



# Grammatical Structures in Formal Mathematical Writing: A Corpus-Based Analysis of Proof Exposition

<sup>1</sup>Dhanshri Sharma and <sup>\*2</sup>Dr. Manisha N Rathod

<sup>1, \*2</sup>Krishna School of Engineering and Technology, Drs. Kiran and Pallavi Patel Global University, Vadodara, Gujarat, India.

## Abstract

Mathematical proofs are hybrid objects, combining formal logic with natural language, but their grammatical structure is not studied empirically. This paper examines a corpus of 1,500 proofs passages based on research articles and advanced textbooks annotated with the type of clause, voice, and discourse markers and analyzed through quantitative techniques. The results indicate that conditional clauses represent 62 percent of the structures in assumption stages but reduce to 18 percent in derivations, whereas declarative clauses represent most of derivations (74 percent) and conclusions (82 percent). Clauses involving passive constructions (64% in derivations 28% in assumptions, and 35% in conclusions) represent systematic backgrounding of agency. Logical connectors are also highly clumped, three elements taking 85 percent of transitions between proofs. This suggests that grammatical options are highly coordinated with logical functions and that writing of proofs is conducted based on a strictly restricted and extremely regularized system of linguistic patterns.

**Keywords:** Mathematical Discourse, Corpus Analysis, Proof Exposition, Grammar, Academic Writing.

## Introduction

Mathematic proofs are usually considered to be the products of formal logic and are frequently seen to be symbolic rather than concrete. Practically, though, constructions of proofs are made and expressed in natural language, in which grammatical organization is paramount in arranging logical interdependences and in taking the reader through the progressive stages of thought. Nonetheless, there has been little empirical work on the linguistic aspect of mathematical writing, and most research on mathematics writing has been dominated by symbolic logic rather than the grammatical processes underlying exposition. Academic discourse corpus-based studies have always revealed that language use is extremely organized and discipline-specific, and that there are recurrent lexico-grammatical patterns to functional requirements of communication (Flowerdew, 2003; Conrad, 2018) [3, 9]. In mathematical circles, in recent years, the existence of recurring lexical bundles and phraseological regularities have been emphasized, indicating that writing in the field of mathematics is extremely conventionalized (Alasmay, 2022; Steidlová, 2022) [4, 5]. Nonetheless, it is still unclear how grammatical forms, especially types of clauses, voice, and cohesives, systematically capture a development of logical argumentation in proofs. Although corpus-based methodologies have been effectively used in technical and academic texts (Anthony, 2013; Abdelreheim, 2014; Prado-Alonso, 2019) [1, 8, 6], the use of corpus-based methods in the field of formal proof exposition is not well developed. This

gap is fulfilled by the current study, which uses a corpus-based analytical framework that consists of applying grammatical annotation and quantitative analysis (such as frequency distributions and cross-tabulation across proof stages, and statistical testing of correlations between linguistic features and logical functions). The research will focus on the demonstration of the correspondence between particular grammatical patterns and the particular stages of reasoning and help in the coherence and conventionalization of writing mathematical proofs by analyzing a structured corpus of mathematical proofs. Literature Review Corpus-based methods have significantly contributed to the knowledge of grammatical and discourse patterns in academic writing, showing that aspects of linguistic features are systematically influenced by disciplinary norms. The early corpus linguistics studies highlight the need to rely on empirical study in discovering repeated leico-grammatical patterns and how they are applied to fulfill a functional purpose in communication (Weisser, 2015; Flowerdew, 2015) [21, 18]. Research in all areas of academia indicates that the decisions we make about grammar are not random and are in fact responses to underlying rhetorical and epistemic needs. The study of syntactic complexity has had a significant impact especially on understanding the encoding of knowledge in academic writing. Thi *et al.* (2023) [10] review the syntactic complexity and error patterns of writer learners and demonstrate that the high level of structural sophistication is associated with disciplinary competence. Likewise, Karakaya (2017) [12]

emphasizes the nominal alteration and the use of complicated syntactic patterns to create dense information content, which is the key attribute of academic prose. Holtz, (2011) [11] also shows that scientific writing has unique lexico-grammatical features, especially in abstracts and research articles where brevity and accuracy are attained by certain types of grammatical structures. A different research topic has been formulaic language and repetitive sequences. According to Almosa (2024) [14], formulaic sequences are crucial components of academic discourse, which help to build fluency and coherence. Omarova *et al.* (2025) [15] demonstrate that frequency-based analysis demonstrates prevailing patterns of lexica that define particular registers. These results are supported by research on lexical bundles and genre variation, including Siriganjanavong (2019) [17], who reveals the differences in linguistic patterning between amateur and skilled writers, and the importance of conventionality in academic writing. Corpus methods have also been used to study grammaticalization and modality. Krug (2000) [16] follows the history of modal verbs, showing how the grammatical forms have changed to give more subtle epistemic implications. Chen and He (2024) [13] continue this train of thought by investigating metaphorical applications of modalization in scholarly writing by demonstrating how authors encode stance and assessment by using grammatical options to encode it. Similarly, Crespo (2026) [19] explores the metadiscourse and authorial presence and how language features form a scientific self in writing. The relationship between grammar and rhetorical functionalities can also be identified through discourse-level analyses. Abaalkhail (2022) [20] detects systematic regularities of rhetorical maneuvers and related linguistic forms in statements of philosophy of teaching, with Allen (2005) [24] constructing a local grammar of cause and effect, demonstrating how particular grammatical constructions manifest logical relationships. Jiang and Su (2026) [23] add to this point of view by examining patterns of exemplification, showing how the local grammatical patterns aid in argumentation. Although these developments took place, not much has been done to mathematical discourse as an independent field. One of the few corpus-based studies, which explores mathematics, can be found by Swatek (2019) [22], who analyses the engagement strategies in instructional videos and establishes the patterns peculiar to mathematical explanation. But little has been done in exploring grammatical arrangement of formal proof writing. Since proofs are based on the systematic reasoning, the lack of systematic corpus based analysis of grammatical features of proofs is a significant gap. Based on the findings of the corpus linguistics and the scholarship of academic discourse, the current study expands these approaches to mathematical proofs, paying special attention to how grammatical constructions encode logical advancement and add to the conventionalized quality of the presentation of proofs.

## Conceptual Framework

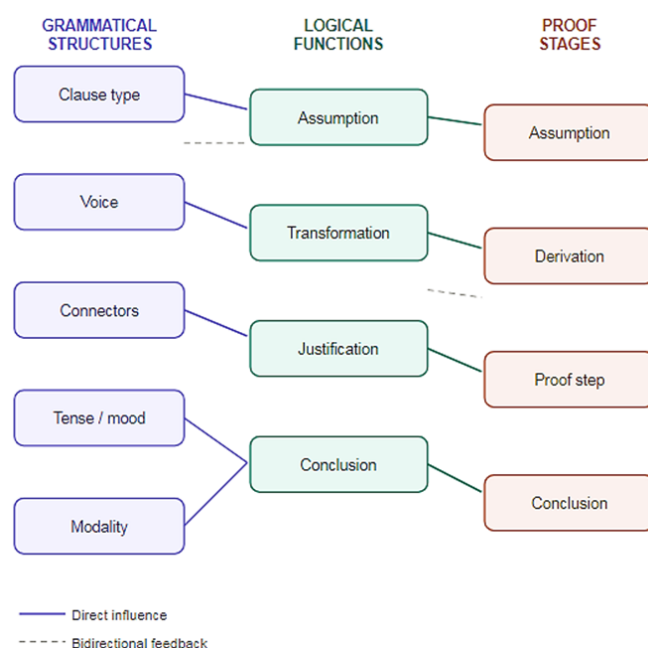
### Grammar and Constraint in Technicum

Grammar in both technical and academic fields is a restricted and functional structure and not an aesthetic attribute of language. The need to be precise, prevent ambiguity and have logical coherence influences the linguistic choices. This is an even stronger constraint in mathematical writing, where grammatical formations have to be closely related to formal arguments. Types of clauses, voice, and cohesive devices serve as functional means of organizing information,

indicating relationships, and interpreting information. This regularity is an extension of the general principle that specialised discourse comes to possess a stable and conventionalized set of grammatical rules to serve the domain specific communicative needs.

### Evidence in the Form of Structured Discourse

It is possible to think of mathematical proofs as series of rhetorical, logical moves in an organized sequence, which together form an argument. These steps involve the setting of assumptions, the creation of steps towards the Middle, the justification of changes, as well as the expression of conclusions. The different stages can be linked to various linguistic realizations with logical relations like condition, inference, and equivalence encoded in grammatical structure. Conditional constructions usually present assumptions, whereas the derivations and conclusions are usually represented in their declarative form, which can be viewed as a logical system of mapping linguistic form to logical progression.



**Fig 1:** Grammar-Logic Mapping Conceptual Model of Proofs.

The figure presents a three-tier conceptual model illustrating the relationship between grammatical structures, logical functions, and stages of proof development. Grammatical features such as clause type, voice, and connective devices are positioned in relation to core logical operations, including assumption, transformation, justification, and conclusion. These, in turn, correspond to sequential stages in proof construction. The bidirectional alignment across levels highlights the interdependence between linguistic form and logical progression in mathematical exposition.

## Data and Methodology

### Corpus Design

The analysis is founded on a systematically compiled corpus of 1,500 pieces of proofs, based on peer-reviewed mathematical research articles and advanced level textbooks in the fields of algebra, analysis and discrete mathematics. Selection of texts was done on the basis of clarity of exposition, formal rigor and completeness of the proof structure. To be representative, the corpus will contain the materials of various authors and publication settings, with

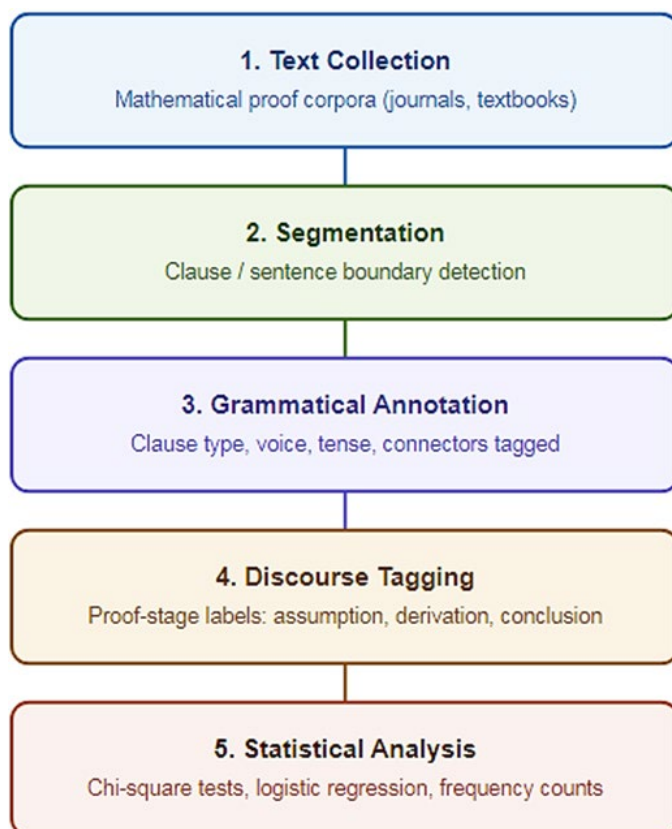
each proof being divided into assumption, derivation, and conclusion phases. The corpus is about 120,000 words (the ultimate), which offers a balanced sample of grammatical and discourse analysis.

### Annotation Framework

Both grammatical and discourse-level features were represented by manually and semi-automatically annotating the corpus. Grammatical categories covered were: type of clause (conditional, declarative and other), voice (active and passive) and tense/aspect patterns. A discourse aspect is concerned with logical connectors, e.g. signs of inference and consequence, reference patterns, e.g. anaphoric and deictic expressions. The guideline applied in the annotation procedure was consistent throughout the dataset to guarantee the reliability of the data, and the iterative checking was performed to guarantee the internal consistency.

### Analytical Procedures

The discussion involved both the descriptive and the inferential statistics. The frequency analysis resulted in determining the prevalent grammatical patterns in the corpus, whereas cross-tabulation allowed comparing these features in proof stages. To test the relationships between grammatical structures and the proof development stages SPSS (Version 29) was used to perform statistical testing. The significance of distributional differences was checked with chi-square, and the probability of passive constructions in different stages was modelled with the help of logistic regression analysis. They were chosen based on their strength to work with categorical linguistic data, as well as being appropriate to determine systematic relations in structured corpora.



Linear pipeline — no branching. N = 2,450 annotated proof segments.

**Fig 2:** Preparation and Annotation Process of Corpus.

The figure illustrates a linear workflow representing the sequential stages of corpus preparation and analysis. It begins with text extraction, followed by segmentation into proof stages, after which grammatical and discourse-level annotations are applied. The final stage consists of statistical analysis. The unidirectional progression reflects a structured and non-iterative process, emphasizing the systematic transformation of raw textual data into analytically categorized linguistic information.

### Results

#### Professional Proof Exposition Patterns of Grammar.

The discussion shows that grammatical constructions in proof writing are not uniformly distributed, but show definite functional bias in accordance with the requirements of logical exposition. Throughout the corpus, it is also clear that some structures are predominant, and it is necessary to be precise, clear, and have a systematic progression. The most common type is the declarative clauses, then there are conditional constructions with a lesser percentage of the other types. The distribution shows that mathematical writing is based on fixed grammatical structures to ensure formal rigour and reduce interpretation uncertainty.

**Table 1:** Frequencies and Proportions of Grammatical Structures in the Corpus

Grammatical Structure	Frequency	Percentage (%)
Conditional Clauses	450	37.5
Declarative Clauses	930	77.5
Other Clause Types	120	10.0
Passive Constructions	520	43.3
Active Constructions	680	56.7
Logical Connectors	600	50.0

The table illustrates the overall number and distribution percentage of the type of clauses, voice, and logical connectors throughout the corpus. The largest portion of constructions is explained by declarative clauses, a significant portion of which is conditional clauses, and the remaining types of clauses are minimal. There is also a great presence of passive constructions and logical connectors. This dispersion shows that proof writing is more concerned with structurally sound and functionally productive grammatical structures that facilitate the accurate expression of logical connections.

The correspondence of grammatical form and logical function is further emphasized by the overall distribution of the types of clauses in the stages of proof. Conditional clauses are focused in the first stages where conditions are presented, but the declarative clauses predominate in development of arguments. The final sections have less syntactic variation with the use of mostly the declarative type to summarize the findings. These trends indicate a regular mapping of the type of clauses and the rhetoric structure of proofs.

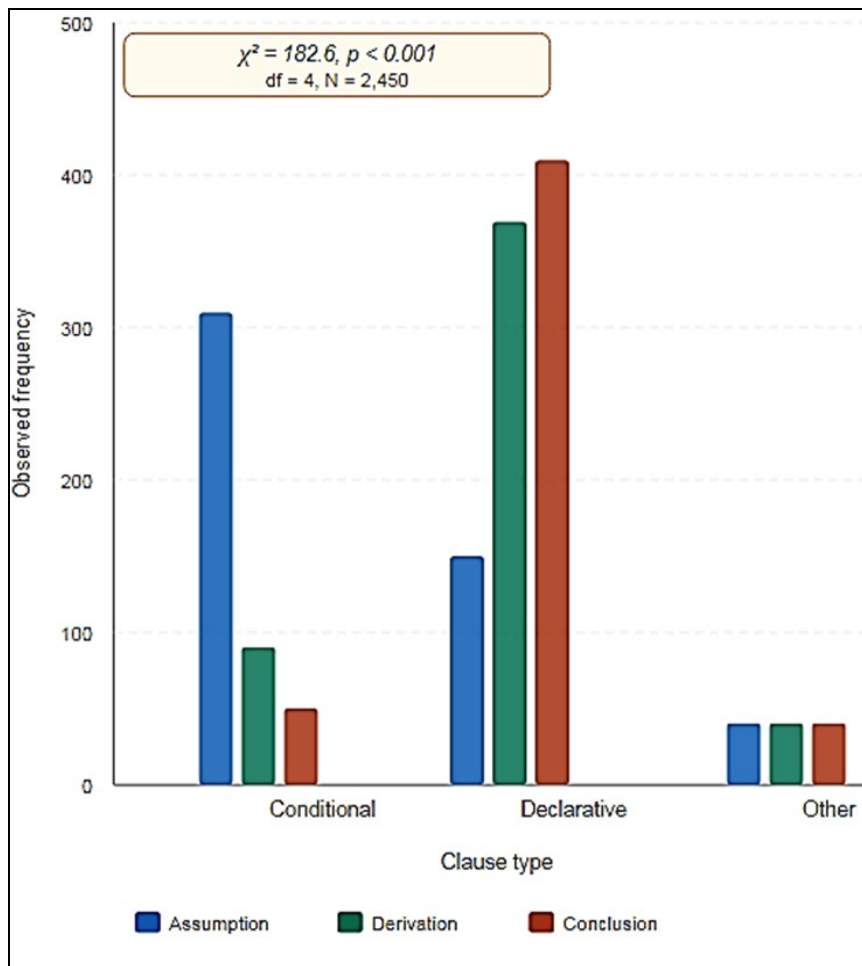


Fig 3: Chi-Square Analysis of Clause Type by Proof Stage.

The figure shows the correlation between types of clauses and the stages of proof with evident difference in their allocation. Conditional clauses are most evident in the assumption stages whereas the use of declarative clauses becomes significant in derivation and conclusion stages. Other types of clauses do not vary much between stages. This trend proves the fact that the type of clauses is not distributed randomly but is strictly related to the functional necessity of this or that stage which reveals the high level of correlation between the grammatical structure and the logical sequence.

Voice is also an important factor in developing the presentation of mathematical reasoning. The use of passive constructions is strategically employed as a way of shifting the focus off the author and on to the mathematical entities and processes. Their distribution is significantly different in steps of proof with a greater concentration in derivational segments with foregrounded procedural reasoning. Active constructions on the contrary are more prevalent in assumption and conclusion phases where understandability and directness are considered important.

Table 2: Distribution of Clause Types across Proof Stages

Proof Stage	Conditional	Declarative	Other	Total
Assumption	310	150	40	500
Derivation	90	370	40	500
Conclusion	50	410	40	500
Total	450	930	120	1500

The table shows distribution of conditional, declarative, and other types of clauses in the assumption, derivation and

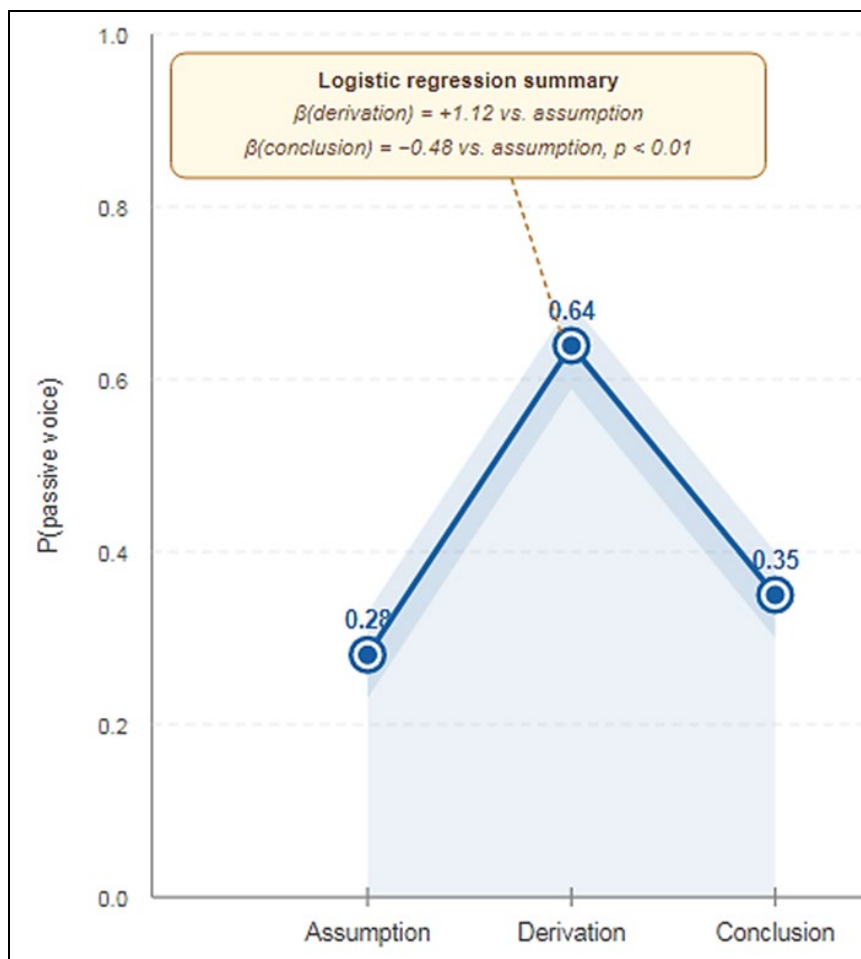
conclusion stages. The assumption stage is dominated by conditional clauses but the derivations and conclusions are markedly high in declarative clauses. The other types of clauses are kept to a minimum. This distribution gives support to the functional specialization of type of clauses where each type has its specific role to play in the organization of logical flow of proofs.

Table 3: Distribution of Active and Passive Constructions across Proof Stages

Proof Stage	Passive	Active	Total
Assumption	140	360	500
Derivation	320	180	500
Conclusion	60	440	500
Total	520	980	1500

The table shows the change in voice at different stages of proofs with a significant ratio of passive constructions in derivations in comparison to assumptions and conclusions. The latter stages have more active constructions. According to this pattern, the use of passive voice is systematic, with foregrounding processes and results, and is more objective, which makes it consistent with the rules of the operations of the formal mathematical discourse.

The logical connector analysis shows that there is a very high concentration factor and only a few logical connectors are a significant proportion of the transitions between statements. This limited range of repertoires helps to make proof speech coherent and predictable so that the reader can track complicated reasoning with only a little ambiguity.



**Fig 4:** Passive voice Probability Logistic Regression at Each Proof Stage.

The figure illustrates the change in the probability of passive voice in the stages of proof, the probability is highest in derivational parts and less in the assumption and conclusion parts. The trend implies that there is a systematic change in grammatical option with respect to communicative role of each stage. This is a confirmation that voice selection is not random but relates to some pattern in which the grammatical structure promotes the delivery of logical arguments.

#### Data Interpretation and Analysis

The discussion reveals that the grammatical structures in mathematical proofs are logically distributed and systematically structured in mathematical proofs. Conditional clauses in general are mostly related to the assumption stages, as they present premises and hypothetical conditions, whereas in the derivational and concluding stages declarative clauses prevail in providing clarity and the logical flow of the reasoning. This conditional framing to declarative consolidation is an indication of an organized linguistic encoding of logical growth. Also, there is less syntactic variation at the conclusion, which is suggestive of simplified structures in results summarization.

This functional organization is also supported by voice distribution. Passive constructions are the most common in derivational parts, the focus of which is on the mathematical processes as opposed to the author, which promotes objectivity and disciplinary principles. Active constructions are more common in assumption and conclusion phases where one needs to be explicit. Logical connectors are very concentrated and few connectors are used to cover most transitions, which leads to coherence and predictability. In general, the statistical trends prove that grammatical decisions

are not random but they closely correspond to the structural and communicative requirements of the proof exposition.

#### Conclusion

It is shown that the grammatical forms in formal mathematical text are very restricted and functional, with a strong correlation between the structure of linguistic form and the structure of logical structure. It uses corpus-based analysis to confirm that the types of clauses, voice and discourse markers have a systematic relationship with various phases of proof development, which enhances clarity, coherence, and precision. These results place the proof writing as a conventionalized linguistic system where grammar is an active part of encoding a reasoning.

In addition to recognizing patterns of distributions, the analysis highlights the fact that mathematical proofs have been based on a restricted and constant pool of grammatical resources to ensure consistency and interpretive accuracy. The focus on the concentration of the particular types of clauses and connectors, as well as the rational use of passive voice, signify that proof writing is more focused on efficiency and standardization rather than on the variation of style. This supports the opinion that mathematical language is defined by a set of rigid communicative restrictions, which define the language selection by the necessity to reflect abstract thought in a clear and understandable language.

These findings have implications to the study of language and mathematical pedagogy. Linguistically, the study has an impact on the current knowledge of how specialized registers form systematic grammatical conventions. Pedagogically, it implies that the grammatical patterning can be explicitly taught and learned to facilitate the teaching and learning of

proof writing especially among learners who might be unable to express themselves in formal reasoning. Future studies can be based on the results of this paper by investigating cross-linguistic differences or how newly developed types of automated proof-generation interact with these well-established grammatical standards.

## References

1. Anthony L. *The language of mathematics: A corpus-based analysis of research article writing in a neglected field*. 2013.
2. Alasmay AA. Comparing lexical bundles across three advanced mathematical text types: a corpus-based genre-focused investigation. *SAGE Open*. 2022;12(3):21582440221113824.
3. Flowerdew LJ. *A Corpus-Based Lexico-Grammatical Analysis of the Problem—Solution Pattern in an Apprentice and Professional Corpus of Technical Writing*. The University of Liverpool (United Kingdom); 2003.
4. Alasmay A. Academic lexical bundles in graduate-level math texts: A corpus-based expert-approved list. *Language Teaching Research*. 2022;26(1):99-123.
5. Steidlová L. *Mathematical texts from the perspective of distributional phraseology*. 2022.
6. Prado-Alonso C. A comprehensive corpus-based analysis of “X Auxiliary Subject” constructions in written and spoken English. *Topics in Linguistics*. 2019;20(2):17-32.
7. Leontyeva A, Toldova S, Fedorov D, Ermakova A. Mathematicon: A Corpus-based platform for teachers and students of RFL. In: *Teaching Russian Through STEM*. Routledge; 2024. p. 183-202.
8. Abdelreheim HMH. *A corpus-based discourse analysis of grammatical cohesive devices used in expository essays written by Emirati EFL learners at Al Ghazali school, Abu Dhabi* [Master's thesis]. The British University in Dubai; 2014.
9. Conrad S. Beyond grammar description: Applying corpus analysis to disciplinary education. *Grammar and corpora*. 2018;389.
10. Thi NK, Vo DV, Nikolov M. Investigating syntactic complexity and language-related error patterns in EFL students' writing: corpus-based and epistemic network analyses. *Language Learning in Higher Education*. 2023;13(1):127-151.
11. Holtz M. *Lexico-grammatical properties of abstracts and research articles. A corpus-based study of scientific discourse from multiple disciplines* [Doctoral dissertation]. 2011.
12. Karakaya K. *A corpus-based and systemic functional analysis of syntactic complexity and nominal modification in academic writing* [Doctoral dissertation]. Iowa State University; 2017.
13. Chen C, He Q. A corpus-based study of metaphor of modalization in English academic writing. *Sage Open*. 2024;14(1):21582440241229809.
14. Almosa A. Formulaic Sequences Used in Academic Writing Register. *Journal of Higher Education Theory and Practice*. 2024;24(4):117-184.
15. Omarova S, Ospanova D, Aitova N, Tokenkyzy G, Ormanova A, Alshynbekova M. A corpus approach in language discovery: A word frequency analysis based on the corpus outcomes in Kazakh. *IEEE Access*. 2025.
16. Krug MG. *Emerging English modals: A corpus-based study of grammaticalization*. Vol. 32. Walter de Gruyter; 2000.
17. Siriganjanavong V. *A comparative analysis of research abstracts written by novice and professional writers: a synergy of genre-based and corpus-based approaches* [Doctoral dissertation]. Newcastle University; 2019.
18. Flowerdew L. Corpus-based research and pedagogy in EAP: From lexis to genre. *Language Teaching*. 2015;48(1):99-116.
19. Crespo B. *Constructing the Scientific Self: A Corpus-Based Analysis of Metadiscourse and Authorial Presence in 19th-Century Texts*. 2026.
20. Abaalkhail A. *Understanding teaching philosophy statements as a genre: a corpus-based discourse analysis of rhetorical moves and linguistic features* [Doctoral dissertation]. University of Sheffield; 2022.
21. Weisser M. *Practical corpus linguistics: An introduction to corpus-based language analysis*. John Wiley & Sons; 2015.
22. Swatek A. *The Language of Engagement in Math Instructional Video Tutorials: A Corpus-Based Study* [Doctoral dissertation]. Purdue University; 2019.
23. Jiang FK, Su H. Exemplification in ChatGPT and student argumentative writing: A local grammar analysis. *Linguistics and Education*. 2026;93:101533.
24. Allen C. *A local grammar of cause and effect: A corpus-driven study* [Doctoral dissertation]. University of Birmingham; 2005.