allow quantitative tracking of language change. New areas such as *corpus-based language teaching* (data-driven learning) have leveraged corpora to create teaching materials and help learners observe authentic usage patterns. Furthermore, corpus methods have been adopted in domains outside traditional linguistics. A striking example is the emergence of *Law and Corpus Linguistics*, where legal scholars analyze large collections of legal texts or general language corpora to interpret the meaning of words in statutes and constitutions. This approach has even been cited in court cases, demonstrating corpus linguistics' broad influence. Technological improvements have made corpora more accessible than ever: user-friendly interfaces and powerful query tools (like Sketch Engine, LancsBox, or online concordancers) enable researchers and the public to tap into corpora of billions of words. Today, corpora are available for hundreds of languages, and many are annotated with multiple layers of information (syntax, semantics, pragmatics). The line between corpus linguistics and computational language processing has further blurred with the rise of big data and natural language processing techniques; for instance, the huge text datasets used to train AI language models (such as GPT) can be seen as ultra-large corpora, and insights from corpus linguistics inform how such data are sampled and cleaned. Corpus linguistics is now a well-integrated part of the linguist's toolkit, and its empirical insights have encouraged a more evidence-based approach to studying language. From a marginal, sometimes contested practice in the 1960s, it has evolved into a mainstream paradigm that has fundamentally changed how we document and analyze languages.

Conclusion

The history of corpus linguistics showcases a remarkable evolution from humble beginnings to a central role in modern language research. Starting with scholars who manually compiled examples from revered texts, the corpus approach was revolutionized by the advent of computers, as exemplified by the Brown Corpus in 1961. Over subsequent decades, corpus linguistics expanded in scope and scale – embracing new languages, spoken data, and enormous text collections. It moved from the fringe of linguistic theory to a position of broad acceptance, influencing not only linguistics but also fields like lexicography, education, translation, and computational language technology. The milestones discussed – from early concordances and the first million-word corpora to today’s billion-word databases – highlight how advances in technology and methodology went hand in hand with changes in linguistic thought. Corpus linguistics has proven that analyzing real-world language data can yield insights into usage, variation, and structure that introspection alone cannot provide. As we look ahead, the continuing growth of digital text and multimedia ensures that corpus-based inquiry will remain at the forefront of linguistic research, driving further innovation in our understanding of language. In summary, what began as an unconventional idea of gathering “mere data” has become an indispensable scientific approach, fundamentally enriching the study of language.

References

1. Kučera H., Francis W. N. Computational Analysis of Present-Day American English. Providence: Brown University Press, 1967. 424 p.
2. Quirk R. Towards a description of English Usage // Transactions of the Philological Society. 1960. 59 (1): 40–61.
3. Leech G. N. The state of the art in corpus linguistics // English Corpus Linguistics: Studies in Honour of Jan Svartvik / Ed. by K. Aijmer, B. Altenberg. London: Longman, 1991. P. 8–29.
4. Sinclair J. M. Corpus, Concordance, Collocation. Oxford: Oxford University Press, 1991. 179 p.
5. Biber D., Conrad S., Reppen R. Corpus Linguistics: Investigating Language Structure and Use. Cambridge: Cambridge University Press, 1998. 312 p.
6. McEnery T., Hardie A. Corpus Linguistics: Method, Theory and Practice. Cambridge: Cambridge University Press, 2012. 280 p.



7. Авезов С. КОРПУСНАЯ ЛИНГВИСТИКА: НОВЫЕ ПОДХОДЫ К АНАЛИЗУ ЯЗЫКА И ИХ ПРИЛОЖЕНИЯ В ОБУЧЕНИИ ИНОСТРАННЫМ ЯЗЫКАМ //International Bulletin of Applied Science and Technology. – 2023. – Т. 3. – №. 7. – С. 177-181.
8. Sobirovich S. A. A PRAGMATICALLY ORIENTED APPROACH TO GENERATIVE LINGUISTICS //CURRENT RESEARCH JOURNAL OF PHILOLOGICAL SCIENCES. – 2024. – Т. 5. – №. 04. – С. 69-75.