



**THE LINGUISTIC PROBLEMS FACED BY TURKISH SCHOLARS
IN WRITING FOR PUBLICATION**

Batuhan Selvi

DOCTORAL DISSERTATION

DEPARTMENT OF FOREIGN LANGUAGE TEACHING

GAZİ UNIVERSITY

GRADUATE SCHOOL OF EDUCATIONAL SCIENCES

FEBRUARY, 2021

TELİF HAKKI VE TEZ FOTOKOPİ İZİN FORMU

Bu tezin tüm hakları saklıdır. Kaynak göstermek koşuluyla tezin teslim tarihinden itibaren 6 (altı) ay sonra tezden fotokopi çekilebilir.

YAZARIN

Adı : Batuhan

Soyadı: : Selvi

Bölümü : Yabancı Diller Eğitimi/İngiliz Dili Eğitimi

İmza :

Teslim Tarihi :

TEZİN

Türkçe Adı: Bilimsel Yayın Yapma Sürecinde Türk Akademisyenlerin Yaşadığı
Dilsel Sorunlar

İngilizce Adı: The Linguistic Problems Faced by Turkish Scholars in Writing for
Publication

ETİK İLKELERE UYGUNLUK BEYANI

Tez yazma sürecinde bilimsel ve etik ilkelere uyduđumu, yararlandıđım tüm kaynakları kaynak gösterme ilkelerine uygun olarak kaynakçada belirttiđimi ve bu bölümler dışındaki tüm ifadelerin şahsıma ait olduđunu beyan ederim.

Yazar Adı Soyadı: Batuhan SELVİ

İmza:

JÜRİ ONAY SAYFASI

Batuhan Selvi tarafından hazırlanan “The Linguistic Problems Faced by Turkish Scholars in Writing for Publication” adlı tez çalışması aşağıdaki jüri tarafından oy birliği / oy çokluğu ile Gazi Üniversitesi İngiliz Dili Eğitimi Anabilim Dalı’nda Doktora Tezi olarak kabul edilmiştir.

Danışman: Prof. Dr. Hacer Hande UYSAL GÜRDAL

İngiliz Dili Eğitimi, Hacettepe Üniversitesi

Başkan: Prof. Dr. Kadriye Dilek AKPINAR

İngiliz Dili Eğitimi, Gazi Üniversitesi

Üye: Doç. Dr. Zekiye Müge TAVİL

İngiliz Dili Eğitimi, Gazi Üniversitesi

Üye: Dr. Öğr. Üyesi Hatice ERGÜL

İngiliz Dili Eğitimi, Hacettepe Üniversitesi

Üye: Dr. Öğr. Üyesi Aysel SARICAOĞLU AYGAN

İngiliz Dili ve Edebiyatı, Ankara Sosyal Bilimler Üni.

Tez Savunma Tarihi: 01./02/2021

Bu tezin Eğitim Bilimleri Enstitüsü Yabancı Diller Eğitimi Anabilim Dalı’nda Doktora Tezi olması için şartları yerine getirdiğini onaylıyorum.

Prof. Dr. Yücel GELİŞLİ

Eğitim Bilimleri Enstitüsü Müdürü

ACKNOWLEDGMENT

I would like to express my sincere and deepest gratitude to the following people who helped and supported me in completing this dissertation.

First and foremost, I would like to first express my deepest gratefulness to my supervisor, Prof. Dr. Hacer Hande UYSAL GÜRDAL for her understanding, kindness and support during my doctoral education. I must thank her for her invaluable help, guidance, patience, encouragement, suggestions, and productive feedback which helped me to complete this dissertation.

I also would like to thank the committee members Prof. Dr. Kadriye Dilek BACANAK, Assoc. Prof. Dr. Müge TAVİL, Ass. Prof. Dr. Hatice ERGÜL and Ass. Prof. Dr. Aysel SARICAOĞLU AYGAN for taking their invaluable time to read my dissertation. They guided me and provided eminent suggestions and comments on my dissertation.

I am grateful to the people at the ELT Department of Firat University for their support and understanding. A lot of colleagues and friends have also encouraged and supported me in this difficult journey for years. First, I owe special thanks to Assoc. Prof. Dr. Aysel ŞAHİN KIZIL for providing invaluable suggestions and recommendations on corpus linguistics and research methodology. She has always expressed her thoughtful interest in my study, and has always guided me gain novel perspectives. I also want to thank Ass. Prof. Dr. Ömer Faruk CANTEKİN for helping me during the interrater reliability studies. Furthermore, I would like to thank Osman MEŞE, Fatih SANSAR, and Murat SANSAR for both their moral support and valuable ideas.

I also express my special gratitude to TUBİTAK for supporting me with 2211- PhD Scholarship Program throughout my doctoral education.

Finally, I would like express my greatest gratitude and love to my family for their endless support. Also, I would like to thank my brother Sadık Metehan SELVİ for his support throughout my graduate education. Last but not least I, I am also greatly indebted to my

beloved wife, Özgül SELVİ. I sincerely and with all my heart would like to thank for her endless love, support and understanding.



**BİLİMSEL YAYIN YAPMA SÜRECİNDE TÜRK
AKADEMİSYENLERİN YAŞADIĞI DİLSEL SORUNLAR
(Doktora Tezi)**

Batuhan Selvi

GAZİ ÜNİVERSİTESİ

EĞİTİM BİLİMLERİ ENSTİTÜSÜ

Şubat 2021

ÖZ

Dünya genelinde pek çok akademik kurum, akademisyenlerin işe alma, terfi ve ödüllendirme süreçlerinin bir parçası olarak, yüksek etki oranına sahip uluslararası dergilerde yayın yapmasını zorunlu hale getirmiştir. Ayrıca, dünya genelince birçok üniversitede uluslararası indeksli dergilerde yayın yapmak doktora öğrencilerinin mezun olabilmeleri için bir ön koşul haline gelmiştir. Bu anlamda, ağırlıklı olarak İngilizcenin hâkim olduğu uluslararası arenada yayın yapmak isteyen akademisyenlerin ve araştırmacıların sayısı hızla artmaktadır. Bununla birlikte, birçok bilim insanı, yayın bilimsel yayın yapma amacıyla yazma sürecinde bir takım sıkıntılar ve zorluklar yaşamaktadır. Bu nedenle, bilimsel çıktılarının hem niceliğini hem de niteliğini geliştirmelerine yardımcı olmak için, akademisyenlerin yayın süreçlerinde karşılaştıkları sorunları tespit etmek çok önemlidir. Bu çalışma, Türk bilim insanlarının bilimsel yayın yapma sürecinde makalelerini yazarken yaşadıkları dil problemlerini incelemeyi amaçlamaktadır. Bu nedenle bu çalışma kapsamında Multidisipliner Bilimsel Yayın Yapma Amacıyla Yazma Derlemi (The Multidisciplinary Corpus of Writing for Publication or MCWP) oluşturulmuştur. Derlem iletişim, ekonomi, eğitim, mühendislik ve tıp olmak üzere beş farklı disiplinden Türk akademisyenler tarafından yazılmış 216 adet düzenlenmemiş makaleden oluşmaktadır. Ayrıca, Türkçe derlemin karşılaştırılması amacıyla anadili İngilizce olan bilim insanlarının yazmış olduğu 163 yayınlanmış makaleden oluşan bir referans derlem hazırlanmıştır. Bu çalışma kapsamında üç dil boyutu incelenmiştir: sözcüksel çeşitlilik, sözdizimsel karmaşıklık ve yazım sürecinde yapılan dilsel hatalar. Sözcük çeşitliliği tür / türce oranı, standartlaştırılmış tür / türce oranı ve hareketli ortalama tür / türce oranı olmak üzere üç farklı ölçü kullanılarak araştırılmıştır. Sözdizimsel karmaşıklık analizinde Lu'nun (2010) İkinci Dil Sözdizimsel Karmaşıklık Analiz (L2SCA) programı kullanılmıştır. Hem çalışma hem de referans derleminde yer alan her bir makale 5 farklı kategoride 14 ölçü temel alınarak incelenmiştir.

Son olarak, dilbilgisi hataları, 54 alt alana bölünmüş 8 geniş kategoride analiz eden Louvain Hata Etiketleme Taksonomisi kullanılarak incelenmiştir. Tür / türce oranı sonuçları, MCWP'de yer alan akademisyenlerin daha az sözcüksel varyasyon kullanmış olmasına rağmen, hem MCWP hem de referans derlemde yer alan akademisyenlerin, makalelerinde oldukça çeşitli sözcükleri kullandıklarını ortaya çıkarmıştır. Sözdizimsel karmaşıklık analizi sonuçları, ana dili İngilizce olan akademisyenlerin, Türk akademisyenlere göre iki ölçü dışında tüm indekslerde daha karmaşık yapılar kullandığını göstermiştir. Bilgisayar destekli hata analizi, Türk bilim insanlarının bir metinde ortalama 110.15 hata ve 1000 türcede ortalama 20.08 hata yaptığını göstermiştir. Ayrıca en problematik üç kategorinin sırasıyla Dilbilgisi, Noktalama ve Kesit olduğu görülmüştür. Ayrıca bilgisayar destekli hata analizi, Türk bilim insanlarının özellikle artikel, fiil, noktalama işaretleri ve doğru sözcüğü kullanmada zorluk yaşadıklarını ortaya koymuştur. Bu çalışmanın sonucunda yazarların ana dilinin sözdizimsel karmaşıklıkta önemli bir faktör olduğu ve Türk bilim insanlarının İngilizce araştırma makalesi yazmanın retorik yönü hakkında yeterli bilgiye sahip olmadıkları ve bunun da metinlerinde daha fazla hataya yol açtığı sonucuna varılmıştır.

Anahtar Kelimeler : Yayın Yapma Amacıyla Yazma, Akademik Yazma, Sözdizimsel Karmaşıklık, Bilgisayar Destekli Hata Analizi, Bilimsel Yayın Yapma Amacıyla İngilizce
Sayfa Adedi : 261
Danışman : Prof. Dr. Hacer Hande UYSAL GÜRDAL

**THE LINGUISTIC PROBLEMS FACED BY TURKISH SCHOLARS
IN WRITING FOR PUBLICATION
(Ph.D Thesis)**

Batuhan Selvi

GAZI UNIVERSITY

GRADUATE SCHOOL OF EDUCATIONAL SCIENCES

February 2021

ABSTRACT

Many academic institutions worldwide have compelled academicians to publish their works in high-impact international journals as a part of institutional hiring, promotion and reward. Besides, at many universities around the world publishing in international indexed journals has even become a prerequisite for doctoral students to obtain their degree. In this sense, the number of academic members and research students, willing to secure publication in international arena which is overwhelmingly dominated by English, are increasing. However, many scholars experience a number of troubles and difficulties in writing for publication processes. Therefore, in order to help them improve both the quantity and quality of their scientific outputs, it is very important to identify the problems scholars face in their publication process. This study aims to investigate the linguistic problems experienced by Turkish scholars in their writing for publication process. For this reason, The Multidisciplinary Corpus of Writing for Publication was compiled as a part of this study. It consists of 216 unedited research articles written by Turkish scholars from five disciplines: communication, economics, education, engineering, and medicine. In addition, a reference corpus, consisting of 163 published research articles written by Native Speakers of English was compiled for the comparison of the Turkish corpus. Three linguistic dimensions were examined: lexical diversity, syntactic complexity, and errors they made in their manuscripts. Lexical diversity was investigated using three different measures: type/token ratio, standardized type/token ratio, and moving average type/ token ratio. In the analysis of syntactic complexity, Lu's (2010) L2 Syntactic Complexity Analyzer (L2SCA) was used. Each article in both the study and the reference corpus was investigated 14 measures in 5 different categories. Finally, grammatical errors were investigated using Louvain Error Tagging Taxonomy in which grammatical errors were analyzed in 8 broad categories which were broken down into 54 subdomains. The results of the type/token ratio analyses revealed

that scholars in both the MCWP and the reference corpus used a wide variety of vocabulary in their manuscripts even though the MCWP had slightly less lexical variation than the reference corpus. The results of syntactic complexity demonstrated that native scholars used significantly more syntactically complex structures in all of the indices except for two (CP/C and CP/T) in which Turkish scholars used comparable amounts. In this regard, Turkish scholars were found to have remarkable similarities with upper intermediate Japanese and Chinese learners in terms of syntactic complexity. The computer-aided error analysis showed that Turkish scholars made 20.08 errors per 1000 tokens with an average of 110.15 errors per text. It was also found that the most problematic three categories were Grammar, Punctuation, and Register, respectively. In addition, the examination of the errors revealed that Turkish scholars experienced difficulty, especially in articles, verbs, punctuation, and using the correct lexical item. It was concluded as a result of this study that the first language of the writers is a significant factor in syntactic complexity and Turkish scholars suffered from sufficient knowledge of the rhetorical aspect of writing a research article in English, which lead to more errors in their manuscripts.

Key Words : Writing for Publication, Academic Writing, Syntactic Complexity, Computer-aided Error Analysis, English for Research Publication Purposes

Page Number : 261

Supervisor : Prof. Dr. Hacer Hande UYSAL GÜRDAL

TABLE OF CONTENTS

TELİF HAKKI VE TEZ FOTOKOPİ İZİN FORMU	i
ETİK İLKELERE UYGUNLUK BEYANI.....	ii
JÜRİ ONAY SAYFASI.....	iii
ACKNOWLEDGMENT.....	iv
ÖZ	vi
ABSTRACT	viii
TABLE OF CONTENTS.....	x
LIST OF TABLES	xvi
LIST OF FIGURES	xx
CHAPTER I.....	1
INTRODUCTION.....	1
1.1. Background to the Study.....	1
1.2. The Significance of the Study.....	3
1.3. The Purpose of the Study	5
1.4. Definitions	7
CHAPTER II	9
LITERATURE REVIEW.....	9
2.1. Writing in a Second Language.....	9
2.1.1. A Brief History of Second Language Writing	10

2.1.1.1. <i>The Product Approach</i>	10
2.1.1.2. <i>Current – Traditional Rhetoric</i>	12
2.1.1.3. <i>The Process Approach</i>	13
2.1.1.4. <i>The Socio-cultural Approach</i>	14
2.1.1.5. <i>Genre Approach</i>	15
2.1.1.6. <i>English for Specific Purposes</i>	16
2.2. English for Academic Purposes	18
2.3. English for Research Publication Purposes	20
2.3.1. Types of Publications	20
2.3.1.1. <i>Book</i>	20
2.3.1.2. <i>Research Article</i>	21
2.3.1.3. <i>Conference Proceedings</i>	23
2.3.1.4. <i>Thesis and Dissertation</i>	24
2.3.2. The Theoretical Frameworks in ERPP	25
2.3.2.1. <i>Genre Approach</i>	25
2.3.2.2. <i>Discourse Analysis</i>	27
2.3.2.3. <i>Social Constructivism</i>	29
2.3.2.4. <i>Situated Learning Theory</i>	30
2.3.2.5. <i>Social Capital Theory</i>	31
2.3.3. Problems of Multilingual Scholars in Writing for Publication	32
2.3.3.1. <i>Narrow-mindedness</i>	33
2.3.3.2. <i>Rhetorical Problems</i>	33
2.3.3.3. <i>The Nature of Writing for Publication in English</i>	34
2.3.3.4. <i>Lack of Connections with Core Academic Communities</i>	35
2.3.3.5. <i>Bias Against Multilingual Scholars</i>	36
2.3.3.6. <i>Lack of Sufficient Funds to Conduct Research</i>	36

2.3.3.7. <i>Language Problems</i>	37
2.4. Linguistic Problems	38
2.4.1. Linguistic Accuracy	39
2.4.2. Linguistic Complexity	42
2.4.2.1. <i>Frequently Used Measures in Syntactic Complexity</i>	43
2.4.2.2. <i>Syntactic Complexity Studies in Second Language Writing</i>	46
2.5. Corpus Linguistics	50
2.5.1. What is Corpus?	51
2.5.1.1. <i>Mode of Communication in Corpora</i>	53
2.5.1.2. <i>Types of Corpora</i>	54
2.5.2. Learner Language.....	57
2.5.2.1. <i>Learner Corpora</i>	59
2.5.2.1.1. <i>Key Concepts in Learner Corpora</i>	60
2.5.2.2. <i>Research Orientations</i>	62
2.5.2.3. <i>Methodological Issues</i>	63
2.5.2.3.1. <i>Contrastive Interlanguage Analysis</i>	63
2.5.2.3.2. <i>The Integrated Contrastive Model</i>	65
2.5.2.3.3. <i>Computer-aided Error Analysis</i>	66
2.6. Corpora and Written Academic English	70
2.7. The Sociolinguistic Background in Turkey	73
CHAPTER III.....	77
METHODOLOGY	77
3.1. Methodological Background	77
3.2. The Collection of the Corpus	79
3.2.1. Data Collection	79
3.2.2. Design Criteria	81

3.2.3. Error Taxonomy	85
3.2.4. Error Tagging Procedure	88
3.2.5. Piloting	89
3.2.6. Inter-Rater Coding Reliability in the Study	90
3.3. Reference Corpus	90
3.4. Quantitative Analyses	91
3.4.1. Type / Token Ratio.....	92
3.4.2. Syntactic Complexity	92
3.4.3. Statistical Analysis	93
3.4.3.1. <i>Mann Whitney U</i>	94
3.4.3.2. <i>Independent Samples T-test</i>	94
3.4.3.3. <i>ANOVA</i>	94
CHAPTER IV	97
RESULTS.....	97
4.1. The Multidisciplinary Corpus of Writing for Publication	97
4.1.1. Type/Token Ratio.....	97
4.1.2. Syntactic Complexity	98
4.1.3. Computer-Aided Error Analysis	103
4.2. Communication Subcorpora	108
4.2.1. Type/Token Ratio.....	108
4.2.2. Syntactic Complexity	109
4.2.3. Computer-Aided Error Analysis	111
4.3. Economics Subcorpora	116
4.3.1. Type/Token Ratio.....	116
4.3.2. Syntactic Complexity	117
4.3.3. Computer-Aided Error Analysis	120

4.4. Education Subcorpora	124
4.4.1. Type/Token Ratio.....	124
4.4.2. Syntactic Complexity	125
4.4.3. Computer-Aided Error Analysis	127
4.5. Engineering Subcorpora.....	132
4.5.1. Type/Token Ratio.....	132
4.5.2. Syntactic Complexity	132
4.5.3. Computer-Aided Error Analysis	134
4.6. Medicine Subcorpora.....	139
4.6.1. Type/Token Ratio.....	139
4.6.2. Syntactic Complexity	140
4.6.3. Computer-Aided Error Analysis	141
4.7. Comparison of Turkish Subcorpora	146
4.7.1. Type/Token Ratio.....	146
4.7.2. Syntactic Complexity	147
4.7.2. Computer-Aided Error Analysis	152
CHAPTER V.....	157
DISCUSSION.....	157
5.1. Lexical Diversity	157
5.2. Syntactic Complexity	159
5.3. Computer-Aided Error Analysis	166
CHAPTER VI.....	177
CONCLUSION.....	177
6.1. Summary of the Study	177
6.2. Theoretical and Pedagogical Implications	180
6.3. Suggestions for Future Research	182

REFERENCES	185
APPENDICES	235
Appendix 1. List of Communication Research Articles in the Reference Corpus	236
Appendix 2. List of Economics and Administrative Sciences Research Articles in the Reference Corpus	237
Appendix 3. List of Education Research Articles in the Reference Corpus	238
Appendix 4. List of Engineering Research Articles in the Reference Corpus.....	241
Appendix 5. List of Medicine Research Articles in the Reference Corpus.....	246
Appendix 6. Skewness Kurtosis Values of Amount of Coordination for Subcorpora	252
Appendix 7. Skewness Kurtosis Values of Amount of Subordination for the Subcorpora.....	253
Appendix 8. Skewness Kurtosis Values of Degree of Phrasal Sophistication for the Subcorpora.....	254
Appendix 9. Skewness Kurtosis Values of Length of Production Unit	255
Appendix 10. Skewness Kurtosis Values of Overall Sentence Complexity for the Subcorpora.....	256
Appendix 11. Skewness Kurtosis Values of Amount of Coordination	257
Appendix 12. Skewness Kurtosis Values of Amount of Subordination	258
Appendix 13. Skewness Kurtosis Values of Degree of Phrasal Sophistication.....	259
Appendix 14. Skewness Kurtosis Values of Length of Production Unit	260
Appendix 15. Skewness Kurtosis Values of Overall Sentence Complexity.....	261

LIST OF TABLES

Table 1. <i>IMRAD Model</i>	23
Table 2. <i>CARS Model</i>	27
Table 3. <i>Developmental, Lerner and Lingua Franca Corpora</i>	55
Table 4. <i>Major English Language Corpora</i>	56
Table 5. <i>Monolingual and Comparable Corpora</i>	57
Table 6. <i>International Corpus of Learner English Design Criteria</i>	61
Table 7. <i>Learner Corpora Associated with Error Tagging Systems</i>	69
Table 8. <i>The Distribution of Research Articles in The Multidisciplinary Corpus of Writing for Publication Corpus</i>	80
Table 9. <i>Design Criteria of The Multidisciplinary Corpus of Writing for Publication</i>	85
Table 10. <i>Main Error Categories</i>	86
Table 11. <i>The Distribution of Research Articles in the Reference Corpus</i>	91
Table 12. <i>Definitions of Syntactic Complexity Measures</i>	93
Table 13. <i>Type/Token Ratio of the MCWP and Reference Corpus</i>	98
Table 14. <i>The Comparison of the T-Test Results of the Syntactic Complexity Measures</i> .	102
Table 15. <i>The Computer-Aided Error Analysis Results of the MCWP</i>	103
Table 16. <i>The Distribution of Form Errors in the MCWP</i>	103
Table 17. <i>The Distribution of Grammar Errors in the MCWP</i>	105
Table 18. <i>The Distribution of Lexis Errors in the MCWP</i>	106
Table 19. <i>The Distribution of Register Errors in the MCWP</i>	106
Table 20. <i>The Distribution of Style Errors in the MCWP</i>	106

Table 21. <i>The Distribution of Word Errors in the MCWP</i>	107
Table 22. <i>The Distribution of Lexicogrammar Errors in the MCWP</i>	107
Table 23. <i>The Distribution of Punctuation Errors in the MCWP</i>	108
Table 24. <i>Type/Token Ratio Results of Communication Subcorpora</i>	108
Table 25. <i>Means of Syntactic Complexity Measures for Communication Corpora</i>	109
Table 26. <i>Mann Whitney U Results of Communication Corpora</i>	111
Table 27. <i>The Computer-Aided Error Analysis Results of TR_COM Subcorpus</i>	112
Table 28. <i>The Distribution of Form Errors in TR_COM Subcorpus</i>	112
Table 29. <i>The Distribution of Grammar Errors in TR_COM Subcorpus</i>	113
Table 30. <i>The Distribution of Lexis Errors in TR_COM Subcorpus</i>	114
Table 31. <i>The Distribution of Register Errors in TR_COM Subcorpus</i>	114
Table 32. <i>The Distribution of Style Errors in TR_COM Subcorpus</i>	114
Table 33. <i>The Distribution of Word Errors in TR_COM Subcorpus</i>	115
Table 34. <i>The Distribution of Lexicogrammar Errors in TR_COM Subcorpus</i>	115
Table 35. <i>The Distribution of Punctuation Errors in TR_COM Subcorpus</i>	116
Table 36. <i>Type/Token Ratio Results of Economics Subcorpora</i>	116
Table 37. <i>Means of Syntactic Complexity Measures for Economics Corpora</i>	118
Table 38. <i>Mann Whitney U Results of Economics Corpora</i>	119
Table 39. <i>The Computer-Aided Error Analysis Results of TR_EAS Subcorpus</i>	120
Table 40. <i>The Distribution of Form Errors in TR_EAS Subcorpus</i>	120
Table 41. <i>The Distribution of Grammar Errors in TR_EAS Subcorpus</i>	121
Table 42. <i>The Distribution of Lexis Errors in TR_EAS Subcorpus</i>	122
Table 43. <i>The Distribution of Register Errors in TR_EAS Subcorpus</i>	122
Table 44. <i>The Distribution of Style Errors in TR_EAS Subcorpus</i>	123
Table 45. <i>The Distribution of Word Errors in TR_EAS Subcorpus</i>	123
Table 46. <i>The Distribution of Lexicogrammar Errors in TR_EAS Subcorpus</i>	124

Table 47. <i>The Distribution of Lexicogrammar Errors in TR_EAS Subcorpus</i>	124
Table 48. <i>Type/Token Ratio of Economics Subcorpora</i>	125
Table 49. <i>Mean Values and T-test Results of Education Subcorpora</i>	126
Table 50. <i>The Computer-Aided Error Analysis Results of TR_EDU Subcorpus</i>	127
Table 51. <i>The Distribution of Form Errors in TR_EDU Subcorpus</i>	128
Table 52. <i>The Distribution of Grammar Errors in TR_EDU Subcorpus</i>	129
Table 53. <i>The Distribution of Lexis Errors in TR_EDU Subcorpus</i>	130
Table 54. <i>The Distribution of Lexis Errors in TR_EDU Subcorpus</i>	130
Table 55. <i>The Distribution of Style Errors in TR_EDU Subcorpus</i>	130
Table 56. <i>The Distribution of Style Errors in TR_EDU Subcorpus</i>	131
Table 57. <i>The Distribution of Lexicogrammar Errors in TR_EDU Subcorpus</i>	131
Table 58. <i>The Distribution of Punctuation Errors in TR_EDU Subcorpus</i>	132
Table 59. <i>Type/Token Ratio of Engineering Subcorpora</i>	132
Table 60. <i>Mean Values and T-test Results of Engineering Corpora</i>	133
Table 61. <i>The Computer-Aided Error Analysis Results of TR_ENG Subcorpus</i>	135
Table 62. <i>The Distribution of Form Errors in TR_ENG Subcorpus</i>	135
Table 63. <i>The Distribution of Grammar Errors in TR_ENG Subcorpus</i>	136
Table 64. <i>The Distribution of Lexis Errors in TR_ENG Subcorpus</i>	137
Table 65. <i>The Distribution of Register Errors in TR_ENG Subcorpus</i>	137
Table 66. <i>The Distribution of Style Errors in TR_ENG Subcorpus</i>	137
Table 67. <i>The Distribution of Word Errors in TR_ENG Subcorpus</i>	138
Table 68. <i>The Distribution of Lexicogrammar Errors in TR_ENG Subcorpus</i>	138
Table 69. <i>The Distribution of Punctuation Errors in TR_ENG Subcorpus</i>	139
Table 70. <i>Type/Token Ratio of Medicine Subcorpora</i>	139
Table 71. <i>Mean Values and T-test Results of Medicine Subcorpora</i>	141
Table 72. <i>The Computer Aided Error Analysis Results of TR_MED Subcorpus</i>	142

Table 73. <i>The Distribution of Form Errors in TR_MED Subcorpus</i>	142
Table 74. <i>The Distribution of Grammar Errors in TR_MED Subcorpus</i>	143
Table 75. <i>The Distribution of Lexis Errors in TR_MED Subcorpus</i>	144
Table 76. <i>The Distribution of Register Errors in TR_MED Subcorpus</i>	144
Table 77. <i>The Distribution of Style Errors in TR_MED Subcorpus</i>	144
Table 78. <i>The Distribution of Word Errors in TR_MED Subcorpus</i>	145
Table 79. <i>The Distribution of Lexicogrammar Errors in TR_MED Subcorpus</i>	145
Table 80. <i>The Distribution of Punctuation Errors in TR_MED Subcorpus</i>	146
Table 81. <i>Type/Token Ratios of Turkish Subcorpora</i>	147
Table 82. <i>ANOVA Results of Length of Production Unit Category</i>	148
Table 83. <i>ANOVA Results of Amount of Subordination Category</i>	149
Table 84. <i>ANOVA Results of Amount of Coordination Category</i>	150
Table 85. <i>ANOVA Results of Degree of Phrasal Sophistication Category</i>	151
Table 86. <i>ANOVA Results of Overall Sentence Complexity Category</i>	151
Table 87. <i>The Comparison of Computer-Aided Error Analysis Results</i>	152
Table 88. <i>The Comparison of Form Errors</i>	152
Table 89. <i>The Comparison of Grammar Errors</i>	153
Table 90. <i>The Comparison of Lexis Errors</i>	154
Table 91. <i>The Comparison of Register Errors</i>	154
Table 92. <i>The Comparison of Style Errors</i>	155
Table 93. <i>The Comparison of Word Errors</i>	155
Table 94. <i>The Comparison of Lexicogrammar Errors</i>	156
Table 95. <i>The Comparison of Punctuation Errors</i>	156

LIST OF FIGURES

<i>Figure 1.</i> Product approach.....	11
<i>Figure 2.</i> Process approach model.....	14
<i>Figure 3.</i> Categories of ESP	17
<i>Figure 4.</i> The hour-glass model of article structure	22
<i>Figure 5.</i> Conventional IMRD Model	24
<i>Figure 6.</i> Taxonomy of Second Language Complexity.....	42
<i>Figure 7.</i> Sample of Concordance Search	53
<i>Figure 8.</i> Data types used in SLA research	58
<i>Figure 9.</i> The main concerns of mainstream language teaching	58
<i>Figure 10.</i> Focus on learner output.....	59
<i>Figure 11.</i> Contrastive interlanguage analysis.....	63
<i>Figure 12.</i> CIA ²	65
<i>Figure 13.</i> Integrated contrastive model.....	66
<i>Figure 14.</i> Data Collection Procedure	80
<i>Figure 15.</i> The Complete Error Tagging Taxonomy.....	87
<i>Figure 16.</i> An example of the UAM corpustool error tagging screen.....	89
<i>Figure 17.</i> The comparison of mean values of the length of production unit measures.....	99
<i>Figure 18.</i> The comparison of mean values of the the amount of subordination measures	99
<i>Figure 19.</i> The comparison of mean values of the amount of coordination measures	100
<i>Figure 20.</i> The comparison of mean values of the degree of phrasal sophistication measures	101

Figure 21. The comparison of mean values of the overall sentence complexity measure 101



CHAPTER I

INTRODUCTION

1.1. Background to the Study

Rapid and extensive developments in information and communication technologies (ICT) have been significantly influencing our lives since the information revolution took place in the second half of the 20th century. The information revolution has gradually brought about a shift in the fundamental human resources from land to scientific knowledge (Dhia, 2006). The changes and the evolution during this period have turned information into a basic commodity that is used as an indicator of development, prosperity, and education of countries on both national and international scale (Carneiro, 2016). As a result, universities, as the specialized institutions where scientific knowledge is produced and disseminated and where the next generation is educated, have gained critical importance. Therefore, countries, aiming to secure a position in the competitive global market, started to make massive investments in universities for carrying out scientific activities and boosted research programs (Min, 2014).

In order to aid sustainable development, governments around the world have established specific criteria for increasing the competitiveness of their universities to gain an international reputation and attract students (Jiang, Borg & Borg, 2017). In this reputation race (Hazelkorn, 2011), the success of universities is generally measured through global ranking systems in which international publications and citations are two major components, putting pressure on universities to produce publishable research, especially in the high-ranked international journals (Gonzales, Martinez & Ordu, 2014). For that reason, many academic institutions worldwide have compelled academicians to publish their works in high-impact international journals as a part of institutional hiring, promotion, and reward systems (Baldwin & Chandler, 2002; Canagarajah, 1996; Englander & Uzuner-Smith, 2013; Flowerdew, 1999a; Lillis & Curry, 2010). Besides, at many universities around the world,

publishing in international indexed journals has even become a prerequisite for doctoral students to obtain their degrees (Ho, 2017). Since most of the journals included in the highly credible international databases such as Web of Science Core Collection (Clarivate Analytics, 2016) and Ulrich's Periodicals Directory (Wang, Hu & Liu, 2017) are published in English, English-medium publishing has become a must for scholars all over the world for international reputation and higher citation rates (Bocanegra-Valle, 2014). In this sense, there is no doubt that English has secured its position as the primary language of scholarly publications and international research (Bardi, 2015; Flowerdew, 2000; Swales, 1990), involving more than 5.5 million scholars and 2,000 publishers all over the world (Lillis & Curry, 2010).

Starting from the 1980s, the dominant position of English in academia has attracted considerable attention in the English language teaching (ELT) field as multilingual authors have increasingly been forced to learn academic English to publish in international journals. On the one hand, some scholars enthusiastically welcomed the spread of English in facilitating globalization and modernization (Grabe, 1988); others, on the other hand, regarded the dominance of English as a form of cultural and linguistic imperialism (Phillipson, 1992). Since then, questions regarding the standards of academic English writing and the challenges faced by multilingual authors in scholarly publishing have been raised (Koyalan & Mumford, 2011). These questions have mainly revolved around the Native Speaker (NS) - Non-Native Speaker (NNS), center – peripheral, and novice vs. expert dichotomies and addressed such issues as the linguistic capabilities, rhetorical resources (Lillis & Curry, 2006) and non-discursive resources (Canagarajah, 1996) required for international scholarly publication (for a detailed review see Uzuner, 2008).

The literature on writing for publication has successfully explored a wide range of key areas where multilingual scholars encounter difficulty. It has been argued that scholars from outside the inner circle (Kachru, 1985) face various challenges and find themselves at a disadvantageous position in publishing their articles in English (Burgess, Gea-Valor, Moreno & Rey-Rocha, 2014) since they have “the triple disadvantage of having to read, do research and write in another language” (Van Dijk, 1994, p. 276). In this vein, empirical evidence suggests that such language-related problems of the multilingual authors as using limited vocabulary (Flowerdew, 1999a), inappropriate use of idioms (Kaplan & Baldauf, 2005), and complexity of syntax (Muresan & Perez-Llantada, 2014) may pave the way for the rejection of their manuscripts (Hewings, 2006).

It has also been documented that being competent in English is not enough to achieve publication since multilingual scholars have to conform to the standard Anglophone rhetorical and publication norms (Bennett, 2010; Swales, 2004), which is a serious challenge for effective writing as rhetoric is a culture-specific phenomenon (Uysal, 2012a). The other area of difficulty for multilingual scholars is reported as the potential prejudice of the journal gatekeepers, namely publishers, editors, and reviewers who are mostly native English speakers. Scholars all over the world, including Hong-Kong (Flowerdew, 1999a), China (Li & Flowerdew, 2007), Germany (Ammon, 2000), and Hungary (Curry & Lillis, 2004) have reported that they thought that the journal gatekeepers are biased towards them because of their linguistic backgrounds, which, in turn, may lead multilingual scholars to be stigmatized as their work are treated as of lower quality (Flowerdew, 2008). Drawing attention to the inequalities between the scholars in the center and those in the periphery, Canagarajah (1996) also emphasized the role of non-discursive resources that interfere in securing publication. He identified a number of non-discursive requirements and classified them into three categories: material, financial, and social. These requirements include access to libraries, research funds, and technological devices.

The aforementioned brief overview reveals that although multilingual scholars experience a number of difficulties, they have to overcome them in order to secure international publication since their academic careers are closely linked to publishing in English. As social and cultural factors shape the way people write (Atkinson, 2003), and thus distinctive writing practices exist in each society (Uysal, 2008), it is very important to identify the problems scholars face in their publication process, taking their local environment into account (Duszak & Lewkowicz, 2008). Therefore, second/foreign language writing professionals, one of whose roles is to guide scholars to acquire the necessary academic writing skills (Uysal, 2012b) and to secure publication, should first identify the problems in writing for publication in their local contexts so that they will develop pedagogical solutions to remedy these problems. It seems that only in this way can they help scholars overcome their difficulties in writing and publishing.

1.2. The Significance of the Study

Being the language of scientific publication, the increasing trend to publish in English has also been growing in Turkish academia. In 2019, Turkish academics published 49,930 research articles, indicating an almost 8% increase compared to 2018 and a nearly 11%

increase compared to 2017 (SCImago, 2020). The main reason for this growth is linked to state policies and planning aiming at increasing the production of the scientific output by indicating a shift from teaching to research, which requires the academic staff to publish in international journals in English (Uysal, 2014).

Recent changes in hiring and promotion policies in higher education institutes add a further emphasis on the issue. According to The Turkish Regulation on Academic Incentive Payments issued in the Official Gazette dated 31/12/2016 (YOK, 2017a), the point that can be obtained for scientific activities carried out in SSCI, SCI-Expanded or AHCI indexed journals is twice as much as that for activities in other international and national journals. Similarly, more universities are increasingly adopting a policy that requires doctoral students to publish in an international indexed journal before graduation. Moreover, the Turkish government has recently determined research universities whose aims include excellence not only in research but also in the production, transfer, and sharing of knowledge within the scope of the Mission Differentiation and Specialization Project, which was launched by the Council of Higher Education (YOK, 2017b). These developments propose that the role of scientific publication in well regarded international journals will be more important in Turkey in the following years. This points out an urgent need for scholars to improve the quality and quantity of their scientific outputs. However, the empirical evidence reveals that Turkish writers have unique rhetorical and argumentative differences and preferences (Uysal, 2008), which are likely to hinder their success of writing in English. Therefore, in order to improve the scientific output of Turkish scholars and satisfy the aforementioned criteria, more studies aiming to investigate their academic writing conventions are needed.

Although writing for publication literature includes well-documented information about a wide range of disciplines in many countries all over the globe, including both central and peripheral countries, these studies suffer from several weaknesses.

First, most of the studies in the literature have focused on a limited number of disciplines, such as English Language Teaching. However, it is fair to claim that academic writing cannot be limited to a particular field since each discipline has different discursive requirements. Therefore, more comprehensive studies, including various fields and disciplines, are needed in the issue of writing for publication.

Second, it seems that a great number of works have focused on graduate and undergraduate level students. Therefore, more studies addressing the writing of academicians are needed since the results of graduate and undergraduate level academic writing studies may not be

generalized to the academia since writing needs and purposes of scholars and students are quite different (Uzuner, 2008).

Third, the studies in writing for publication literature typically involve participants from mainly European countries (Lillis & Curry, 2010) or China (Li, 2007). Nevertheless, in order to respond to the changing demands and challenges of academia in terms of writing for publication, especially writing in English, more studies investigating the scholars from different cultural and scientific backgrounds are also needed since writing for publication is a different phenomenon with its distinctive burdens and challenges.

Finally, it seems that most of the studies in the literature have investigated published research articles. However, analyzing only published articles through corpus research may not be an effective method since published research articles, which are shaped, and reach their final form after several drafts, are *brokered* by several shareholders including editors, reviewers, and language specialists during the publication process (Lillis & Curry, 2006). Therefore much of the information a research article provides as to the challenges and problems of the international scholars in their publication quest may be lost in the publication process. For that reason, the first drafts or unedited research papers written by scholars and should also be investigated (Wu, Mauranen & Lei, 2020).

As for the Turkish context, the academic writing literature in Turkey is limited to English Language Teaching departments and preparatory English classes, and often included graduate or undergraduate level ELT students as the participants (Kan, 2017). Almost no research has addressed the issue of writing for publication and challenges scholars face while trying to satisfy the international requirements listed above. All these indicate a serious gap regarding the academic writing studies in the Turkish context as well. Aiming at identifying the linguistic difficulties faced by Turkish scholars during their publication processes by examining unedited drafts of their scientific texts, this study is an attempt to fill this gap in the literature.

1.3. The Purpose of the Study

No studies in the literature have focused on Turkish scholars' writing challenges in the publication process yet. To contribute to this lack in the field, the present study investigated Turkish scholars' writing for publication practices to identify the linguistic challenges they face, which often affect their scientific products. For that reason, the following linguistic

features were the target of investigation: lexical diversity, syntactic complexity, and computer-aided error analysis. The significance of this study lies in its data collection method, which aimed to collect the drafts of research papers written prior to submission. Adopting a corpus-based approach, this study obtained detailed information from the research papers and, in turn, identified the similarities and differences of writing practices between Turkish scholars and native scholars and what kind of linguistic challenges Turkish scholars face in this process. With these powerful data collection and analysis methods, this study identified a number of linguistic challenges of Turkish scholars in writing for publication, which potentially has crucial implications for writing for publication literature and facilitating the scholarly writing process of international scholars. To do that, a corpus, including unedited research articles written by Turkish scholars from five different disciplines, were compiled. In addition, a comparison corpus, including native speakers' research articles published in high-ranked journals, was also compiled in order to reveal the differences and identify the unique problems of Turkish scholars. Finally, the Turkish corpus was investigated in terms of errors Turkish scholars made in their manuscripts, using computer-aided error analysis. This study, meanwhile, argues that identifying the linguistic challenges of the scholars would shed light on the future training programs in academic writing in English and thus would help policy makers and teachers improve the quality of the academic writing courses and programs.

In this sense, the present study addressed the following research questions:

- 1) What are the lexical diversity levels in research articles written by Turkish scholars and English Native Speakers?
 - a) whether and to what extent these levels differ?
 - b) whether and to what extent these levels are affected by discipline?
- 2) What are the features of syntactic complexity in research articles written by Turkish scholars and English Native Speakers?
 - a) whether and to what extent these features differ?
 - b) whether and to what extent these features are affected by discipline?
- 3) What are the most common types of errors in research articles written by Turkish scholars?
 - a) whether and to what extent these errors are affected by discipline?

1.4. Definitions

- a. Social Sciences Citation Index (SSCI): SSCI is a commercial multidisciplinary citation index that covers over 3,000 social sciences journals across more than 50 disciplines.
- b. Science Citation Index SCI: SCI is a highly selective subset of journals that are typically the most consistently high impact titles in many scientific disciplines.
- c. Science Citation Index – Expanded (SCI-Expanded): SCI- Expanded is the larger version of the Science Citation Index. It covers more than 6,500 notable and significant journals across 150 disciplines.
- d. Arts and Humanities Citation Index (AHCI): ACHI is a citation index, with abstracting and indexing for more than 1,700 arts and humanities journals and coverage of disciplines that includes social and natural science journals.



CHAPTER II

LITERATURE REVIEW

This chapter provides a review of literature on L2 writing and writing difficulties. Chapter 2 includes an overview of the historical development of second language writing research, the theoretical foundations of academic writing, the relationship between L1 and L2 writing, the academic writing difficulties and challenges of the academicians, and the previous research on writing for publication.

2.1. Writing in a Second Language

Second language writing has recently been the focus of an extensive amount of research. Now, studies on second language writing can be found in nearly every issue of applied linguistics or educational journals. However, only after the middle of the 20th century, formal investigations into second language writing emerged as a consequence of the increasing number of international students registered in higher education institutes of English-speaking countries (Hinkel, 2011). After that time, studies aiming to provide instructional models and materials for second language learners to improve their writing skills flourished.

As a result of a large body of research on particularly L2 writing, it has been understood that writing in a second language differs considerably from writing in a first language in that it is “strategically, rhetorically, and linguistically different in important ways from L1 writing” (Silva, 1993, p. 669). In his seminal paper, Silva (1993) defines these differences under two categories: Composing processes and written text features. While the former covers three subprocesses, which are planning, transcribing, and reviewing, the latter involves fluency, accuracy, quality, and structure. For the first category, Silva reported that second language writers tend to spend less time on planning; need considerable time and effort to produce

given texts; have problems in fluency and productivity; use a dictionary and review less. Furthermore, Silva stated for the second category that second language writers are less fluent, less effective, and make more errors.

In addition, a lack of proficiency in the second language is asserted as another source of differences between second language writing and first language writing (Bardovi-Harlig, 1995; Cumming, 1989). Other differences include mastery of genres in the second language (Swales, 1990) the amount of experience in the first language writing, and the differences between the first and second language writing culture (Connor, 1996). Although it is claimed that most of these differences are evident, especially with second language writers having low L2 proficiency, and thus depending mostly on their first language during the writing (Zimmerman, 2000), there exists a significant individual variation as well.

In sum, empirical findings indicate a significant difference between writing in the first language and writing in a second language. Therefore, in order to thoroughly understand this phenomenon and get valuable insights, it will be beneficial to review second language writing history.

2.1.1. A Brief History of Second Language Writing

Having emerged around the 1950s, second language writing is relatively a young research field (Silva, 1990) and has grown, to some extent, in line with the developments in the teaching of L1 English writing (Yağız, 2009). The history of second language writing can be divided into several periods depending on the dominant teaching approach. It is important to note that these periods should be regarded as being complementary rather than being sequential (Ferris & Hedgecock, 2005). The overview presented here is an integration of the previous historical reviews (Matsuda, 2003; Silva & Leki, 2004) in the literature.

2.1.1.1. The Product Approach

Influenced by structural linguistics and behaviorism, the distinctive feature of this period was the practice of teaching writing as a habit formation through imitation (Silva, 1990). Focusing on the accurate use of grammar rules on sentence level, this approach regarded the writer as the manipulator of the learned grammatical structures and the reader as the proofreader of the produced output.

The lessons start with the introduction of a model text, and students are encouraged to imitate the structures of the given text. Then, students are expected to produce a very similar text. Thus, it can be said that this model prioritizes the mechanical skills of writing. Proponents of this model argue that the product approach enables learners to advance their writing skills, especially accuracy, error control, and essay structure skills, and paves the way for the development of motivation and self-confidence in learners, which in turn help learners promote their writing (Gorrell, 1981). This traditional approach usually consists of 4 stages (Hyland, 2003).

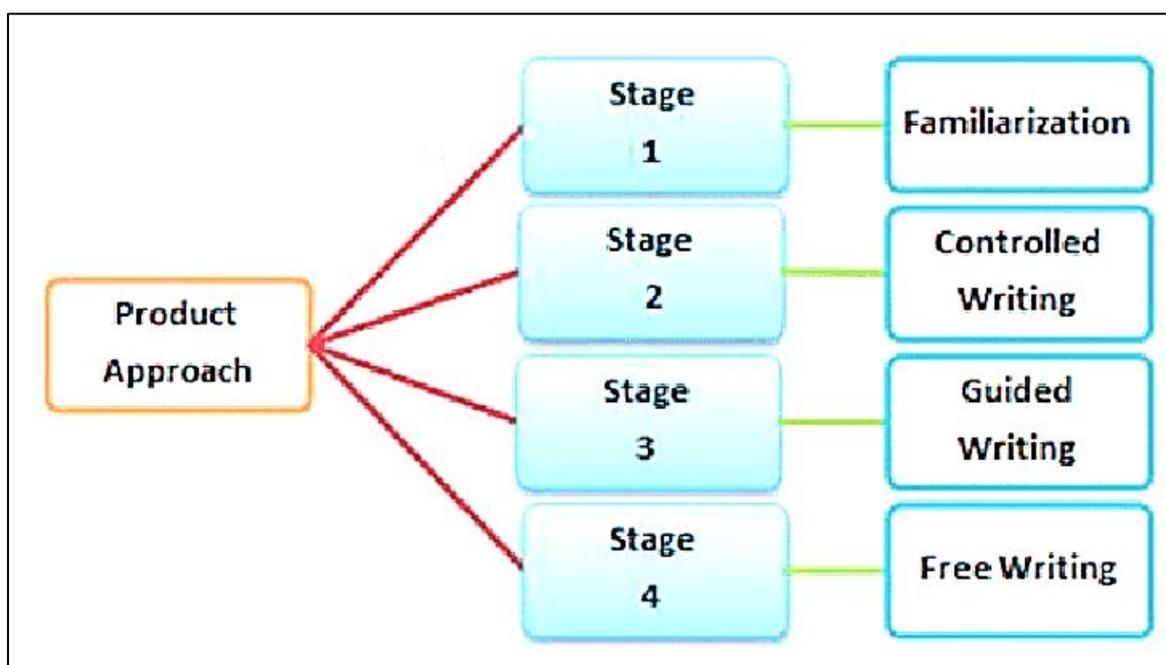


Figure 1. Product approach

In Stage 1, the model text is presented, and its distinctive aspects are emphasized. After that, learners are given controlled exercises regarding the elements of the model text and practice these features in Stage 2. Then, learners start to organize ideas to mimic the model text in Stage 3, which is considered the most critical stage of this approach. Finally, in Stage 4, learners are expected to produce a final error-free product.

In sum, the product approach to writing can be summarized as having the following features: i) imitation of a model text, ii) an emphasis on linguistic features rather than the ideas, iii) producing only one draft, iv) focusing on and practicing the features of the model text, and v) emphasis on a final product (Steele, 1992).

2.1.1.2. Current – Traditional Rhetoric

The 1960s witnessed an increase in the number of international students in American universities, which led to an awareness of learner needs and achievements. As the product approach and the controlled composition model was only limited to sentence-level structures, it fell short to meet the expectations of college-level composition classes in terms of teaching how to write longer paragraphs and essays as well as fulfilling the communicative function of writing (Hutchinson & Waters, 1987). Therefore, a need to “shift the focus of attention from the grammatical to the communicative properties of language” (Allen & Widdowson, 1974, p. 1) arose. The attempts to facilitate paragraph formation pointed to the conclusion that controlled and freewriting should be integrated. The solution was offered by the current – traditional rhetoric. The current – traditional rhetoric approach focuses on such elements as the topic, supporting, and concluding sentences in a paragraph. It also gives importance to the adequate development of paragraphs, organizational units, and patterns (Silva, 1990). However, it was criticized for being linear and prescriptive (Leki, 1991). Silva (1990) accused this approach of turning writing into a skill of arrangement by “fitting sentences and paragraphs into prescribed patterns” (p.14).

Kaplan’s theory of contrastive rhetoric (1966) also influenced this period. Examining the relationship between the culture and rhetoric, Kaplan concentrated on the rhetorical problems second language writers experienced. Influenced by a range of theories such as linguistic relativity, text linguistics, rhetoric, genres, and literacy, Kaplan argued that L2 writers use their L1 rhetorical traditions in the organization of their paragraphs and thus L1 rhetorical practices are transferred to L2 writing, which causes interference (Kaplan, 1966; Connor, 1996; Uysal, 2008). Since rhetoric was seen as language and culture-specific, paragraph structures were also regarded as culture and language-specific (Matsuda, 2003). Addressing the logical organization, this period stressed paragraph and essay development, which required learners to analyze a model text and apply the rhetorical patterns in their writing (Silva, 1990). However, it was criticized for being static in that it depends primarily on L2 writers linguistic, cultural, and educational backgrounds as the sole features promoting their second language writing improvement (Matsuda, 1997).

2.1.1.3. The Process Approach

In the 1980s, composition studies investigated how learners carry out their tasks and the ways writers write in order to identify the processes of poor and good writers, which led to the conclusion that writing is a “non-linear, exploratory and generative process whereby writers discover and reformulate their ideas as they attempt to approximate meaning” (Zamel, 1983, p. 165). Thus, writing started to be regarded as a complicated cognitive activity comprising such skills as choosing topics, generating ideas, cohesion, and coherence (Xiaoxiao & Yan, 2010). In this regard, second language writers are expected to use their independent thinking skills during the accomplishment of a writing task. Rather than focusing on the product, the process approach emphasized the process underlying writing, which was thought to be recursive (Hyland, 2015b). Focusing on meaning-making, the process approach allows learners to control their writing and choose their topics. In this way, learners were encouraged to focus on the content and creating meaning, which facilitated learning (Brown, 2001).

The process approach to writing in English received broad support from the instructors and researchers over the former approaches to writing. The process approach prioritizes the process of the writing instead of the final product since it is thought that “a good product depends on the good process” (Sun & Feng, 2009, p. 151). Second language writers have opportunities of thinking about their writing topics, prepare drafts, review and edit what they have written as well as obtaining feedback from their instructors and peers during the writing task. The role of feedback during the feedback is of crucial importance. The supporters of this approach claim that second language writers feel confident and develop positive attitudes when they receive feedback during the process (Stanley, 1993). A sample model of the process approach to writing is presented in Figure 2. However, the process approach underemphasized the role of product and regarded it as secondary and thus failed in more academic contexts (Horowitz, 1986).

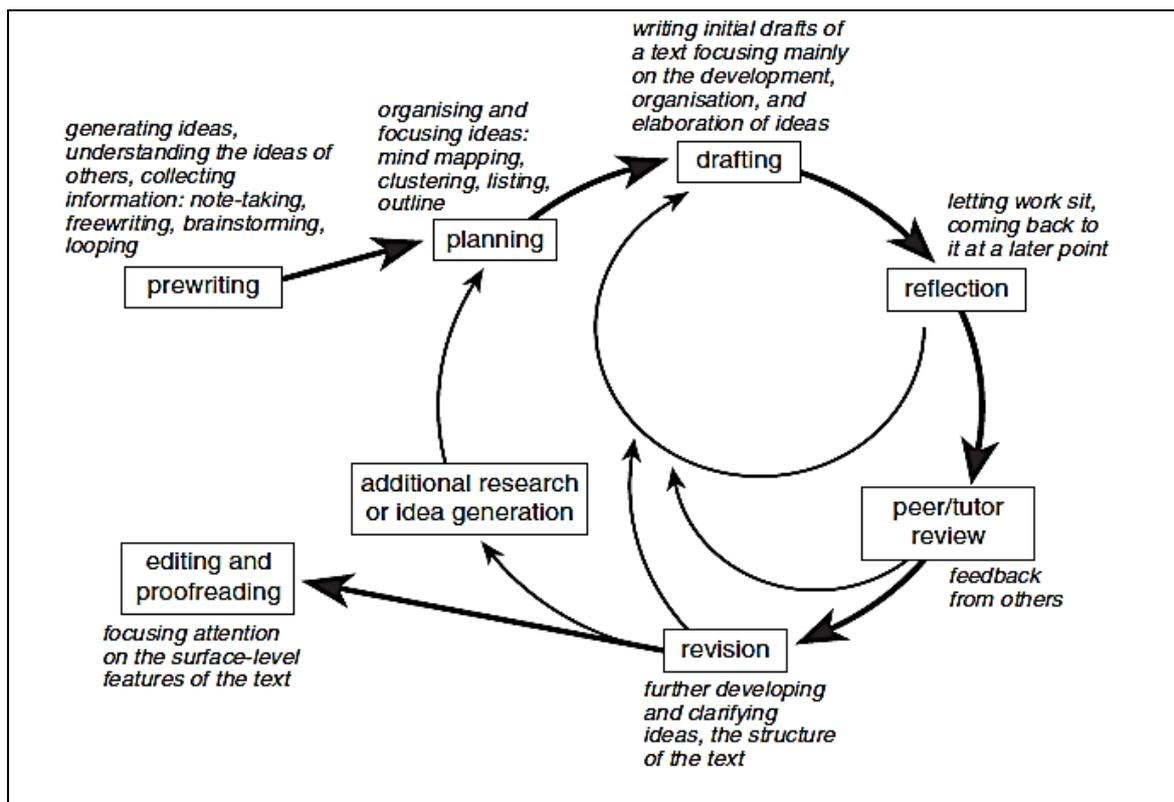


Figure 2. Process approach model

2.1.1.4. The Socio-cultural Approach

The socio-cultural approach to writing is profoundly inspired by Vygotsky's (1978) cognitive development theory. According to Vygotsky (1978), the main subject of the socio-cultural approach is that knowledge is inherently social and is formulated through a course of collaboration, interaction and communication in social situations. The socio-cultural approach has three main underlying principles: mediation, the activity theory, and the zone of proximal development (ZPD). As stated by Vygotsky (1978), each precisely human psychological process is mediated by psychological tools, including arithmetic systems, language, and semiotic systems. These psychological tools are taught to children throughout their joint activity with adults. As a result, these tools are adopted by children and then operate as mediators for the children to advance into higher psychological processes (Karpov & Haywood, 1998). The second crucial concept is ZPD. It is proposed that people have learning capacities that are recognized through the assistance of a more capable other in the course of action (Wigglesworth & Storch, 2012). Such assistance is *scaffolded* to the needs of the children (Wood, Bruner, & Ross, 1976), allowing them to advance their mental power by means of supporting them to access ZPD. Scaffolding, here, refers to the help or assistance given during the course of interaction by a more capable participant such as a

teacher or peer (Donato, 1994). The last concept is activity theory. Put forward by Leontiev (1978); it suggests that three levels exist in an activity: the motives that bring about the activity, the actions generated by objectives to succeed the action, and the conditions (or operations) under which the activity is performed.

Writing calls for students to accomplish some sort of activities and respond to feedback in pairs. This paves the way for students to involve in a lifelike social situation and offers them occasions to establish meaningful contact with each other, the piece of writing, and feedback. During this process, students mediate their actions through language, and communicate with the feedback they receive from their partners, and get chances to scaffold one another's assistance and to examine the feedback (Wigglesworth & Storch, 2012). It is predicted that such a process could result in improved learning for the reason that the discussion activities students take part in with their partners may lead students to notice features of language which they had failed before, and involved in conversations about language by means of which they possibly improve their knowledge on the language. Therefore, learning occurs as a result of the discussions and conversations about the language students utilize and with respect to the feedback they are provided. This concept of the socio-cultural theory to writing is reflected in error correction and peer-review in which one of the peers supports the other to attain advanced level of competence (Villamil & DE Guerrero, 2006). As a result, the proponents of this approach argue that the peer review process will help students improve their writing abilities (van Zundert, Sluijsmans & van Merriënboer, 2010).

2.1.1.5. Genre Approach

Having been inspired by Halliday's (1994) Systemic Functional Linguistics, which deals with the association between language and its social functions, the genre-based approach to writing emerged as a response to the disadvantages of the process approach in the 1990s. Genre, in its basic form, can be described as "abstract, socially recognized ways of using language" (Hyland, 2007a, p. 149). The central premise of the genre is that participants of a group typically identify similarities in the texts that are often used among the group and can employ their recurrent practices with these texts to read, comprehend, and possibly write them in a rather easy way (Hyland, 2003). Furthermore, genre depends on a communication situation occurring in a social context; and thus, the different social contexts are likely to result in different genres. In that sense, the genre is a goal-oriented course of communication for its participants in a particular communication situation taking place in a specific social

context. In other words, individuals do not merely write; instead, they write in order to complete different requirements in different contexts, which implicates differences in the means of using language (Halliday, 1994).

The genre-based approach to writing concentrates on the conventions of a specific type of text and make an effort to assist writers in producing a text based on purpose, audience, and organization (Paltridge, 2001). For this reason, obvious and precise instruction is delivered to clarify the linguistic features and rhetorical patterns in relation to the particular genre students attempt to take part in as well as purposes, organizations, and linguistic features of genres (Hyland, 2004a). The aim of this instruction is due to the fact that writing is an application on the basis of expectations, and the probabilities of a reader to understand the writer's aim would be better when the writer pay attention to predict the expectations of the readers regarding the type of writing (Hyland, 2007a). Therefore, teaching the genres and their organization explicitly allows writers to produce texts in the required way (Hammond, 1987).

2.1.1.6. English for Specific Purposes

The roots of English for Specific Purposes can be traced back to the end of World War II (Hutchinson & Waters, 1987). However, it was not until English appeared as the language of globalization that ESP came into the forefront of the educational agenda. Since English was regarded as the new lingua franca, people all over the world tried to leverage English to overcome problems they faced in several areas such as business, dissemination of knowledge, information sharing, and cultural communication (Teodorescu, 2010). Furthermore, the economic developments and sociopolitical conditions after the World War II required nations to raise well-educated and qualified people to compete and communicate in the international arena, which resulted in an increasing amount of international students with diverse needs in English speaking countries (Dudley-Evans, 2001; Hutchinson & Waters, 1987). To address these problems, ESP emerged as a scientific discipline.

English for Specific Purposes emerged as a response to the process approach. Contrary to the process approach, which emphasized the role of the writer, ESP prioritized the role of the reader, more specifically, the role of the academic discourse community (Silva & Leki, 2004). The writing was seen as an essential medium for entering disciplinary communities, and thus the focus in writing instruction was on the academic discourse genres and academic tasks (Tardy & Jwa, 2016). The first ESP teaching materials were mostly based on scientific

and technical writing, prioritizing technical vocabulary. Therefore, the main aim of ESP is “to teach the technical vocabulary of a given field or profession” (Smoak, 2003, p. 23). Consequently, ESP tries to make learners to be familiar with the general forms in tertiary education by creating realistic university writing conditions and informing learners of the genres specific to disciplines. In this regard, ESP relies heavily on discourse analysis, whose focus is on communicative features rather than grammatical or lexical aspects of writing and text analysis, which are mostly research papers (Johns, 2013; Maleki, 2008). Described as a learner-centered approach, ESP defined the writer's role as the one writing to satisfy the demands of the academic discourse community to achieve academic success. ESP is usually divided into two main categories, which also have their sub-categories: English for Occupational Purposes and English for Academic Purposes (Figure 3).

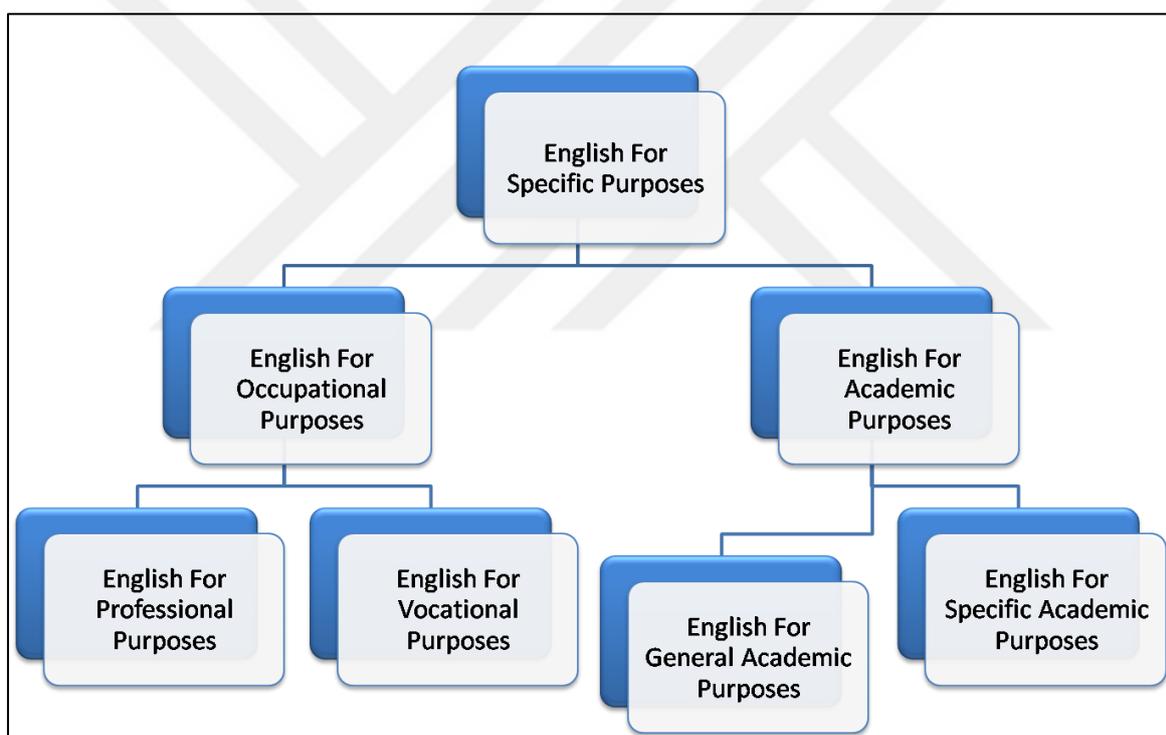


Figure 3. Categories of ESP

Throughout the decades, there has been an increasing interest in ESP, paving the way for a large body of scientific studies. These studies have implied a set of characteristics of ESP and pedagogical implications. Dudley-Evans & St John (1998) defines the distinctive features of ESP as follows:

- ESP is described to satisfy particular needs of learners;
- ESP benefit from the original methodology and activities of the field it deals with;
- ESP focuses on the language (grammar, lexis, and register), skills, discourse, and genres suitable to the teaching materials.
- ESP might be linked to or intended for particular disciplines;

- ESP might make use of a separate approach from that of general English instruction;
- ESP tends to be designed for intermediate or advanced adult learners. (pp. 4-5)

2.2. English for Academic Purposes

Academic writing refers basically to any kind of writing students need to achieve the tasks in a college or university. Although these tasks may differ widely, from writing an essay to preparing a research proposal, they all have a similar aim and criterion. Students entering higher education institutes have to show “the ability to use discipline-specific rhetorical and linguistic conventions to serve their purposes as writers” (Berkenkotter, Huckin & Ackerman, 1991, p. 19). Therefore, academic writing involves representing an understanding and showing expertise regarding a specific discipline (Irvin, 2010, p. 8). In that sense, academic writing has some differences compared to other types of writing, and has its own body of standards and applications (Oshima & Hague, 2006).

The main distinction of academic writing lies in three features: audience, tone, and purpose. First, academic writers address specific audiences who are most likely the instructors or lecturers in the university to convince that they have an in-depth knowledge of the topic or subject under investigation. In addition, academic writers should pay attention to tone in their writing. They should adopt a formal and impersonal tone, usually achieved through the use of the third person and ignoring subjective judgments, with highly technical and specialist vocabulary (Hyland, 2006; Nga, 2009). In addition, structures and sentences used in academic writing are more complicated than in spoken English. Thus, it is expected to involve more complement clauses, subordinate clauses, passive voice, attributive adjectives, and extended sequences of prepositional phrases (Biber, 1991). Finally, the purpose of academic writers should be precise and clear. The primary aim of the writers is to convince the readers that the standpoints they adopt are decent and proper. To do that, writers should be explicit both in the organization of the text and the ideas to assist the readers in comprehending the text (Biber & Gray, 2010).

International academic writers should learn rhetorical features such as writing practices and processes in order to produce texts of high quality. Unlike other texts, academic writing has a distinguishing structure, which is usually comprised of an introduction, a body, and a conclusion. Giving the first impression and informing the readers of the scope of the paper, the writer announces the topic and the motivation of the essay in the introduction part (Awelu, 2011). These parts should include the outline of the paper, the conditions for the body, and the definitions of key terms (Jones, 2015). The body, the most important part of

the essay, should be separated into progressive coherent, and cohesive paragraphs whose role is to support the topic of the paper. The paragraphs should also be logical and argumentative, and they should present evidence and examples (Whitaker, 2009). Finally, the results of the paper, solutions regarding the problem under investigation, and suggestions are presented in the conclusion part, where the main points are also summarized. What is vital in this part is that the writers should avoid repeating themselves (Anderson & Poole, 2009).

Academic writing plays a significant role in students' lives in academia as they have to spend a great deal of time writing papers during their academic life (Fukao & Fujii, 2001). For that reason, informing students of the necessities of good academic writing is quite important (Leibowitz, 2000). In this vein, several scholars have tried to identify the characteristics of academic writing and the points academic writers should focus. One of them, Thaiss and Zawacki (2006) synthesize the following characteristics in their comprehensive study in which they investigated the perspectives of university lecturers regarding academic writing:

- The writing should include strong evidence of being purposeful, willing to consider new ideas, and including technical knowledge.
- Priority should be given to the reason over sentiment or sensual perception.
- The piece of writing should address readers who are logical, seeking information, and aiming to elicit a rational response. (pp. 6-7)

Similarly, Thonney (2011) argues that students should produce their writing according to the following criteria:

- respond to what others have written about the topic they are concerned,
- state the value of their work and announce the plan for their papers,
- acknowledge that others might disagree with the position they have taken,
- adopt a voice of authority,
- use academic and discipline-specific vocabulary,
- emphasize evidence, often in tables, graphs, and images. (p.348)

Furthermore, Jalongo (2002) claims that academic writers need to consider the followings:

- Use specialized terminology and jargon in the discipline.
- Avoid ambiguous statements.
- Omit clichés and use expressions of your own.
- Use concrete details to highlight the main points and put forward ideas.
- Avoid markedly long or technical writing, unnecessary reiteration, and ambiguous language. (p.82)

Finally, Whitaker (2009) urges the academic writers that they should show considerable expertise in thinking, reasoning, researching, and evaluating abilities instead of describing the whole things they know.

2.3. English for Research Publication Purposes

The pressure to publish, especially in English-medium journals, has become a critical part of academic life for academicians (Lillis & Curry, 2010). The number of academic members and research students willing to secure publication in the international arena, which is overwhelmingly dominated by English, is increasing. The growing interest in international scholarly publication has resulted in the emergence of a new field of research: English for Research Publication Purposes (ERPP) (Flowerdew, 2015). In the special issue of the *Journal of English for Academic Purposes* devoted to ERPP, Cargill and Burgess (2008, p. 75) define ERPP as “a branch of EAP addressing the concerns of professional researchers and post-graduate students who need to publish in peer-reviewed international journals.” It has appeared that the dependency on English in academia poses severe challenges for especially multilingual scholars whose L1 is not English (Cho, 2004; Flowerdew, 2013; Lillis & Curry, 2010). In that sense, the aim of ERPP is to investigate the issues that the spread of English in the academic world has brought up. Although the term ERPP has been coined recently, the attempts to identify the multilingual scholars in academic writing for publication date back to the late 1980s (Uzuner, 2008). Since then, the research on the issue has flourished, and now ERPP is considered a distinct field of research.

2.3.1. Types of Publications

The publication is one of the most critical requirements in the modern academic world. Although publication in its general term covers a wide range of types, academic publishing, whose aim is to disseminate scientific knowledge and research, deals with a limited number of types. Over the decades, academics have gained a wealth of experience regarding what to publish and how to publish. Thus, each academic discipline has developed its own standards and norms for publication. However, there are some aspects that are accepted as general throughout the disciplines. This section aims to introduce four main types of academic publications, which are book, research article, conference proceedings, and thesis, and their common characteristics.

2.3.1.1. Book

Books may vary in their shapes, sizes, and purposes. Hartley identifies types of books as follows: popular science books, edited books, conference collections, handbooks, and

textbooks. Within ERPP framework, special focus is given to academic textbooks. According to Bhatia (1998), textbooks “make accessible established knowledge in a particular discipline to those readers who are being initiated into a specific disciplinary culture” (p.17). Nevertheless, regardless of the types, there is a general procedure for writing and publishing a book.

First, scholars seeking publication should ponder a suitable publisher. Some publishing houses announce the types of books and topics they publish. Therefore, it should be better to check the publishers first to have an idea of the market. Then, scholars had better send a letter to editors of the publishing house to state their intends to write or publish a book and question the suitability of the purpose and topic. At that point, publishers tend to demand a proposal, including an outline or sample chapter(s) of the book. Publishers usually wish to publish books that are “of high quality, original, with no or few competitors” in order to ensure that they will sell the books they will publish in the market easily (Woods, 2005, p. 148). If the editors consider the proposal to be worth publishing, the proposal is gone through the review process. Finally, on the basis of the reviewer’s’ comments, the final decision as to acceptance or rejection of the proposal is made (Cargill, Charvat & Walsh, 1996).

2.3.1.2. Research Article

In his seminal study, Swales (1990) introduced the research article genre. Since there is a wide variety of disciplines in the academic world and the requirements of these disciplines differ, it is not much possible to make a specific definition or research article. According to Swales, a research article is:

a written text (although often containing non-verbal elements), usually limited to a few thousand words, that reports on some investigation carried out by its author or authors. In addition, the RA will usually relate the findings within it to those of others and may also examine issues of theory and/or methodology. It is to appear or has appeared in a research journal or, less typically, in an edited book-length collection of papers. (Swales, 1990: 93)

Swales proposed that scholars write research articles to publish them in peer-reviewed journals, and this means the research article should undergo a reviewing process to be accepted for publication. Although the organization of the research article may differ from journal to journal and one discipline to another, a frequently used organizational model is the hour-glass model (Hill, Soppelsa & West, 1982) (Figure 4).

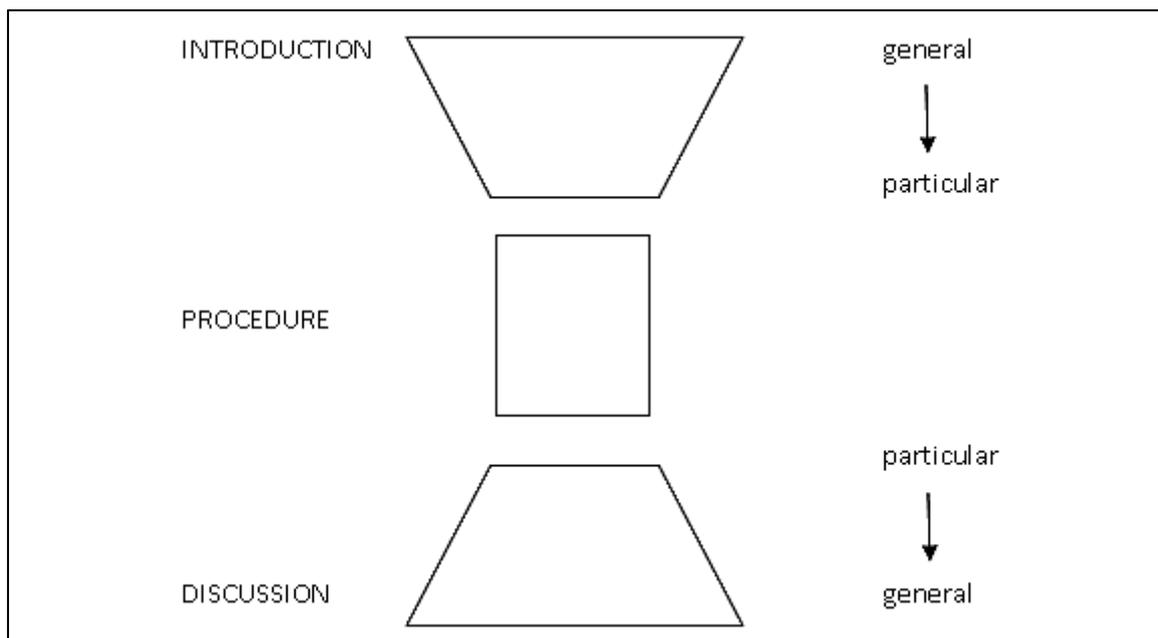


Figure 4. The hour-glass model of article structure

According to this model, a shift from a general topic to the particular focus of the study should take place at the beginning. Then, at the end of the research article, another shift from the findings to implications should occur. However, it should be remembered that this model presents a macrostructure and may not be suitable for and used by all disciplines (Swales, 2004).

Another structure for the research article is IMRAD (or IMRD) model, which is usually used in experimental papers. This model comprises of four main sections: Introduction, Methods, Results, Discussion. However, there are also different alternatives as to the structure of the IMRAD model (Glasman-Deal, 2010). Table 1 presents different variations of the IMRAD model.

Table 1

IMRAD Model

Introduction			
Methods			
Results or Data Analysis	Results or Data Analysis	Results and Discussion	Results or Data Analysis
Discussion	Discussion	-	Discussion and Conclusion
Conclusion	-	Conclusion	-

In the introduction section, writers are expected to introduce the topic of the study, identify the gaps in the literature, present the significance of the research, and finally outline the aim of the study. The function of this section is to convince the reader that the study is noteworthy and significant. In the methods section, the writers describe the method(s) they used in the study. The results section is where the finding of the study is presented. Typically, visuals and graphs such as tables, figures, charts, and diagrams are used to show the results of the study. Finally, the discussion section is where the findings of the research article are discussed and related to the existing literature in the discipline.

2.3.1.3. Conference Proceedings

A conference paper is the type of text that is presented at a conference. It can be defined as ‘the essential launching pad for nearly all scholarly careers’ (Gould, 1995, p. 37). Conference proceedings, on the other hand, is the compilation of papers presented at a particular conference. Conference papers introduce the latest trends and developments in a specific academic community. Therefore, conference papers and proceedings serve as a valuable source of up-to-date inquiries and tendencies in a discipline. The literature has revealed that at least a third of the papers presented at a conference become published papers in less than two years (Drott, 1995; Stolk, Egberts & Leufkens, 2002; Weller, 2002). Conference papers can be published in different ways: as books, articles, special issues, or full textbooks.

Conference papers are a means of engaging with the current research before the publication of the study. Therefore, a typical conference paper includes the same sections as those in the research article. However, conference papers are less formal than research articles. The main aim of the conference paper is to introduce new ideas and practices or present a study that is not yet finished. Conference paper allows especially novice scholars to connect with more knowledgeable peers and thus help them become part of networks in the discipline.

2.3.1.4. Thesis and Dissertation

A Ph.D. dissertation or a MA thesis has distinctive features in academia; that is, they are the first step of novice scholars into the academic community. Therefore, these kinds of writing should be well written, indicate professional knowledge, critical thinking, thoroughness, and scientific value (Swales, 2004). The Ph.D. dissertation is mainly written in two formats: monograph format and article compilation format. The former, written up as a comprehensible, concise text, is the more traditional type which is commonly used in soft sciences, while the latter preferred mostly in hard sciences comprises of previously published or publishable texts which are piled up as to form a coherent text (Swales, 2004). Similar to other types of publications, there are some disciplinary variances in the structure of Ph.D. dissertations. However, one format, which is called traditional structure, dominates thesis writing. The conventional model is an enlarged version of the traditional IMRD model adopted in research articles (Figure 5).

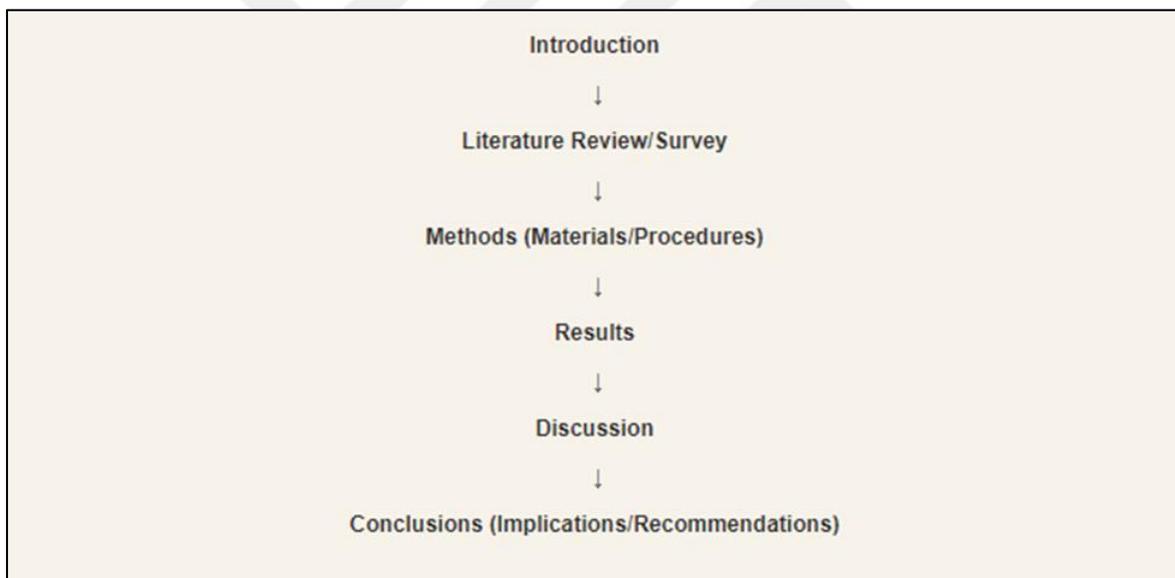


Figure 5. Conventional IMRD Model

Dissemination of doctoral studies is promoted in the academic communities. Even most of the universities and intuitions require doctoral studies to publish an article derived from their doctoral studies as a prerequisite for graduation. Therefore, Ph.D. students should write their dissertations in a style that is suitable to be published or be convertible into a publishable text.

2.3.2. The Theoretical Frameworks in ERPP

Throughout its history, the research on ERPP has used the following frameworks: genre approach, discourse analysis, social constructivism, situated learning theory, and social capital theory.

2.3.2.1. Genre Approach

Introduced by John Swales (1981) with his famous Creating a Research Space (CARS) model, the term genre has been widely used in the field of second language writing since the 1980s (Paltridge, 2014). Although a number of definitions have been produced for this concept, it basically refers to “a term for grouping texts together, representing how writers typically use language to respond to recurring situations” (Hyland, 2008b, p. 44). In addition, Bhatia (2004) defines genre as:

language use in a conventionalized communicative setting in order to give expression to a specific set of communicative goals of a disciplinary or social institution, which give rise to stable structural forms by imposing constraints on the use of lexico-grammatical as well as discursal resources (p. 23).

In this sense, Tardy (2011) lists the general characteristics of genre as follows:

- Genres are primarily a rhetorical category,
- Genres are socially situated,
- Genres are intertextual, not isolated,
- Genres are carried out in multiple – and often mixed – modes of communication,
- Genres reflect and enforce existing structures of power. (p.55)

Three standpoints exist within genre analysis. These can be listed as: New Rhetoric school (Bazerman, 1988), Sydney School (Halliday & Martin, 1993), and ESP school (Swales, 1990). The first viewpoint, New Rhetoric, considers the use of language as “recurrent social actions” in a context of situation (Freedman & Medway, 1994, p. 3). According to the Sydney School, language use is connected to the context, and thus is regarded as a structured “social process” concerned with achieving a specific aim (Martin, 1992, p. 505). Finally, ESP school of genre analysis views language use as “a class of communicative events” (Swales, 1990, p. 58). In addition, expanding Swales’s ideas, Bhatia (2008) recently suggested that ESP school of genre analysis should surpass mere linguistic analysis, and proposed a multidimensional approach to genre theory. Taking into account socio-cognitive, socio-critical, textual (genre

as a structured process), and ethnographic (genre as a social action) perspectives, his approach aims incorporate socio-cognitive and sociocultural examination in order to obtain a more compound interpretation of how these perspectives play a part in the production of genres in separate disciplines.

The genre approach is based on the idea that writing is contextual; that is, social context shapes the writing, and thus different contexts require different kinds of writing (Badger & White, 2000). It also argues that members of a community can easily identify the similarities in the texts they are familiar with, and they are able to read, understand and write the texts using their past experiences and existing schemas (Thompson, 2001). The primary criterion for identifying genres is the communicative purpose, which is recognized and reinforced within a community (Hyland, 2003). Therefore, the genre approach strives for investigating how different genre accomplishes their communicative purpose. In that sense, the research on genre approach has primarily focused on the linguistic structures of the text, the rhetorical actions, and the social context in which the texts occur in order to help learners produce texts that suitably conform to native-speakers norms (Dudley-Evans, 2002). Therefore, this approach is not limited to the mere delivery of writing instruction to exercise structures or grammar; it also emphasizes context and audience (Hyland, 2004a).

The main aim of this approach in the EAP paradigm is to teach characteristics of specific academic genres, such as research proposals, grant applications, essays, and research articles, and investigate the linguistic and discursive features of the text. The motivation behind this aim is the idea that such an investigation might offer writers a model to produce their own texts. However, genres are notably different with regard to their disciplinary features, particularly in the means disciplinary knowledge is set about, the way arguments are put forward, and the type of proof and confirmation they regard well-formed and sound in the related discipline (Bhatia, 2004). Therefore, this should be taken into account that genres are *prototypes* that differ among disciplines and even among the samples within the very same genre rather than being *rule-governed patterns* (Flowerdew, 1993). The genre approach offers a framework for both scholars and students to investigate the practices of disciplines and to employ this understanding in teaching writing and material development (Hyland & Hamp-Lyons, 2002). Typical genre analysis methods consists of text analysis, move structure analysis, comparative genre analysis, diachronic genre analysis, genre system analysis, and critical genre analysis (Tardy & Swales, 2014).

Of the all identified genres and genre traditions (for a detailed review, see Hyon, 1996), ESL tradition and academic writing genres have attracted considerable attention in second language writing (Hyland, 2002a). Consequently, numerous studies have investigated forms, linguistic features, and rhetorical structures of genres from several texts, genres, disciplines or language backgrounds. Much of the research on ESL and academic writing have focused on discourse analysis aiming to identify the rhetorical actions in the texts (Paltridge, 2014). The topics investigated in this research paradigm include, but are not limited to, research articles (Swales, 1987), functions that grammatical structures serve in academic writing (Harwood, 2005), the differences between L1 backgrounds of the writers (Jaroongkhongdach, Watson Todd, Keyuravong, & Hall, 2012), and the differences among different languages (Martín & León Pérez, 2014).

Table 2

CARS Model

Move	Step
Move 1: Establishing a territory	Step 1: Topic generalizations of increasing specificity
Move 2: Establishing a niche	Step 1A: Indicating a gap or Step 1B: Adding to what is known Step 2: (optional) Presenting positive justification
Move 3: Presenting the present work	Step 1: (obligatory) Announcing present research descriptively and/or purposely *Step 2: (optional) Presenting RQs or hypothesis *Step 3: (optional) Definitional clarifications *Step 4: (optional) Summarizing methods Step 5: (PISF**) Announcing principal outcomes Step 6: (PISF) Stating the value of the present research Step 7: (PISF) Outlining the structure of the paper

2.3.2.2. Discourse Analysis

The term Discourse Analysis (DA) was introduced by Zellig Harris in 1952 and refers to an approach to language analysis that investigates language patterns across texts and the social and cultural context in which these texts were produced (Paltridge, 2012). Discourse needs to be arranged systematically in a coherent manner that is practicable in an interactional context. Therefore, discourse analysis considers texts as wholes instead of treating them as sentences. It is argued that particular situations call for particular ways of language use, and these typical situations have characteristic linguistic features and share certain meanings (Flowerdew, 2014). In this vein, DA aims at examining these meanings and their recognition in language. DA practitioners focus on how people organize their speech or writing by

investigating the order that they typically use, which varies according to cultures (Gee, 2014). Used as an umbrella term, discourse analysis in ERPP uses the genre analysis, corpus-based discourse analysis, and contrastive rhetoric as the foundations for research (Flowerdew, 2013). The underlying assumption for using DA in ERPP research is that having more knowledge about the structure of the RA will facilitate the teaching of academic writing for publication.

Since the seminal work of Swales on the introductions of RA (1981), most of the studies on ERPP have focused on examining the moves in research article (RA) genre in linguistic and rhetorical terms (Dudley-Evans, 2002). Not only has the introduction section further been investigated (Nwogu, 1989; Paltridge, 1994; Samraj, 2002), but also have the studies extended to other parts of RA, that is, the discussion sections (Dudley-Evans, 1994; Lewin & Fine, 1996; Peacock, 2002), the abstracts (Hyland & Tse, 2005; Salager-Meyer, 1992; Santos, 1996), the results section (Kanoksilapathan, 2007; Thompson, 1993; Yang & Allison, 2003), metadiscourse (Hyland, 1998), politeness (Myers, 1998) and the citation practices (Dahl, 2004b).

Another strand of research within DA in ERPP addressed variation in academic, and professional or contexts (Bhatia, 2004). In this context, putting forward that different types of texts possess different linguistic features and thus each kind of text represent systematic variation patterns, Biber (2012) suggested that this variation should be examined under the heading of register, which he described as “text varieties of a language associated with particular situations of use” (p.191). In this sense, register consists of three main elements: the situational context which deals with the depiction of the conditions in which text is produced and received, the typical linguistic features which comprises all of the grammatical and lexical characteristics that are distinctive of the variety of the text, and the functional relationships between the first two components (Biber & Conrad, 2009, p. 6–11). Register analysis studies have revealed significant variation in genres within academic and professional domains (Bhatia, 2004).

Other studies have been conducted within the contrastive rhetoric (CR) paradigm, which mainly concentrates on the differences between L1 and L2 (Ahmad, 1997; Burgess, 2002; Melander, Swales & Fredrickson 1997). Within L2 writing research, the contrastive rhetoric appeared as “an area of research that identifies problems in second language writing and attempts to explain them by referring to the rhetorical strategies of the first language” (Connor, 1996, p. 5). The main idea of CR studies in ERPP paradigm is that multilingual

writers transfer some characteristics of RA in their L1 into L2 writing, and thus teaching should focus on these possible areas of transfer. A significant example of CR with Turkish participants is Uysal's (2008) study. She examined 18 Turkish writers texts to find out whether common rhetorical pattern existed between their L1 and L2 written production. The findings revealed a bidirectional transfer of some rhetorical preferences from Turkish. However, it was concluded in the study that a number factors such as participants' L2 proficiency, the educational context, topic, and audience may be the reason of the transfer.

2.3.2.3. Social Constructivism

Contrary to discourse analysis, which investigates the product, social constructivism (for detailed information, see Bruffee, 1986) is applied when the process of writing for publication is investigated. The Social Constructionist position regards entities such as knowledge thought, and reality as the products of communities of like-minded peers (Kukla, 2000). Within social constructivism, writing is seen as a social practice that consists of an implicit/explicit written communication between the writer and the reader in a particular community (Nystrand, 2006). In this sense, according to social constructivists, the primary purpose of writing for an individual is to become a new member of a particular community. Accordingly, Nystrand (1982) put forwards that “the special relations that define written language functioning and promote its meaningful use . . . are wholly circumscribed by the systematic relations that obtain in the speech community of the writer” (p. 17). In addition, Atkinson (2003) argued that when teaching, learning and the written language are associated in order to carry out various sorts of socio-cognitive activities, this enhances the field second language writing and makes it more comprehensive, profounder and more pertinent. Hence, second language writing should concentrate more on the implicit cultural and social practices functioned among different academic and sociocultural contexts (Atkinson, 1997). In this context, novice writers should have a sufficient amount of knowledge regarding to socio-cultural features of the particular community they are likely to encounter with.

In order to become a member of the community, writers need to: i) learn the nature and requirements of writing, ii) learn the ways to get access to and continue the relationship with the pertinent networks, and iii) know the rhetorical discourse of these particular communities (Flowerdew, 2013). Therefore, the social constructivist approach to research into academic writing for publication deals with the publishing process features such as negotiating with editors and reviewers, and acceptance or rejection of the manuscripts (Gosden, 1992, 1995,

2003) since academic writing for publication is seen as “the process by which novices are socialized into the academic community, the recognized route to insider status” (Hyland, 2007b, p. 88) due to the fact that providing feedback to novice writers requires them to adopt the perspective of readers, which, in turn, may possibly increase awareness of audience (Nystrand, 1986).

2.3.2.4. Situated Learning Theory

Similar to Social Constructivism, Situated Learning Theory considers writing as a situated social process. The basic assumption is that learning occurs in *Community of Practices* (CoP) (Lave & Wenger, 1991) by involving in ongoing activities and practices. Wenger (1998) mentions that a community of practice comprises three aspects: (1) mutual engagement, (2) shared repertoire, and (3) joint enterprise. Mutual engagement signifies that members in the community take part in activities and practices and comprehend them with others. The shared repertoire comprises means of performing things that members in the community have taken on and practiced. Finally, joint enterprise denotes “the result of a collective process of negotiation that reflects the full complexity of mutual engagement” (Wenger, 1998, p. 77). Novices in a community of practice are either on or outside the periphery. They thus are considered as legitimate and peripheral participants in the quest for full participation in a community of practice (Lave & Wenger, 1991). Therefore, the participants in CoPs interact through an apprenticeship with more experienced peers and experts, share information and experiences, and learn from each other (Cho, 2005; Flowerdew, 2013). According to Lave and Wenger (1991), newcomers in a particular field use legitimate peripheral participation (LPP) to act in a centripetal manner in their CoP, in the end adopting the roles of more experienced members. In that sense, novices obtain knowledge, adjust their inventiveness, and finally build a distinctive style, skills, identity, and discourse by means of this relationship with experts. In this aspect, ongoing participation in the practices of the communities allows newcomers to progress from peripheral to full participation (Uzuner, 2008).

Emerged from the notion of cognitive apprenticeship, the underlying rationale of LPP is that newcomers primarily involve in low-risk and relatively simpler peripheral tasks that are, however, legitimate and productive for the objectives of the community (Flowerdew, 2000). Through the peripheral tasks such as writing an introduction or literature review or preparing a research paper for an examination, novice participants learn more about the nature of the

community and gradually begin to be more centrally located within the community's social practices (Armstrong, 2015). Thus, rather than considering knowledge as a product, it is seen as a process, which becomes valid when activated within the specific community (Korthagen, 2010). Therefore, this approach emphasizes the importance of the support for writing for publication, drawing attention to the critical role of experts, who are mainly academic members and supervisors, in scaffolding, enculturating, and mentoring novices in their journey to full participation in the target communities, becoming experts in the field.

In an attempt to provide an alternative explanation for communities of practices, Swales (1990) put forwards the notion *discourse community*, referring to a group of people who share a set of social practices that are directed towards some purposes. In this sense, academic discipline can be considered as an academic discourse community. A discourse community is heterogeneous, heteroglossic, multivoiced, and multidimensional; that is, different members of a discourse community may concentrate on different layers of the discipline through discussing scientific results and theories (Bazerman, 1992).

Swales (1990, p. 24–27) lists six criteria for a discourse community: i) common goals, ii) participatory mechanisms, iii) information exchange, iv) community-specific genres, v) a highly specialized terminology, and vi) a high level of expertise. Writing, especially writing for publication as a primary means of disseminating knowledge and providing feedback, is regarded as the primary way of communication among the members of a discourse community since it may have members from a wide range of geographical regions (Pecorari, 2002). In order to become successful members of particular discourse communities, individuals need to learn the social practices that support these criteria (Flowerdew, 2000). Therefore, participation in a discourse community depends on the newcomers' abilities to satisfy these norms.

2.3.2.5. Social Capital Theory

In his seminal capital theory (1986), French sociologists Pierre Bourdieu stressed the corporal characteristics of social life and drew attention to the function of practice and personification in social life. He proposed three forms of capital, namely economic, cultural, and social. Bourdieu defines social capital as “the aggregate of the actual or potential resources which are linked to possession of a durable network of more or less institutionalized relationships of mutual acquaintance and recognition - or in other words, to membership in a group” (p. 248-249). According to him, social capital belongs to the

individual, not to the public. It allows one to apply the power on a group and is accessible for those putting efforts to get it through gaining power and status and establishing goodwill (Bourdieu, 1986). In that sense, social capital theory shares many common features with Situated Learning Theory.

The social capital theory has recently attracted a great deal of attention in ERPP since it underlines the genuine social relations and systems that support an individual to access and use the resources (Stanton-Salazar, 1997). Within ERPP, it refers to the relationship with students, lecturers, supervisors, professors, and members of a community. Several researchers have used this notion to examine the influence of social capital and its forms on expert and novice scholars' academic success and development, investigating disciplinary enculturation (Cho, 2009; Curry & Lillis, 2010).

The aforementioned notions provide insightful frameworks for understanding the complex nature of academic writing for publication. Research and publication communities of scholars in specific academic fields are good examples of communities of practice or discourse communities. Communities of practice, by definition, differ in their organization and their practices; that is, the expected norms of writing in one discipline are probably different in another (Ho, 2017). In addition, it seems that academic writing for publication takes the form of legitimate peripheral participation since novice scholars learn by actively participating in the process of writing (Flowerdew, 2013). They write research articles, submit them to journals, and receive feedback from the editors and reviewers, through which they become members of their discourse community.

2.3.3. Problems of Multilingual Scholars in Writing for Publication

The problems and challenges faced by multilingual scholars in academic writing for publication process have long been a topic of discussion. On the one hand, some scholars argue that Native English speakers have an advantage since they acquire the language during their childhood naturalistically while multilingual users have to invest time, effort, and money to learn English, which may potentially build barriers to writing in English (e.g., Flowerdew, 2008). On the other hand, others claim that the challenges faced by multilingual scholars are also shared by the novice native speakers since writing is a social practice through which novice native writers are also socialized (e.g., Hyland, 2016). Regardless of the result of the debate, it is clear that multilingual scholars have some specific issues and challenges in academic writing for publication. These problems will be discussed below.

2.3.3.1. Narrow-mindedness

Narrow-mindedness refers to a limited perspective focused on the local area. This problem, also known as parochialism, occurs when multilingual scholars fail to link their studies into the international research context. In order to be published in high-ranked journals in international indexes, a research article is required to be relevant to and contribute to the center-based international community since too localized research may fall short behind meeting the expectations of the discourse communities in the Center who consider such studies as less pertinent and persuasive (Flowerdew, 2001). Therefore, parochialism is stated as one of the reasons that lead to the rejection of the paper.

A journal editor in Flowerdew (2001) expressed that contributions from peripheral locations showed a tendency to be too localized, and multilingual scholars failed to indicate the significance of their study and how their study would contribute to the international community. It is also argued that multilingual scholars have to explicitly indicate that their locally-based research will contribute to the international community of scholarship and convince members of the center-based scholarly communities that their study goes beyond the local context in terms of their results and significance (Curry & Lillis, 2004; Duszak & Lewkowicz, 2008). It is argued that the main reason for this problem is the lack of adequate presentation of the topic instead of the content since local concerns may provide insights for the members of the community to design future studies (Belcher, 2007).

2.3.3.2. Rhetorical Problems

As is discussed in the contrastive rhetoric section, all cultures have their distinctive writing styles and rhetorical conventions (Connor, 1996). However, harsh protective applications for English scientific conventions exist in writing for publication in academic journals, especially in high-ranked ones. It is a well-known fact that editors and reviewers are relatively intolerant of deviation from the standards of academic writing rhetorics. Such an attempt is treated as inaccurate and finally results in the rejection of the paper. Even journals and editors who try to criticize the governing paradigms and publish such articles experience extensive confrontation from the discourse community. Therefore, multilingual scholars consistently find themselves in a position where they have to conform to center-based rhetorical norms (Bennett, 2011). This absolute prerequisite requires multilingual scholars to become familiar with the Anglophone rhetoric and write their research papers according to this norm (Berkenkotter & Huckin, 1995).

Anglophone rhetoric demands scholars to adopt an exegetical and critical approach and to organize their manuscripts in a dialogical manner while writing their research papers (Uzuner, 2008). Because of the differences in rhetorical practices and patterns across languages and cultures, multilingual scholars face challenges in such rhetorical issues as the structure of a research article, clarity of ideas, and the degree of metadiscourse marking (Flowerdew, 1999b). The literature has revealed that the stylistic difference between multilingual scholars' cultural values and rhetorical practices and those of the English speaking discourse communities frequently bring about rhetorically unsound or weak texts. In this regard, English academic rhetoric constitutes an important source of struggle and disadvantage for multilingual scholars.

Anglophone rhetoric is based on a phenomenon called creating a research niche (Swales, 1990). In this tradition, writers have to follow certain rhetorical moves to write successful research articles. Introduction, literature review, and discussion sections are the main areas where these rhetorical moves should be applied, and thus multilingual scholars face serious challenges in writing these parts (Flowerdew, 1999; Li, 2007; Shaw, 1991; St. John, 1987). The research showed that the stylistic differences due to the differences across cultures rather than the linguistic barriers are the main reasons for the rhetorical problems of multilingual scholars (Flowerdew, 2001; Li, 2002; Li & Flowerdew, 2007; Swales, 1990). The cultural values of multilingual scholars seem to prevent them from using argumentative strategies necessary to persuade the readers of the significance of their research, which leads to failure to conform to the mainstream international community (Cho, 2004; El Malik & Nesi, 2008; Tardy, 2005). When local styles that differ from the mainstream rhetorical structure are identified, the manuscripts of multilingual scholars tend to be rejected since these differences hinder them from prioritizing the value, significance, and credibility of their work (Curry & Lillis, 2004). Therefore, lack of knowledge regarding the rhetorical moves in scholarly writing has a negative impact on the ability to produce acceptable papers and consequently decreasing the publication rates in the academic community (Bazerman, 1985; Swales, 1990).

2.3.3.3. The Nature of Writing for Publication in English

Writing a research paper for publication is a burdensome activity itself. The burden is even increased when this activity is carried out in a second language since multilingual scholars need extra time and effort for reading and reporting research results in a second language

(Flowerdew, 2008). In this vein, most of the studies have emphasized the tiresome and burdensome nature of writing for publication and have pointed to the additional burdens of multilingual scholars in this process (Casanave, 1998; Curry & Lillis, 2004; Flowerdew, 2000; Liu, 2004). The poor writing strategies multilingual scholars employ were cited in the literature as one of the reasons for the time-consuming nature of writing in English. It was found that most of the scholars devote an enormous amount of effort and time translating their manuscripts into English after writing them in their L1 first (Gosden, 1996; Li, 2007; St. John, 1987). Regardless of the techniques or strategies they employ, writing a research article in English is a challenging task for multilingual scholars. The time consuming and tiresome nature of writing for publication in English not only results in delays in publications due to constant reviews and corrections but also prevent multilingual scholars from engaging in new projects, which leads to a reduction in their motivation and academic productivity (Curry & Lillis, 2004; Flowerdew, 1999).

2.3.3.4. Lack of Connections with Core Academic Communities

Multilingual scholars, especially those in the periphery, usually get limited opportunities to compete for publication with those in the center discourse communities (Canagarajah, 2003). There are a number of restraints arising from geographically isolated positions of the multilingual scholar. Lack of connections with core academic communities in the center may impede multilingual scholars to publish in international high ranked academic journals, though having a limited network or relationship does not necessarily result in rejection of multilingual scholars' submissions (Belcher, 2007). For that reason, a number of studies have investigated whether having contact with members of center-based scholarly communities make publishing manuscripts in international journals easier. An earlier attempt was made by Casanave (1998). She investigated four bilingual Japanese scholars who returned to Japan after a graduate-level education in the US and concluded that scholars having connections with the center communities were more likely to attain publication in the center-based journals. Flowerdew (2000) also investigated the process Oliver, a recently returned doctoral graduate, underwent in trying to publish in English and found that the privileged position of Oliver compared to other multilingual scholars enabled him to publish his manuscript. In the same vein, Curry and Lillis (2004) proposed that being remote from the center scholarly communities was a limiting factor and likely to reduce the probability of involving fully in the international research network. Similarly, investigating the accepted

and rejected manuscripts submitted to English for Specific Purposes Journal, Belcher (2007) revealed that 83% of networked submissions were accepted for publication while only 24% of off-networked submissions were published. It may be argued in the light of this evidence that having limited or no connections with the center academic communities constitute another challenge for multilingual scholars seeking international publication.

2.3.3.5. Bias Against Multilingual Scholars

As a consequence of being isolated from and having a limited relationship, if they have any, with the center academic communities, diverse bias regarding multilingual scholars exists. Therefore, potential bias due to issues of linguistic background, race, affiliation, and class is another obstacle multilingual scholars have to overcome. The literature on writing for publication seems to focus specifically on the reviewer and editorial biases. On the one hand, some studies (Li, 2002; Liu, 2004) provided no evidence of prejudice against the submissions of multilingual scholars. On the other hand, other studies (Belcher, 2007; Cho, 2004; Flowerdew, 2001; Gosden, 1992; Hewings, 2006; Li, 2006) reported evidence supporting the argument that potential bias against multilingual scholars exists. For example, editors in Gosden's study (1992) stated that they did not take submissions from particular regions into consideration because of their bias. However, journal editors in other studies claimed that they treated each paper equally without discrimination against NNS manuscripts; scholars submitting manuscripts with severe grammatical mistakes or style or form improper to norms of the center-based communities are labeled as unintelligible or even stupid (Ammon, 2000; Curry & Lillis, 2004). In sum, the literature emphasizes the biased treatment of multilingual scholars' submissions and that most of the editors and reviewers request multilingual scholars to have their manuscripts checked by a native speaker, which is an impossible demand for most of the multilingual scholars due to lack of availability of native speakers in their regions and high cost of editing services (Li & Flowerdew, 2007).

2.3.3.6. Lack of Sufficient Funds to Conduct Research

Although Canagarajah identified a lack of sufficient funds to conduct research as one of the non-discursive challenges of multilingual scholars in 1996, little research was done to investigate writing for publication from an economic perspective. It was found that countries devoting more funds to research, which are mainly from Kachru's inner circle, dominated

the international publication (Man, Weinkaif, Tsang & Sin, 2004). The issue of lack of funding was also identified by Aydinli and Mathews (2000) as one of the explanations for the limited contribution of academicians in the field of international relations from peripheral regions to mainstream international journals. Moreover, Salager-Meyer (2008) proposed that the contribution of the private sector to research funding is twice as much as that of state funding in highly developed nations. In contrast, research is funded substantially by the states in developing countries. In short, it seems that research outside of Kachru's inner circle is less supported in terms of funds to conduct research, which decreases the possibility of carrying out groundbreaking studies as well as the possibility of publishing in high-ranked international journals (Swales, 1997).

2.3.3.7. Language Problems

It has always been argued that the domination of English in academic circles gives an outstanding advantage to scholars whose mother tongue is English over nonnative scholars. It is also believed that native speakers of English save a considerable amount of time during the writing process since they do not have to deal with the difficulties of writing in a second language and thus find writing for publication relatively easy compared to their nonnative colleagues (Hamel, 2007). Therefore, a large body of research has investigated the language problems of the multilingual scholars in writing for publication since writing for publication to become a member of a particular discourse community is a vigorous attempt. These studies have indicated that a lack of proficiency in English is a unique obstacle for international scholars.

International scholars face problems while writing or communicating with editors (Flowerdew, 2013). Evidence has revealed that lack of success in language demands of core academic communities may result in rejections of the manuscripts (Bordage, 2001; Duszak & Lewkowicz, 2008; Li, 2005). In this regard, Hewings (2006) found that international scholars received more adverse commentaries regarding the clarity of the sentences, word choice, grammar, syntax, and pronoun use. A vast majority of the participants in the surveys investigating the challenges and problems of multilingual scholars expressed that they suffered from language problems in writing and publishing in English (Cho, 2009; Ferguson, Perez-Llantada & Plo, 2011; Huang, 2010). For instance, the participants in Flowerdew's study (1999a) thought that technical problems with the language put them in a disadvantageous position. Similarly, Man, Weinkauf, Tsang, and Sin (2004) found a

relationship between mastery of English and publication rate. Other studies, aiming to identify the specific language problems in academic writing, have shown that multilingual scholars have problems with regard to having less facility of expression and a less rich vocabulary (Flowerdew, 1999a), grammar and use of hedges (Flowerdew, 1999b), complicated syntax and unclear modality (Flowerdew, 2001), improper and inaccurate use of idiomatic expressions (Kaplan & Baldauf, 2005), usage errors (Henry & Roseberry, 2007), problems in cohesion and coherence (Hyland & Salager-Meyer, 2008), difficulties in referencing and citing (Mišak, Marušić & Marušić, 2005), comprehending the comments and notes of the reviewers (Duzak & Lewkowicz, 2008), potential interference of the first language (Gosden, 1996) and semantic and syntactic errors (Mungra & Webber, 2010). In sum, the evidence supports the argument that writing in English places multilingual scholars at a disadvantageous position in which the multilingual scholars consistently try to overcome grammatical and lexical problems. Besides, these serious linguistic barriers lead many multilingual scholars to be excluded from the international research community (Ammon, 2012).

2.4. Linguistic Problems

A brief discussion above clearly shows that although international scholars have to overcome a large number of obstacles, there is also a range of strategies that they can employ during the writing for the publication process. However, the linguistic problems of multilingual scholars are particularly important for second language teaching research as the quality of the research paper is, to a great extent, related to producing an easily readable text with proper vocabulary, sentence, and paragraph organization (Richards & Renandya, 2002). Scholars are required to produce concise, sound, and cogent sentences and paragraphs satisfying both the grammatical rules of English and the norms of academic writing. The linguistic aspects of written academic English also need to achieve a rhetorical purpose such as nominalization, subject-verb number agreement, use of articles, verbal complementation, tenses, relative clauses, and cohesive devices such as conjunction, substituted phrases, synonyms, and lexical repetition (Halliday & Hasan, 1976). Appropriate sentence structure, accurate vocabulary, and appropriate rhetoric enables the reader to recognize the purpose of the text more eagerly. Therefore, scholars, especially multilingual scholars, need to acquire the necessary skills to produce texts with such features. Otherwise, their manuscripts will be labeled as inadequate and rejected. However, acquiring these skills is quite a burdensome

and time-consuming activity and can be achieved best through effective delivery of teaching. In that sense, second language teaching research has much to offer to multilingual scholars for overcoming the linguistic problems they face. Therefore, more studies focusing on language-related problems of multilingual scholars should be conducted.

2.4.1. Linguistic Accuracy

Accuracy, in its basic form, can be described as the precise language use comprising the correct use of vocabulary, pronunciation, and grammar. In this sense, Linguistic accuracy refers to “the ability to be free from errors while using language” (Wolfe-Quintero et al., 1998, p. 33). Linguistic accuracy constitutes a significant part of writing due to the fact that writers address the experts in specific disciplines. Therefore, they are expected to apply linguistic features and rules correctly to a reasonable degree to fulfill their academic and professional purposes (Celce-Murcia, 1991). There is an abundance of studies investigating linguistic accuracy in the literature, most of which have dealt with the effect of written corrective feedback on linguistic accuracy. In her review on the measures of accuracy, Polio (1997) divided the measures of accuracy used in the studies into three categories: holistic scales, number of error-free units, and number of errors. Studies using holistic scales to measure accuracy treated linguistic accuracy as a part of the components in a rating scale (Hamp-Lyons & Henning, 1991; Hedgcock & Lefkowitz, 1992; Tarone et al., 1993). On the other hand, some studies determined the quantity of error-free T-units and/or error-free clauses to measure accuracy (Casanave, 1994; Ishikawa, 1995; Robb et al., 1986). Finally, other studies counted the number of errors, either with or without classification, to measure accuracy (Fischer, 1984; Frantzen, 1995; Kroll, 1990; Zhang, 1987).

The literature in the written academic English covers all studies investigating the use of English in the academy. Therefore, this strand of research includes a wide variety of genres from student essays, lectures, the mailings between instructors and students to book chapters and research articles. In this context, four main research paradigms exist in the written academic English research: genre analysis, contrastive rhetoric, ethnographic studies, and corpora-based studies (Flowerdew, 2014). Of these, the corpora-based studies seem to be the most appropriate way to investigate linguistic accuracy since the corpus linguistics, especially those based on learner language, is quite appropriate to conduct quantitative analyses providing one of the most reliable sources of frequency, which is a feature of language that gives clues about what is possible and what probably occurs in the language

(McEnery & Wilson, 1997). The corpus linguistics also allows researchers to investigate large samples of real language data and thus leads to improved language descriptions and the way it works with a focus on lexical and grammatical patterning. Offering a new empirical aspect, the tools corpus linguistics use in research has transformed the language studies, enabling researchers to make stronger explanations, support their assertions and deal with language more confidently (Hyland, Huat & Handford, 2012). Especially, the introduction of Computer-aided Error Analysis (Dagneaux, Dennes & Granger, 1998) represented a breakthrough in terms of investigating linguistic accuracy. This methodology enabled researchers to analyze learner data in a more detailed and comprehensive way on account of the frequency and statistics generated as a result of the analysis. The researcher can use computer-aided error analysis either to investigate a particular linguistic item or structure or to examine all of the errors in the data.

Computer-aided error analysis studies investigating particular linguistic items have thus far proved that academic writing of international scholars is different from that by native speakers of English in a number of regards, such as syntactic and discourse features (Granger, 1997). For example, Granger and Tyson (1996) reported that nonnative writers overused or underused individual connectors and misused them semantically, stylistically, and syntactically in academic writing. In addition, Granger (1997) found that nonnative writers used fewer participle clauses than native writers and attributed this finding to a lack of stylistic features. Furthermore, Martínez (2005) investigated the use of first-person pronouns of international scholars in the field of biology and identified problems related to phraseological issues, overuse, and underuse. In addition, Bestgen and Granger (2011) investigated spelling errors in data obtained from 3 different L1 backgrounds and identified 1,614 spelling errors. They concluded that the number of spelling errors could be used to assess the quality of the texts produced by second language writers. English articles have also been found to be another structure international scholars have trouble with. For example, Barrett and Chen (2011) examined the article errors of Taiwanese EFL learners using a detailed tagging system. They reported overuse of both the definite and indefinite articles as well as underuse of the zero articles. Similarly, Crompton (2011) examined the article errors of advanced Arabic learners. He first generated frequencies and then compared them with articles in native English. Finally, he reported that the participants made the highest number of errors in the misuse of definite articles for generic reference. Carrió-Pastor and Mestre-Mestre (2014) investigated lexical errors using data consisting of 30 scholarly papers in the

field of engineering. They identified using a word with a similar meaning to another as the most frequent error type. In addition, López and Manalastas (2017) collected a corpus consisting of 60 papers and investigated the grammar, spelling, and punctuation errors of university students at British universities. They found that the participants made more errors in the punctuation category and less in the spelling category.

On the other hand, some researchers tried to investigate all of the errors instead of focusing on a specific item or structure. For example, McDowell (2016) analyzed research articles produced by 13 Japanese scholars in the field of materials science and identified 654 nominal group errors, concluding that Japanese scholars had difficulty particularly in articles, plural -s, and prepositions -of. Similarly, Chuang and Nesi (2006) examined an 88.000 words-written corpus of Chinese ESP learners. They marked up their corpus in the light of a tentative error tagging system they developed particularly for the study. They found that the highest number of errors was detected in the grammar, lexico-grammar, and lexical categories, respectively. They put forward in their study that linguistic features their participant had the most difficulty in were determiners, nouns, verbs, prepositions, punctuation, sentence parts, tense/aspect, modals, conjunctions, and pronouns, respectively. Furthermore, Lopez (2009) carried out a computer-aided error analysis study and tagged the Spanish part of the International Corpus of Learner English using the error types described in Dagneaux et al. (1996). The participants in her study were found to have troubles, especially with lexis, articles, and spelling, respectively. In addition, MacDonald (2016) collected written corpora consisting of data from 304 Spanish EFL students. She used an error tagging system with 170 features in total and identified 15,850 errors in the corpus. It was revealed in her study that the greatest number of errors was tagged in the grammar, lexical, and punctuation categories, respectively. More recently, McDowell and Liardét (2020) developed an elaborated framework for investigating errors in research papers. They collected a corpus consisting of 46,263 words and identified a total of 1,368 errors. As a consequence, they concluded that the participants had difficulty in complex nominal groups, articles, and plural -s.

However, it should be mentioned that the literature into academic writing regards texts written by native speakers of English as the unchallenged linguistic standard, simplifying the role of international scholars to mere learners and ignoring their professional role as writers seeking publication (Wu et al., 2020). Therefore, more studies addressing the professional status of international scholars are needed.

2.4.2. Linguistic Complexity

Started with the study of Wolfe-Quintero et al. (1998), the issue of complexity has drawn a great deal of attention in the field of L2 writing. This construct has been regarded as being quite important in second language studies since the improvement in complexity is considered as a significant indicator of the development in the second language as well as an effective and fundamental descriptor of performance (Bulte & Housen, 2014; Lu, 2010). Being an extremely multifaceted construct which is made up of a number of sub-constructs, components, and features, no clear agreement exists in the second language research literature regarding the definition of this term, which has resulted in confusion as to how to conceptualize complexity (Norris & Ortega, 2009; Pallotti, 2009). Figure 6 illustrates the components and dimensions of the complexity concept.

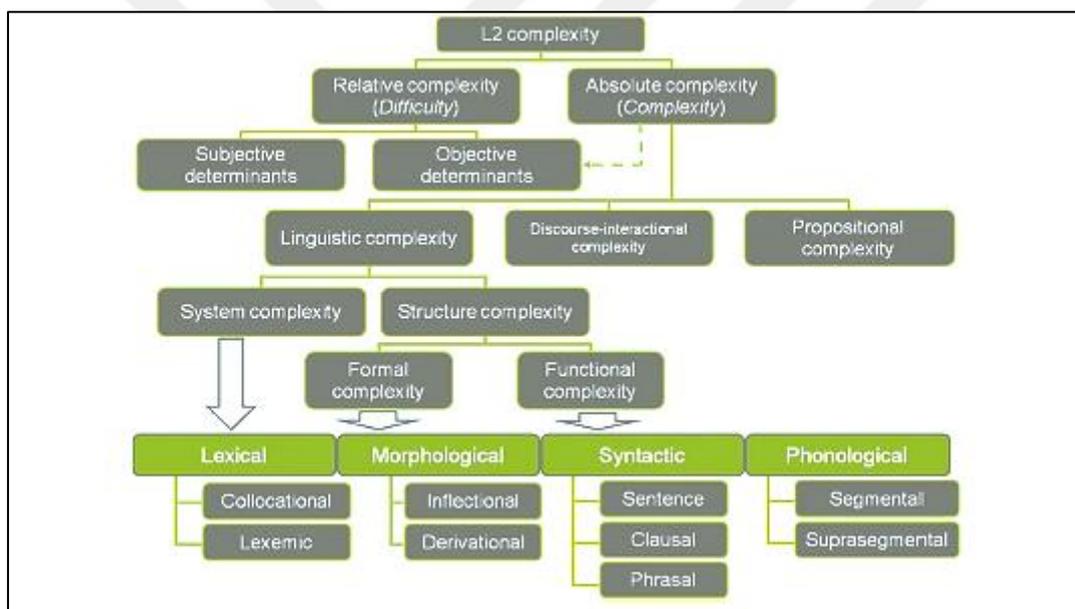


Figure 6. Taxonomy of Second Language Complexity

As a significant constituent of the linguistic complexity construct (Bulte & Housen, 2014), syntactic complexity denotes the variety of forms that appear in linguistic production units and the amount of sophistication of these forms (Ortega, 2003). In addition, Foster and Skehan (1996) state that greater use of syntactic features indicates a development in the language. Similarly, according to Wolfe-Quintero et al. (1998), syntactic complexity is “a manifest in writing primarily in terms of grammatical variation and sophistication” (p.96). These definitions emphasize three significant dimensions of syntactic complexity: elaborateness, variation, and sophistication. Elaborateness refers to the use of longer production units such as sentences and phrases, whereas variation is the use of a wide variety

of syntactic structures from basic to advanced. Finally, sophistication refers to the use of sophisticated forms on which, however, no clear consensus exists (Kyle & Crossley, 2017). In that sense, language development has been thought to be correlated with greater and more diverse use of syntactic features, making the analysis of second language production with regard to its syntactic complexity a widespread research area (Crossley & McNamara, 2014). In order to measure syntactic complexity, a wealth of indices has been offered. A number of researchers have tried to develop reliable and valid tools for measuring the syntactic development of second language learners (Lu, 2011; Ortega, 2003; Wolfe-Quintero et al., 1998;). A considerable amount of these indices are proposed on the basis of formulas, frequencies, and ratios (Norris & Ortega, 2009). These indices can be grouped under the following four broad categories: i) length of the production unit, ii) amount of subordination, iii) amount of coordination, and iv) degree of phrasal sophistication and range of syntactic structures (Ansarifar et al., 2018). However, a number of these measures have attracted researchers in the second language, especially second language writing, more than others (Ortega, 2003). Frequently used syntactic complexity measures will be briefly mentioned.

2.4.2.1. Frequently Used Measures in Syntactic Complexity

Mean length of sentence (MLS), mean length of T-unit (MLT), and mean length of clause (MLC). MLS basically refers to the number of words in a sentence. The ease of identifying a sentence, as it has a clear definition, makes this measure quite operationalizable and thus enables researchers to obtain reliable data quickly. MLS has been reported to be strongly associated with the mean length of T-unit (MLT) (Lu, 2010) and language proficiency (Ortega, 2003).

Put forward by Hunt (1965) for the studies in first language acquisition, the concept of T-unit, or minimal terminable unit, has been used in the field of second language acquisition research since the late 1970s (Larsen-Freeman, 1978). A T-unit can be defined as “one main clause plus any subordinate clause or nonclausal structure that is attached to or embedded in it” (Hunt, 1970, p. 4). T-unit has been the focus of a large number of researchers since then due to the fact that it has clear merits over MLS because appropriate use of punctuation is not a prerequisite for T-unit (Matthiessen & Halliday, 2009). A considerable amount of studies has found a positive correlation between proficiency and MLT, concluding that the proficiency increases with the mean length of T-unit (Ortega, 2003; Wolfe-Quintero et al., 1998;) though there has also been criticism towards this concept (see Bardovi-Harlig, 1992).

However, it has been used as an essential measure of syntactic complexity since its first appearance in language research.

The mean length of clause (MLC) basically refers to the mean of the number of words per clause, which can be described as a group of words involving a subject and a finite verb. It has been reported in the literature that there is a significant positive correlation between MLC and proficiency in a second language in that it is likely to go up when the proficiency level increases (Cumming et al., 2005). However, there exist some studies reporting results that challenge this finding (Knoch, Rouhshad & Storch, 2014).

The ratio of complex T-units per T-unit (CT/T) is another measure used for investigating syntactic complexity. A complex T-unit refers to a T-unit consisting of both an independent and a dependent clause (Lu, 2011). CT/T calculates the number of T-units, including dependent clauses. Although, studies have identified a positive tendency between this measure and development in language, statistically significant results have not yet been reported (Casanave, 1994; Lu, 2011).

The number of T-units per sentence (T/S) can also be used in the analysis of syntactic complexity. This index calculates the quantity of clausal coordination in a given text. In their review, Wolfe Quintero et al. (1998) reported that only Monroe (1975), who examined French a second language found a significant correlation between T/S and proficiency. The nature of this correlation was negative, which means that T/S increases when proficiency in a second language decreases.

Clauses per Sentence (C/S) is used to measure overall sentence complexity. It is a global measure that calculates the ratio of subordination and clausal coordination per sentence. The findings as to this measure seem to be rather controversial. Although Ishikawa (1995) identified a positive correlation between development in a second language, and this measure, a negative correlation between this measure and the school year was reported by Lu (2011).

Another commonly used measure is clauses per T-unit (C/T). C/T calculates the quantity of clausal subordination in a text; however, it does not differentiate types of subordination. Results regarding this measure are also conflicting. In their review, Wolfe-Quintero et al. (1998) reported that 18 studies used this measure up to that time. Significant positive correlation in 6 studies, significant negative correlation in 1 study, and no relationship in 11 studies were found between C/T and proficiency. Furthermore, recent attempts investigating

C/T reported that no significant differences were found between this measure and development in the language (Cumming et al., 2005; Knoch et al., 2014; Lu, 2011).

The number of dependent clauses per clause (DC/C) is another measure used to investigate syntactic complexity. Reminiscent of C/T, it calculates the quantity of clausal subordination. A brief literature analysis reveals that a negative correlation was identified between the school year and this measure in that a greater number of dependent clauses were used when the proficiency in a second language decreases (Lu, 2011).

Another commonly used measure is dependent clauses per T-unit (DC/T). Similar to C/T and DC/C, it also calculates the quantity of clausal subordination in a text. Conflicting results were reported for this measure as well. Homburg (1984) reported a significant positive correlation between language proficiency and DC/T whereas Lu (2011) reported a negative correlation between DC/T and language proficiency. In addition, Vann (1979) reported no relationship.

Coordinate phrases per clause (CP/C) can also be used in the analysis of syntactic complexity. This measure calculates the quantity of phrasal coordination in a text. Lu (2011) reported a positive correlation between this measure and proficiency, concluding that when the proficiency of learners increases CP/C increases as well.

Coordinate phrases per T-unit (CP/T) is used for examining syntactic complexity. This measure and CP/C are quite alike in that both measure phrasal coordination. The only difference is that CP/T is insensitive to phrase types in which the coordination occurs. CP/T and language proficiency were reported to be positively correlated with each other (Lu, 2011).

Another measure is complex nominals per clause (CN/C) CNs are very common in professional texts in English (Nakov, 2013). They refer to a wide variety of syntactic formations such as noun-adjective combinations, nominal clauses, and infinitives or gerunds as the subject (Cooper, 1976). A significant positive correlation was found between all levels in college education apart from years 3-4 with regard to this measure by Lu (2011).

Complex nominals per T-unit (CN/T) is among the most commonly used syntactic complexity measure. Although CN/T and CN/C are quite similar to each other, Wolfe-Quintero et al. (1998) argued that CN/C is more superior in that CN/T. Similarly, Lu (2011) found that CN/C performed better in explaining the differentiation of the students at different levels.

Other commonly used measures are verb phrases per T-unit (VP/T) and verb phrases per clause (VP/C). The use of VP/T in measuring syntactic complexity was put forward Wolfe-Quintero et al. (1998). It calculates the overall number of verb phrases, both finite and non-finite, in a T-unit. VP/C operates similarly, but it calculates the total number of VPs in a clause. However, no relationship was found between VP/T and proficiency in a second language (Lu, 2011).

Finally, last measures are passives per T-unit (P/T), passives per clause (P/C), and passives per sentence (P/S). Wolfe-Quintero et al. (1998) reported only a single study using these measures. Including only active passives, Kameen (1979) found that all of the measures significantly distinguished poor and good writers. It was found that good writers employed passive constructions more than poor writers in their papers.

The brief outline above shows that there is a wealth of measures to investigate syntactic complexity. Nevertheless, in her influential study, Ortega (2003) stated a great majority of 25 studies she reviewed employed three or fewer measures. Supporting Ortega's claim, Bulte and Housen (2018) also argued that a great majority of the studies on syntactic complexity in a second language employ only a limited number of these measures such as MLS and subordination ratios which they label as "popular" (p.2). These findings urge the need for a more multidimensional analysis of syntactic complexity. In that sense, Norris and Ortega (2009) urge that measures of coordination, subordination, and global complexity, in addition to phrasal and clausal sophistication, need to be used in studies investigating syntactic complexity.

2.4.2.2. Syntactic Complexity Studies in Second Language Writing

Syntactic complexity measures are used "at least" for the following three purposes in second language writing research: i) to measure proficiency, ii) to define performance, and iii) to evaluate development (Ortega, 2012, p. 128). Over the years, many studies have been conducted to examine the syntactic complexity of second language learners using the measures above. The original aim of investigating syntactic complexity was linked to monitoring the development of second language learners (Liu, 2018). Nevertheless, the idea that a wide range of syntactic features should be involved in texts of more developed writers paved the way for investigating proficiency. In this sense, it has been stated in the literature that syntactic complexity is an indicator of proficiency in a second language in that more advanced second language learners use more syntactically complex structures (Wolfe

Quintero et al., 1998; Ortega, 2005). In line with this argument, Bulte and Housen (2014) regarded syntactic complexity “as a valid and basic descriptor of L2 performance, as an indicator of proficiency and as an index of language development and progress” (p.43). Studies have revealed that the use of coordination is an indicator of beginner second language learners whereas the use of subordination structures is a descriptor of intermediate proficiency level (Norris & Ortega, 2009). In addition, it has been demonstrated that the features of phrasal complexity only appear at more advanced stages (Bulte & Housen, 2012; Ortega, 2003). These findings have led to the conclusion that second language learners usually develop from a dependence on coordination and subordination features at beginning and intermediate stages to a dependence on the features of phrasal complexity at more advanced stages (Ortega, 2015). In line with this developmental progress, the writing development of second language learners is expected to begin with complexity features at the clausal level, and come to an end with features of phrasal complexity (Biber, Gray & Poonpon, 2011). In this sense, it has been suggested that studies on academic writing, especially in the field of ERPP, should include phrasal complexity measures along with the other complexity measures (Lu, 2011; Lu & Ai, 2015). Recent studies have also investigated the relationship of syntactic complexity with several issues such as topic (e.g., Yang et al., 2015), genre (Lu, 2011; Staples & Reppen, 2016), syntactic development (Bulté & Housen, 2014), and writing quality (Taguchi et al., 2013).

The studies in the literature have distinguished differences between non-native speakers and native speakers in terms of syntactic complexity (Foster & Tavakoli, 2009;). For example, non-native speakers were proved to use a higher amount of coordination and complex phrases and the lower amount of subordination in online discussion (Mancilla et al., 2017). Similarly, Ai and Lu (2013) showed that non-native college students formed shorter sentence, T-units, and sentences; used a lesser amount of subordination and fewer number of nominal phrases compared to native speakers. Furthermore, Neff et al. (2007) stated that Spanish learners of English significantly preferred subordination instead of phrasal sophistication compared to native speakers of English. In that sense, Lu (2011) argued that the best measures among the most commonly used syntactic complexity measures were CN/C and MLC.

Another line of research on this issue is to investigate the effect of the first language on syntactic complexity. In a comprehensive study, Lu and Ai (2015) investigated 1400 argumentative essays written by learners with ten different first language backgrounds. They

compared these essays with 200 argumentative essays produced by U.S. college students using 14 syntactic complexity measures. They reported statistically significant differences between the native and nonnative groups as well as among the non-native group itself in all complexity measures. On the basis of these findings, they suggested that “intergroup variation in syntactic complexity cannot be accounted for by proficiency alone but that learners’ L1 may play a role in the syntactic complexity in their L2 writing as well” (p.24).

Investigating the relationship between writing quality and syntactic complexity, Taguchi et al. (2013) examined linguistic features distinguishing argumentative essays of different proficiency produced by nonnative writers of English. They concluded that noun phrase modification was likely to be a factor in the quality of the essay. Recently, the effect of genre on syntactic complexity has also been the focus of researchers who argues that genre is an important variable influencing syntactic complexity. For example, Staples and Reppen (2016) investigated a corpus of 120 papers produced by three groups of first-year students with different native languages in two genres. They reported that the genre was found to have an effect on the lexicogrammatical choices of the participants. In addition, comparing argumentative and narrative texts, Lu (2011) found that learners produced argumentative essays of higher syntactic complexity. Similarly, Way et al. (2000) reported that participants wrote more syntactically complex descriptive tasks compared to expository tasks.

Studies have also proved evidence that a number of variables related to learner, context, and the task may have an effect on second language syntactic complexity. For example, Sotillo (2000) showed that learners produced more syntactically complex texts in a writing task through computer-mediated communication. Furthermore, Ellis and Yuan (2004) reported that participants given adequate planning time performed better in a narrative essay task when they were compared to those given insufficient planning time. In a more recent attempt, Lu (2011) stated that participants wrote texts of higher complexity in untimed conditions compared to timed conditions. In addition, the instructional setting was also found to influence syntactic complexity. For example, Ortega (2003) found that English as a Second Language learners produced more syntactically complex texts than English as a Foreign Language learners. On the other hand, no relationship has been identified between the complexity of the task and syntactic complexity in second language writing (Adams, Newton & Nik, 2015; Kuiken & Vedder, 2012).

A number of studies have been conducted to investigate the syntactic development of learners over time (Casanave, 1994; Stockwell & Harrington, 2003; Vyatkina, 2013). For

example, in their longitudinal study, Norrby and Håkansson (2007) investigated adult L2 Swedish learners over a period of one year using sentence length, subordination, and nominal vs. verbal style. Collecting both written and oral data, they concluded that learners could be classified into the following types: The Recycler, The Thorough, The Careful, and The Risk-taker on the basis of complexity and morpho-syntax analysis. Similarly, focusing on individual developmental paths and differences, Vyatkina (2013) investigated the syntactic complexity development of two beginner L2 German writers over four semesters. Although she identified a general developmental tendency as to frequency and variety of syntactic complexity features, it was found that the two learners differed from each other after the second semester.

Finally, there exists a recent interest in identifying the needs of novice international scholars in writing for publication process in order to design them an effective training (Li & Flowerdew, 2020). In a respond to this call, drawing attention to the role of syntactic complexity in the pedagogy of ERPP, Lu, Casal and Liu (2020) identified the ways expert scholars use a number of complex structures to successfully attain several rhetorical purposes in a social sciences research article introductions. In addition, Ansarifar et al. (2018) examined the English abstracts of L1 Persian scholars' theses and dissertations and compared them with the abstracts of published research articles written by expert scholars in terms of degree of phrasal complexity. They stated that the lowest phrasal complexity degree was found in theses whereas the highest degree was detected in published research articles, which paved the way for them to conclude that writing for publication turns out to be more multifaceted with increasing writer proficiency. In the same vein, Song and Wang (2019) investigated the English abstracts of dissertations produced by L1 English and L1 Chinese doctoral students and showed that there was significantly less subordination in the dissertations of Chinese students. Furthermore, Wu et al. (2020) examined unpublished research articles produced by ELF scholars from different L1 and disciplinary backgrounds and published research articles of native English scholars in terms of syntactic complexity. They reported that ELF scholars employed longer length of production units, more coordinate phrases and nominals but fewer subordinate clauses compared to native English scholars. As a result, they concluded that ELF scholars used these features in order to be more explicit and concise as well as to follow conventions of academic writing. Finally, Yin, Gao and Lu (2020), examined the syntactic complexity in seven sections of published research articles written by expert and emerging international scholars. They compiled a

corpus consisting of 30 published research articles produced by emerging Chinese scholars and 30 published research articles written by expert Chinese scholars. They reported significant differences between the groups with regard to 14 syntactic complexity measures of research article sections.

In sum, the studies briefly described above indicate that syntactic complexity in a second language can be influenced by a number of variables and may differ for different indices. Therefore, studies into syntactic complexity should regard it as a multidimensional construct and issue take all these variables into account.

2.5. Corpus Linguistics

Before applied in the field of linguistics, the corpora were used to indicate the process of collecting written products of a specific nature. In this sense, the origin of corpora dates back to the 13th century when the Christian Bible was analyzed in depth by the biblical scholars who investigated the occurrence of citations, the words, and the alphabetical ordering in the Bible (O’Keeffe & McCarthy, 2010). Other examples of early corpora include Samuel Johnson’s English dictionary in the 17th century and the Oxford English Dictionary project in the 18th century. The corpora at those times were manually compiled and analyzed. However, an enormous leap in Corpus Linguistics was made in the 1950s when the importance of the collection of real data was emphasized by the American structural linguists who put the focus on data collection and using real language data in linguistics studies. Accordingly, the first computer-based corpora and concordances appeared in the 1950s (Parrish, 1962), and the remarkable development in concordances occurred in the 1970s in the catalog indexing field. However, the use of corpora as a tool in linguistics did not come into until the 1980s and 1990s when the hardware and software revolution in computer science took place (McEnery, Xiao & Tono, 2006). The growth in computer and internet technology made the collection and the analysis of data and the dissemination of knowledge much easier, enabling researchers to pay particular attention to things machines are not able to do (Rundell & Stock, 1986). In this regard, technology played a significant role in the development of corpora in that the use of computer corpus extended to all language-related fields (Granger, 2007). The extensive application of the computer corpus led to the establishment of a brand new discipline: corpus linguistics.

Corpus linguistics simply refers to a linguistic methodology which is based on the utilization of digital collections of real language data; in other words, corpora (Granger, 2002). It is a

“new research enterprise, [. . .] a new philosophical approach to the subject, [. . .] an ‘open sesame’ to a new way of thinking about language” (Leech, 1992, p. 106). Contrary to the Chomskyan approach, the focus in corpus linguistics is on performance and description with both quantitative and qualitative analysis. The main power of corpus linguistics is that it is quite appropriate to conduct quantitative analyses providing one of the most reliable sources of frequency, which is a feature of language that gives clues about what is possible and what probably occurs in the language (McEnery & Wilson, 1997). The corpus linguistics also allows researchers to investigate large samples of real language data and thus leads to improved language descriptions and the way it works with a focus on lexical and grammatical patterning. Offering a new empirical aspect, the tools corpus linguistics use in research has transformed the language studies, enabling researchers to make stronger explanations, support their assertions and deal with language more confidently (Hyland, Huat & Handford, 2012). At the center of corpus linguistics stands corpus. Therefore, corpus and related terms will be explained in the following parts.

2.5.1. What is Corpus?

Corpus can be defined as a collection of texts produced for a specific purpose. The term text in corpus linguistics refers to both written and spoken language. The language of the texts in the corpus needs to be contextual and naturally occurring rather than being in the shape of remote words or sentences (Brazil, 1995). The collection of texts in a corpus is made on a purposeful basis; that is, the aim of the corpus is to typify a language or particular part of the language (Biber, Douglas, Conrad & Reppen, 1998). Although corpora have been used in research for a long time, it has recently become influential due to its transformation from manually analyzed texts into the machine-readable format. Nowadays, there are several corpora containing tens, even hundreds of millions of words, such as British National Corpus (<http://www.natcorp.ox.ac.uk/>) and Corpus of Contemporary American English (<https://corpus.byu.edu/coca/>).

The growth of corpora in size requires the use of computers and some specific software in the analysis since it is impractical to try to analyze large contemporary corpora manually. The software produced to analyze corpora enables researchers to investigate a wide variety of things, from creating lists of word and phrase frequencies and identifying collocations to making use of a range of statistical functions to interpret the compiled texts (Cheng, 2011). Most of the corpora include only naturally occurring data for the aim of corpus linguistics is

to describe the use of language and put forward theories of language based on actual language use (Tognini-Bonelli, 2001).

The basic units of corpora are words which are classified into *types* and *tokens*. The term token indicates the total number of words in a text or corpus involving all repetitions, whereas the term type is the number of the unique word in a text or corpus removing repetitions of the very same words (McEnery & Wilson, 2003). For example, the sentence “a good movie is a movie that you like” has nine tokens but seven types since “a” and “movie” are repeated. Nevertheless, there is a debate regarding what comprises a type. For example, make, make, making, and made are all originated in the same root. That is, they share the same lemma. The question raised at this point is that such words should be calculated as different words (4 words in this example) or as one word since they share the same lemma (Kennedy, 2014). The solution lies in the software used in the analysis. Although most of them count such words as different words, some of them empower researchers to list such words as one word. The type and token frequencies are used for a number of purposes in the analysis, such as identifying keywords (Bondi & Scott, 2010) and phraseology (Granger & Meunier, 2008).

The search in corpora is carried out via special software called concordances. A concordance is a basic program transforming texts into electronic databases which can further be searched. These programs are used to mark the intended information by utilizing special taggers and thus allow users to search for a particular word or word combinations in the corpus (Mackey & Gass, 2015). The searched item is presented in the concordance line with the words around it. In most cases, concordances are displayed in Key Word in Context (KWIC) format. Figure 7 illustrates a sample concordance search.

Query rather 5,021,571 > Negative filter (excluding KWIC) would, should 4,723,167 >	
Negative filter (excluding KWIC) J.*, RB.? 3,697,900 (162.70 per million) ⓘ	
First Previous Page <input type="text" value="62"/> of 184,895 <input type="button" value="Go"/> Next Last	
doc#7732	work does not come from just from me , but rather from continued discussions with others
doc#7732	scientists to follow your methods , but rather to provide readers with a sense of how
doc#7732	details about each observation of experiment ; rather , the results must tell a story and inform
doc#7738	agreed that they would not retreat , but rather hold their ground in order to give the
doc#7760	building a brand that constantly evolves rather than being static , just as the library
doc#7767	taxes at current levels were exaggerated . " Rather than returning to 1995 service levels as
doc#7768	revolve around ' taking rape seriously ' but rather subservience to ruling-class pressure and
doc#7777	your newborn baby and pose for a photo , rather get them to lay down next to the baby or
doc#7786	simply taking children from existing centres rather than increasing participation . One respondent
doc#7801	range) . This returns the value as a number rather than as a string , meaning we no longer
doc#7801	or parseFloat , and the + operator adds rather than concatenates : HTML5 also introduces
doc#7802	is not known for any one cuisine - it is rather , just a place , influenced by time and
doc#7818	differences befitting a non-profit foundation rather than a VC firm . Over the last several
doc#7818	post is not about preventive medicine . Rather , it taps into my neuroscience roots ,
doc#7827	provide evidence that large social networks , rather than large brains , contribute to social
doc#7831	Africa Quirks It is more usual to sit outside rather than inside when visiting a bar in Uganda
doc#7832	Social Media 101 , say event organizers , rather it is a celebration of how social media
doc#7847	to work out who 's a waste of my time , rather than some pseudo crackpot networking theory
doc#7852	</script> If the " S " inside the icon is white rather than blue () , 0 script tags have been
doc#7862	n't . " Then , if they ca n't come in , rather than that person just disappearing off
First Previous Page <input type="text" value="62"/> of 184,895 <input type="button" value="Go"/> Next Last	

Figure 7. Sample of Concordance Search

2.5.1.1. Mode of Communication in Corpora

It is possible to compile a corpus in any mode, though corpora are generally compiled in three modes: spoken, written, and video. Corpora consisting of written production generally do not pose much technical difficulty in the compiling process. The software Unicode Standard, a global encoding standard for identifying characters in written texts in different languages (Unicode Consortium, 2015), enables researchers to store, change and show textual stuff in different languages in a consistent way. However, the compilation of written corpora may be a lengthy and error-prone process if the texts will be scanned or typed from written materials (Kennedy, 2014). The second mode is the spoken language. Construction of a spoken corpus is a time-consuming activity since spoken materials are quite challenging to collect and transcribe (Biber, 2006). Although there are transcribed spoken materials such as parliament speeches and news reports on the Internet, these materials are not considered reliable since they were not prepared for linguistic purposes (Mollin, 2007). Therefore, the spoken corpus data is most of the time collected by making the records of the interactions and transcribing them. The transcriptions are then compiled into a machine-readable spoken corpus. Special software allows researchers to link the transcriptions to the original recording so that users have the opportunity to listen to the searched item in the recording. Examples

of such corpora are the British component of the International Corpus (ICE - GB) and The Origins of New Zealand English (ONZE). Finally, video corpora, which are relatively new in corpus linguistics, focus on the analysis of the visual medium. In such corpora, the video analysis is integrated with textual analysis using special software such as EUDICO Linguistic Annotator, which enables users to match annotations with media streams (Hellwig & Van Uytvanck, 2005). Such corpora make the investigation of the relation between speech and gesture (Carter & Adolphs, 2008) and sign language (Schembri, Johnston & Goswell, 2006) possible. However, it is not a requirement for a corpus to be made up of only one mode. There are many corpora involving linguistic data from multiple modes. The attention should be devoted to the distinction between the modes, which may be substantially different (McCarthy & Carter, 1995).

2.5.1.2. Types of Corpora

With the sprung of digital material, many types of corpora have appeared around the world, and this tendency is not likely to cease since corpus linguistics attracts increasingly more researchers (Meunier & Gouverneur, 2009). Corpus linguistics usually describes two types of corpora: general and specialized. Being larger than the latter, the former is compiled to investigate the examples of language use as a whole. An example of such corpora is the Bank of English. Specialized corpora aim at investigating a particular genre, register or variety of the language use. Michigan Corpus of Academic Spoken English (MICASE) is an example of specialized corpora. In other words, the distinction between the two types of corpora lies in the purpose for which they are built.

In an attempt to provide a bird'-eye view of the available corpora, Lee (2010) identifies three categories: i) major English language corpora, ii) developmental, learner and lingua franca corpora, and iii) non-English corpora and multilingual corpora. The first category refers to monolingual English language corpora, which are categorized into general, speech, parsed, historical, specialized, and multimedia corpora (see Table 3). The corpora in the second category consist of lingua franca speech, non-native or learner language, and native-speaker developmental language (see Table 4). The last category covers the corpora that contain monolingual language data other than English and more than one language data (see Table 5). Of these categories, the second category, especially the learner language corpus, is , are significantly important for English Language Teaching field as well as for this study.

Table 3

Developmental, Lerner and Lingua Franca Corpora

Types of developmental learner and lingua franca corpora	Representative Corpora
Developmental Corpora	Language CHILDESH database and Polytechnic of Wales (POW) Corpus, Louvian Corpus of Native English Essays (LOCNESS; 324,000 words), British Academic Written English Corpus (BAWE; 6.5 million words), Michigan Corpus of Upper-level Student Papers (MICUSP; 2 million words)
ESL/EFL learner corpora	International Corpus of Learner English (ICLE; 3.7 million written words), Louvian International Database of Spoken English Interlanguage (LINDESI; 1 million words), Lancaster Corpus of Academic Written English (LANCAWE), Montolair Electronic Language Learners' Database (MELD; 98,000 words), Chinese Academic Written English (CAWE; 406,000 words), International Corpus of Crosslinguistic Interlanguage (ICCI), Japanese EFL Learner Corpus (JEFLL Corpus; 700,000 words), Learner Business Letters Corpus (Learner BLC; 200, words), Learning Prosody in a Foreign Language (LoaP) corpus; (more than twelve hours of recordings)
Lingua franca corpora	Vienna-Oxford International Corpus of English (VOICE Corpus; 1 million words, 120 hours), The Corpus of English as a Lingua Franca in Academic Settings (ELFA Corpus; 1 million words, 131 hours)

Table 4

Major English Language Corpora

Types of language corpora	Representative corpora
General English Corpora (written, spoken and both)	Brown Corpus of written Academic English, FROWN (Freiburg-Brown Corpus of written American English), Lancaster Oslo-Borgen (LOB) corpus of written British English. FLOB (Freiburg-LOB) corpus of written British English, Wellington Corpus of Written New Zealand English, Australian Corpus of English (ACE), Kolhapur Corpus of Indian English, International Corpus of English (ICE), Bank of English, British National Corpus (BNC), American National Corpus (ANC), Corpus of Contemporary American English (COCA)
Speech Corpora	Spoken English Corpora (SEC), Machine Readable Spoken English Corpus (MARISEC), London-Lund Corpus of Spoken English, Intonation Variation in English (IVE) Corpus, Freiburg Corpus of English Dialects FRED (FRED), Cambridge and Nottingham Corpus of Discourse in English (CANCODE), Switchboard Corpus, Santa Barbara Corpus of Spoken American English (SBCSAE)
Parsed written corpora	Lancaster Parsed Corpora (LPC), Surface and Underlying Structural Analyses of Naturalistic English (SUSANNE) Corpus, ICB-GB (Great Britain), Penn-Helsinki Parsed Corpus of Middle English, York-Helsinki Parsed Corpus of Old English Poetry
Historical corpora	Helsinki Corpus of English, A Representative Corpus of Historical English Registers (ARCHER), Corpus of Historical American English (COHA), Lampeter Corpus, Newdigate Newsletter Corpus, Corpus of Early English Correspondence (CEEC), Corpus of Late Eighteenth-Century Prose, Corpus of later Modern English Prose, Zurich English Newspaper Corpus, Old Bailey Corpus, Corpus of English Dialogues (CED)
Specialized corpora	Michigan Corpus of Academic Spoken English (MICASE), British Academic Spoken English corpus (BASE), Limerick-Belfast Corpus of Academic Spoken English (LIBEL CASE), City University Corpus of Academic Spoken English (CUCASE), British National Corpus (academic component), LOB (category J texts: "learned and scientific writings"), Chermnitz Corpus of Specialised and Popular Academic English (SPACE), Reading Academic Text corpus (RAT), Professionals English Research Corpus (PERC), Wolverhampton Business English Corpus, Business Letters Corpus (BLC)
Multimedia corpora	Santa Barbara Corpus of Spoken American English (SBCSAE), Scottish Corpus of Texts and Speech (SCOTTS), English Language Interview Corpus as a Second-Language Application, Singapore Corpus of Research in Education (SCoRE), multimedia corpus of European teenager talk (SACODEYL project)

Table 5

Monolingual and Comparable Corpora

Types of Corpora	Non-English and Multilingual Corpora	Representative Corpora
Monolingual Corpora	Non-English	CORIS/CODIS (Italian), Czech National Corpus (Český Národní Korpus), Hungarian National Corpus, Hellenic National Corpus (also known as the ILSP Corpus), German National Corpus, Modern Chinese Language Corpus (MCLC), Polish National Corpus, Russian Reference Corpus (BOKR), Slovak National Corpus, the Korean National Corpus (or Sejong Balanced Corpus), PAROLE (Preparatory Action for Linguistic Resources Organisation for Language Engineering), EMILLE (Enabling Minority Language Engineering) Corpus, Corpus Gesproken Nederlands, Balanced Corpus of Contemporary Written Japanese (BCCWJ), Academia Sinica Balanced Corpus, the Peking University corpora, Modern Chinese Language Corpus, Lancaster Corpus of Mandarin Chinese, UCLA Chinese Corpus, Lancaster–Los Angeles Spoken Chinese Corpus (LLSCC), ZJU Corpus of Translational Chinese (ZCTC), the New Corpus for Ireland, Cronfa Electroneg o Gymraeg, Scottish Corpus of Texts and Speech, Oslo Corpus of Bosnian Texts, the ‘Brown’ Corpus of Bulgarian
Parallel and comparable multilingual corpora		The English–Norwegian Parallel Corpus (ENPC), English–Swedish Parallel Corpus (ESPC), Oslo Multilingual Corpus (OMC), IJS–ELAN Slovene–English Parallel Corpus, COMPARA (English and Portuguese), FSU Chinese–English Parallel Corpus, the Babel Chinese–English Parallel Corpus, Kacenska (English to Czech), MULTEXT-East (English to nine different languages) and the HKIEd English–Chinese Parallel Corpus

2.5.2. Learner Language

Learner language has always been a fundamental topic of investigation in second and language education research. Researchers have expressed interest in learner language for two basic reasons: it gives researchers a better interpretation of second language acquisition processes, and it provides teachers with rich and beneficial data for promoting teaching and developing learner tools (Granger, 2007). Researchers have used a wide range of data types in examining learner language. Of these data types, Ellis (1994) emphasizes three main categories: language use data, which represent the attempts of learners to use a second language; multilingual judgements, which shows the feelings of learners regarding the second language and self-report data, which enquires into the strategies of learners (see Figure 8).

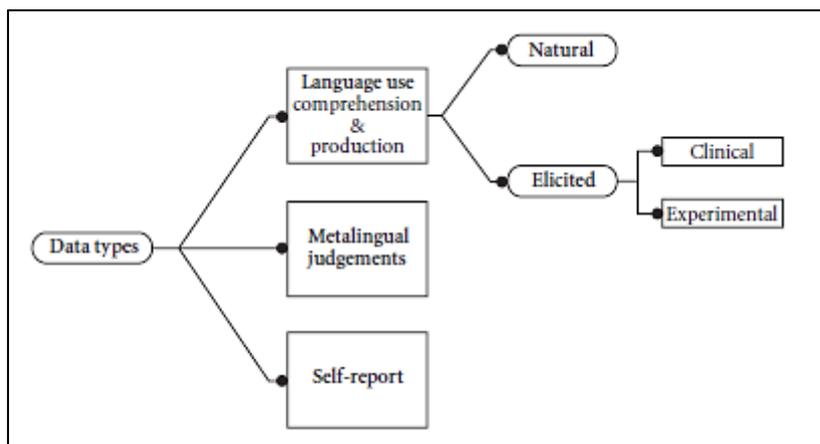


Figure 8. Data types used in SLA research

Most of the studies in second language acquisition research tend to prefer experimental and introspective data and ignore natural language use data since it is quite challenging to control a non-experimental context and natural language use falls short of providing reflecting the whole body of learners' linguistic skills (Larsen-Freeman & Long, 2014). Similarly, Mark (1998) argues that researchers have devoted more attention to some factors affecting language learning and teaching. Mark further states that the three main components were involved in language teaching approaches (see Figure 9).

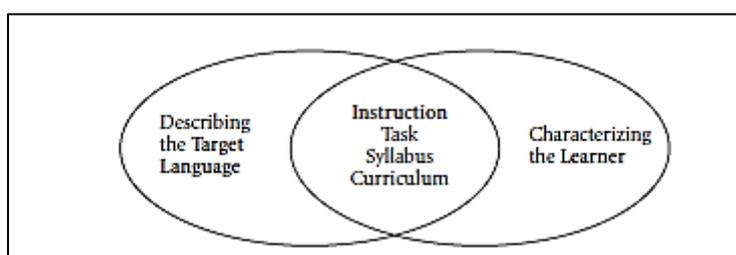


Figure 9. The main concerns of mainstream language teaching

According to Mark, who puts forwards that learner production is absent, the mainstream approaches used in language teaching, learner language, which elucidates the other areas, should be integrated into language teaching and learning research (see Figure 10).

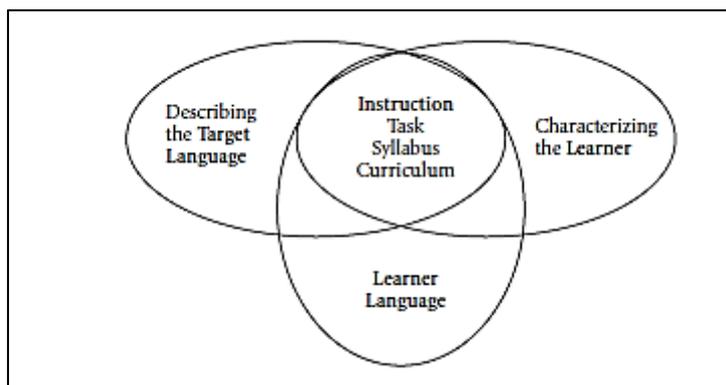


Figure 10. Focus on learner output

It is clear that there is a lack of natural language data used in foreign and second language learning studies. The fact that researchers need more and better quality natural language data to investigate learner language resulted in the introduction of a new type of data source: learner corpora.

2.5.2.1. Learner Corpora

Having emerged in the 1980s as part of corpus linguistics, the learner corpus is a relatively new second and foreign language teaching area. Learner corpora refer to an electronic textual database containing the language of non-native speakers, which is the language of second and foreign language learners (Granger & Leech, 2014). The primary objective of learner corpora is to investigate all aspects of learner language on the basis of corpus linguistics' principles, methods, and tools (Granger, 2002).

The first examples of learner corpora appeared in the heyday of Error Analysis. Corpora at those times were usually small, with a maximum of 2,000 words collected from a limited number of learners (Nesselhauf, 2004a). Besides, learner language samples were diverse in content and character, which made these studies uninterpretable and unrepeatable (Ellis, 1994). The analyses in those early corpora were carried out on the basis of isolated errors and ignore all other aspects of learner language. Thus, the attempts of early corpora failed in providing a full picture of learner performance. Modern learner corpora, on the other hand, are much bigger and more complex than their earlier counterpart in that they present learner language in a computerized form and contextualized way, presenting both the correct and erroneous uses (Dagneaux, Denness & Granger, 1998). Therefore, learner corpora allow researchers to investigate the total interlanguage of learners with software designed for

linguistic analysis. There are a number of key concepts in learner corpora, which will be explained.

2.5.2.1.1. Key Concepts in Learner Corpora

The first notion in learner corpora is authenticity, which is one of the strengths of corpus linguistics. Authenticity in corpus linguistics refers to the data collection process in which the language is collected from real interactions of people as a part of their daily activity (Sinclair, 1996). This notion excludes the data collected in artificial or experimental circumstances. Therefore, it means, for learner corpora, that the data obtained through entirely experimental processes cannot be considered as authentic data. At this point, a complicated problem arises for learner language since it is almost impossible to collect data as genuine as native speaker data in the contexts of foreign and second language learners where a particular amount of artificiality always exists (Granger, 2002). For that reason, authenticity in learner corpora should be regarded somewhat differently from the original term; that is, the data collected through such genuine classroom activities as essay writing and read aloud should be considered authentic (Gilquin, Granger & Paquot, 2007).

The second notion is that learner corpora deal with non-native varieties of English, especially two of them: English as a Second Language (ESL) and English as a Foreign Language (EFL) (Granger, 2002). The term ESL is used to explain the acquisition of English in countries where English is spoken as the native language. EFL refers to the instructed learning of English in settings where English is not the native language. Learner corpora typically contain language samples of these varieties.

The other notion is about textual data in the learner corpora. The learner data are required to be composed of constant sessions of discourse instead of words and sentences presented in isolation. Thus, learner corpora cannot be degraded to so-called corpora of errors, which are the compilation of inaccurate words and sentences obtained from learner texts (James, 1998). In other words, textual data in the learner corpora consist of both correct and incorrect samples of the language.

The importance of design criteria is another notion in learner corpora. Since a great amount of variation exists in learner English, diverse samples of learner language compiled haphazardly cannot be considered as learner corpus data. Therefore, strict design criteria should be adopted in learner corpora. Although the criteria share some aspects with the

native speaker corpora, there are also some unique and specific aspects concerning the learner and the task. Table 6 shows the learner and task-specific design criteria of the International Corpus of Learner English, which contains 3.7 million words (Granger, Dagneaux, Meunier, & Paquot, 2009).

Table 6

International Corpus of Learner English Design Criteria

Learner variables	Task variables
Age	Medium
Learning content	Field
Proficiency level	Genre
Gender	Length
Mother tongue background	Topic
Region	Timing
Knowledge of other foreign languages	Exam
Amount of L2 exposure	Use of reference tools

Another notion is the purpose. Researchers compile learner corpora for specific purposes. For example, they may wish to test a particular theory in the field of second language acquisition, or they may attempt to develop better teaching methods and materials (Granger, 2007).

The last notions are standardization and documentation. Standardization is especially important for the annotated or marked up learner corpora that are corpora improved with interpretative linguistic information such as grammatical tagging or semantic annotation. The process of annotation should be carried out on the basis of specific standards. To help researchers, specific standardized annotation tools and software have been developed. The aim of standardizing the annotation is to make the learner corpora comparable with the native speaker corpora. Standardization also aid researchers in providing transparent analysis and replicable results (Mukherjee, 2004). In this sense, Standardization is also a reliable way to ensure a scientifically strong methodology. Documentation regarding the learner and task variables is required in compiling learner corpora. All details about learners and tasks should be archived and presented to the researchers. Documentation allows researchers to build sub-corpora and thus to make more useful comparisons.

2.5.2.2. Research Orientations

Contrary to general second language acquisition studies, which mainly focus on testing hypotheses, learner corpora adopt a descriptive and explorative research orientation (Gilquin & Granger, 2015). This is compatible with the principles of learner corpora since they are compiled as general sources to deal with a great number of research questions that are established after the data collection. However, there are some corpora addressing more explanatory research questions (Housen, 2002; Tono, 2000). Such corpora aim to revisit some second language theories and hypotheses from the corpus linguistics perspective. This kind of research is valuable in that it connects second language acquisition and learner corpus research (Mackey & Gass, 2011). The other orientation of learner corpus studies is that it is strongly linked to teaching; that is, it emphasizes the role of learner corpora in developing teaching methods and educational materials (Gilquin & Granger, 2015).

From these perspectives, learner corpus studies have tried to handle the features of learner language that was ignored. Providing ways of tracing interlanguage, which was thought to be impossible, learner corpus research revealed that patterns of use, that is, under-use and over-use, form the interlanguage (Mukherjee & Rohrbach, 2006). Therefore, it shifted the research interest from morphology and grammar to lexis and discourse with the help of such corpus software as AntConc (Anthony, 2014) and WordSmith (Scott, 2016), which allowed researchers to generate frequency lists of words, obtain all instances of specific linguistic items and determine their lexico-grammatical patternings (Gilquin & Granger, 2015). With such methodological power, learner corpora research has explored the grey areas of learner language. Studies cover a wide range of area: the use of collocations (Granger & Bestgen, 2014; Nesselhauf, 2005; Paquot, 2018), lexical bundles (Chen & Baker, 2010, 2016; Grabowski, 2015; Shin, Cortes & Yoo, 2018), the use of connectors (Granger & Leech, 2014; Leńko-Szymańska, 2008), discourse markers (Das & Taboada, 2018; Maschler & Schiffrin, 2015; Muller, 2005), stance markers (Fuoli, 2018; Gray & Biber, 2015; Hasselgård, 2009), and involvement features (Adel, 2008).

On the other hand, it seems that grammatical features have not received attention as much as the aforementioned areas in learner corpus research (Biber, Gray & Staples, 2014). This is due to the fact that that tagged or parsed data are required for examining grammatical features, which is quite time-consuming and challenging for researchers who thus tend to use of raw learner data (Gilquin & Granger, 2015). Therefore, the studies of grammatical features in learner language are limited to topics whose investigations are based on raw data.

The grammatical studies involve: causative make (Gilquin, 2012, 2016), what-clefts (Callies, 2009; Collins, 2004), modal auxiliaries (Aijmer, 2002, 2018; Chartrand, 2016), articles (Díez-Bedmar & Papp 2008), and demonstrative pronouns (Gray, 2010; Grosz, 2016; Petch-Tyson, 2000).

2.5.2.3. Methodological Issues

The learner corpus research shares a number of methodological issues with corpus linguistics. However, there are also several methods specific to learner corpora: contrastive interlanguage analysis, the integrated contrastive model, and computer-aided error analysis.

2.5.2.3.1. Contrastive Interlanguage Analysis

The learner corpora compiled on the basis of strict design criteria allow researchers to adopt a contrastive approach (Granger, 2007). Contrary to the contrastive approach of the 1960s in which two languages are compared, this brand new approach investigates the comparison and contrast between the performance between native speakers and non-native speakers in a particular situation (Pery-Woodley, 1990). Identified as a new type of contrastive analysis by Selinker (1989) and Contrastive Interlanguage Analysis by Granger (1996), a new approach has become one of the most frequently used methods in learner corpus studies. Contrastive Interlanguage Analysis focuses on two kinds of comparison: i) between a reference corpus and a learner corpus; in other words, between native speakers and non-native speakers and ii) between different learner corpora or subcorpora, in other words between the non-native speakers (see figure 11).

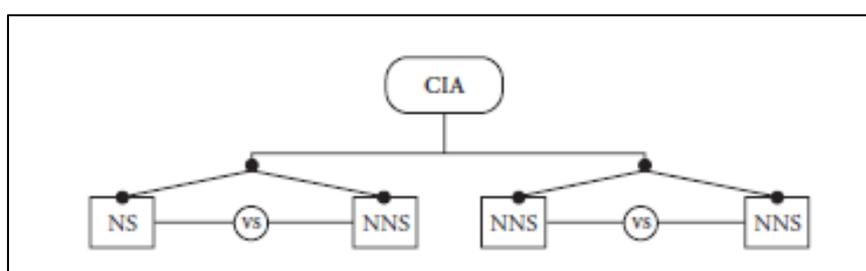


Figure 11. Contrastive interlanguage analysis

The aim of the native speaker and non-native speaker comparison is to discover the distinctive aspects of learner language (Gilquin, 2008). These aspects are revealed by identifying errors and the frequencies of cases where some words or phrases are overused or

underused (Granger, 2015). Before the introduction of contrastive interlanguage analysis, learner language had been analyzed on the basis of errors committed by language learners. Nevertheless, now learner language can be investigated from a much more quantitative perspective with the contribution of frequencies of use, which contributed greatly to the field of language teaching (Hasselgard & Iohansson, 2011).

The significant point in contrastive interlanguage analysis is the use of reference native speaker corpora for the comparison. However, the attention should be focused on the issue of text comparability, which is the use of control corpora having similar features such as genre or topic with the learner corpora (Granger, 2017). The studies where a comparable corpus is not used are likely to fail in providing a correct analysis of learner language leading to a misrepresented view (De Cock, 2002; Granger & Tyson, 1996). It should be noted at this point that the reference corpus does not have to be strictly native or unique as long as the corpus-based reference is stated clearly and in detail (Mukherjee, 2005).

The comparison between the language of non-native speakers aims at improving the existing knowledge about the nature of interlanguage (Granger, 1996). The comparisons of learner language can be made between populations with different native languages. Such a comparison allows researchers to identify the deviating features of the learner language and thus differentiate the sources of these features (Granger, 2017). If the features are found to be restricted to just one native language background, the possible source tends to be interlingual features, in other words, the interference of the native language. However, when the features are found to be shared by a number of populations with different native language backgrounds, then the possible sources are likely to be due to developmental features. Similar to the first type of comparison, a reference corpus for comparison is required for the analysis. Learner corpora involving data from different native language backgrounds are very useful for the comparison of interlanguages since such corpora consist of sub-corpora that share the same design criteria.

In 2015, Granger revisited contrastive interlanguage analysis and designed a new version, which she called CIA² (Figure 12).

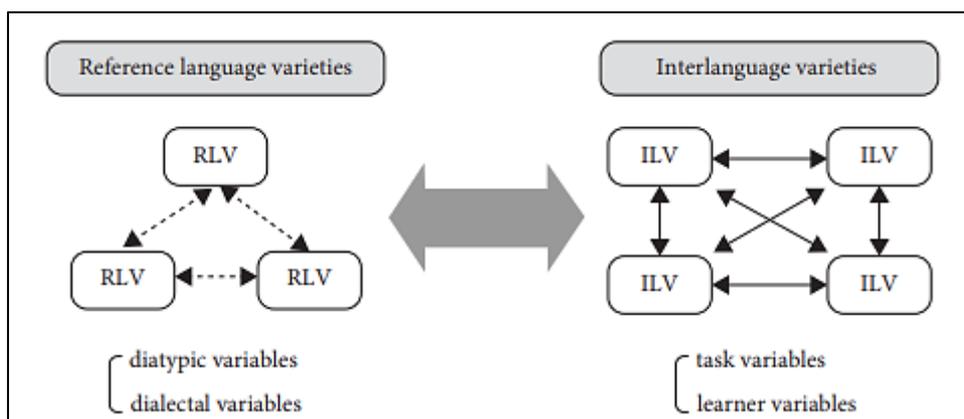


Figure 12. CIA²

In the new model, she promoted the concept of varieties, introducing reference language varieties and interlanguage varieties. The former concept was to suggest that an excessive number of distinct points which can be considered as the reference in analyzing learner language exists (Granger, 2015). Therefore, in addition to the inner circle varieties that are usually used as references, the new version argued that outer circle varieties and corpora consisting of competent second language learners such as English as a Lingua Franca corpus could also be used as a reference. As shown in the diagram, reference language varieties may be used simultaneously and compared. For example, Lee and Chen (2009) used a corpus involving journal article data from experts who were both native and non-native speakers as a reference corpus for the corpus consisting of data from novice writers. Granger also introduced the concept of Interlanguage varieties to describe the learner language. The aim of this new concept was to emphasize the widely varying nature of interlanguage, which should be taken into consideration in learner corpora research. There are a wide variety of variables the effect of that can be investigated. However, there is a lack of research on the effects of most of these variables on learner performance. So, future research on Contrastive Interlanguage Analysis should deal with these variables.

2.5.2.3.2. *The Integrated Contrastive Model*

Being an extremely thorough method to identify the issue of transfer, Integrated Contrastive Model (Gilquin, 2000; Granger, 1996) is the combination of Contrastive Interlanguage Analysis and Contrastive Analysis (Figure 13). Research on the issue of transfer has revealed that insufficient and unreliable data from the contrastive analysis may result in erroneous interpretation of the results and thus lead to faulty attributions of particular phenomena to

transfer (Schachter, 1983). The way to eliminate such a risk is to carry out a cautious contrastive analysis on a bilingual corpus. In this way, researchers are more likely to obtain a comprehensive picture of the features of learners' native language and its differences with the target language (Vanderbauwhede, 2012). The research based on this method allows researchers to prognosticate the possible incidents of transfer and to explain some features of the interlanguage more reasonably (Gilquin & Granger, 2015). For the use of this methodology, see Rasier and Hiligsmann (2007), Demol and Hadermann (2008), and Gilquin, Papp and Dez-Bedmar (2008).

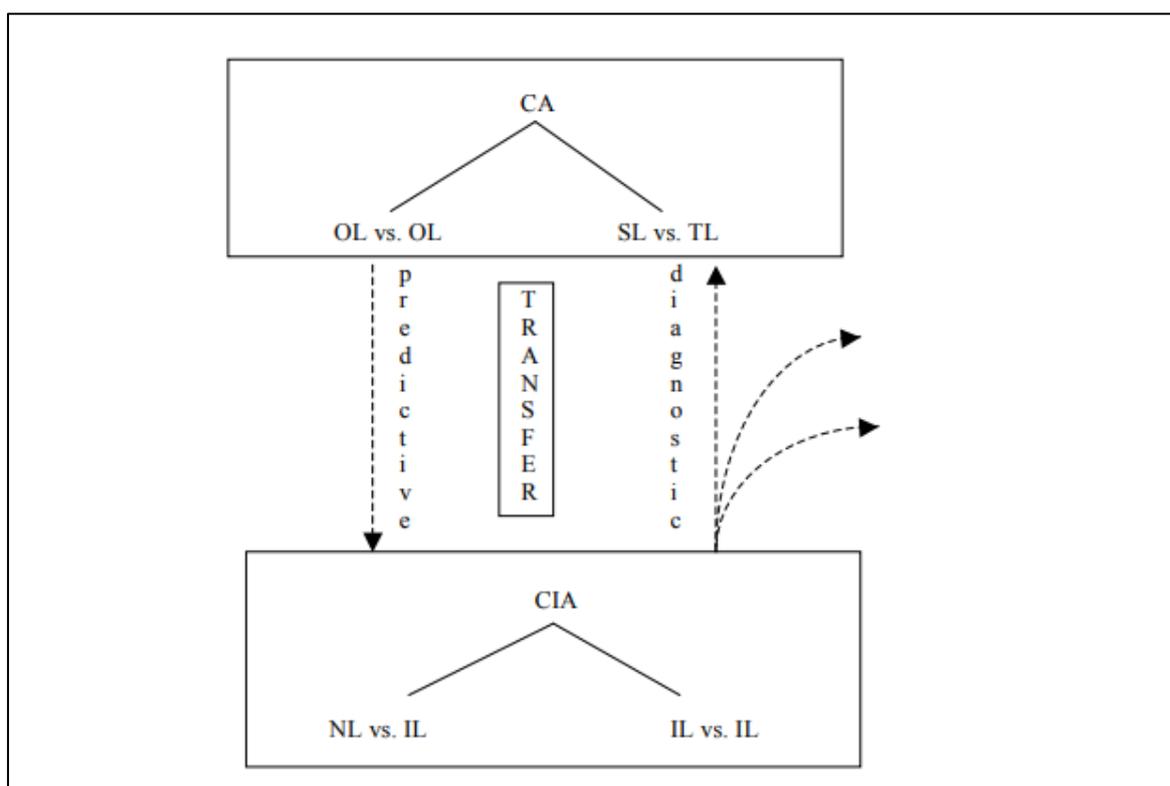


Figure 13. Integrated contrastive model

2.5.2.3.3. Computer-aided Error Analysis

The roots of computer-aided error analysis go back to the 1970s when Error Analysis was popular in the field of language teaching and learning. Although Error Analysis was severely criticized, errors are still regarded as an essential feature of interlanguage, making the study of errors as a worthy endeavor as any other feature of interlanguage (Granger, 2003). Besides, as Ringbom (1987) puts it, error analysis helps researchers to understand the process of second language learning since errors “provide windows onto a system” (Gass, 2013). Therefore, a detailed investigation of learner language, particularly errors, facilitates using the second language more accurately, which is one of the main aims of second

language teaching (Granger, 1999). These requirements and the importance of investigating learner language resulted in the introduction of computer-aided error analysis. With the advancements in technology, especially in computer technology, Error Analysis had its place in the field of second language acquisition one more time. It was thought that although it failed to provide a full picture of the second language data, Error Analysis might still be useful when it was combined with the new technology. In their seminal paper “Computer-aided Error Analysis,” Dagneaux, Dennes, and Granger (1998) argued that Error Analysis “...should be reinvented in the form of computer-aided error analysis...” (p.163). They suggested that the early Error Analysis research had a number of weaknesses and listed the followings as the most serious ones:

- Limitation 1: Error Analysis is based on heterogeneous learner data;
- Limitation 2: Error Analysis categories are fuzzy;
- Limitation 3: Error Analysis cannot cater for phenomena such as avoidance;
- Limitation 4: Error Analysis is restricted to what the learner cannot do;
- Limitation 5: Error Analysis gives a static picture of L2 learning. (p. 164)

According to them, Limitation 1 and Limitation 2 were methodological, whereas others were related to the scope of the Error Analysis. Regarding the data problem, they suggested that clearly described learner language data should be collected in order to provide convincing explanations of the learner errors. With respect to the categories, they classified the error categories used by the traditional Error Analysis research as ill-defined, if any definition was provided. They argued the imprecise nature of categories prevented the researchers from successfully interpret the data. For Limitation 3 and Limitation 4, they argued that the main focus of Error Analysis was on overt errors, and thus correct forms and underused or non-used forms and words were ignored. However, as Harley (1980) puts, learners’ correct use should be considered as important as their errors. Finally, they criticize the Error Analysis for being excessively static and product-oriented though language learning is a dynamic process. Referring to these limitations, they indicated a brand-new direction: the use of computers and corpus linguistics in Error Analysis, which they called computer-aided error analysis. They argued that computer-aided error analysis allows researchers to analyze the correct uses of the learners as well as their errors, making it possible to compare the production of the second language learners and native speakers.

Computer-aided error analysis was introduced by Dagneaux, Denness, and Granger (1998) at Université Catholique de Louvain. The methodology they developed consisted of several steps. First, a native speaker corrected the learner data manually. Then, an error tagging system was developed. Each error was earmarked with an appropriate error-tag, and

corrections were inserted in the file. After the error tagging procedure was completed, with the use of text retrieval software tools, the lists of error types and statistics were obtained. Finally, notable error types were examined on the basis of concordance based linguistic analysis. To increase the efficiency of analysis, they suggested that the analyst should be a non-native speaker with a high level of proficiency in the second language who shares the same native language of learners whose data will be analyzed.

Two methods are generally used in computer-aided error analysis. The first one is choosing a particular linguistic item and searching all erroneous use of the item in the corpus data. Although this method is fairly fast, the investigation is limited to the linguistic items identified as the subject of the study. The second method involves the use of an error tagging system and marking up all of the errors in the learner data. Although the second method is more powerful than the first one in that it provides a fully tagged corpus, it is also more time-consuming. However, with the advances in technology, attempts to develop automatic annotation systems have been made (Izumi et al., 2003; Mason & Uzar, 2000; Tono, 2000). Nevertheless, the common practice in computer-aided error analysis is to use a computer-assisted editor which assists researchers to insert tags in the corpus. When the tagging process is completed, the learner language can be analyzed both qualitatively and quantitatively with the help of special software tools.

Computer-aided error Analysis research has generally used the first method and focused on particular problem areas such as verb tenses (Granger, 1999), lexical problems (Lenko-Szymanska), and article system (Diez-Bedmar & Papp, 2008). Nevertheless, there are also comprehensive error tagging systems whose aim is to examine a large extent of error types. Since such systems have been developed as a part of a project aiming to build a learner's corpus, they have been named after the specific learner corpus. Table 7 shows the learner corpora relating to comprehensive error tagging systems.

Table 7

Learner Corpora Associated with Error Tagging Systems

Learner Corpus	Approximate Size (in words)	L2 Level	L1	Purpose	Mode	Learner Language
CBACLE	1,000,000	Various	Chinese	Academic	Written	English
CLC	20,000,000	Various	Various	Commercial	Written	English
C-LEG	28,000	Advanced	English	Academic	Written	German
FALKO	36,000	Advanced	Unspecified	Academic	Written	German
FRIDA	200,000	Advanced	Various	Academic	Written	French
HKUST	25,000,000	Upper secondary education	Chinese	Academic	Written	English
ICLE	2,000,000	Advanced	Various	Academic	Written	English
JEFLL	700,000	Various	Japanese	Academic	Spoken and Written	English
LLC	10,000,000	Various	Various	Commercial	Written	English
MELD	100,000	Advanced	Unspecified	Academic	Written	English
NICT JLE	2,000,000	Various	Japanese	Academic	Spoken	English
PELCRA	500,000	Various	Polish	Academic	Written	English

Detailed information regarding error tagging systems used in the learner corpora above is not usually available, or such systems have not been made public. On the other hand, there are some well-known and widely-used error tagging systems, which will be briefly explained. These tagging systems are:

- i. the Cambridge Learner Corpus project (CLC),
- ii. the FreeText project,
- iii. the Université Catholique de Louvain (Louvain), and
- iv. the National Institute of Information and Communications Technology Japanese Learner of English (NICT JLE).

The most comprehensive error tagging system is provided by the Université Catholique de Louvain to annotate The International Corpus of Learner English (ICLE). ICLEv2 User Manual (Granger et al., 2009) distinguishes 54 error types in 7 main error domains, which are then divided into sub-categories. This system has been used for scientific purposes and is also commercially available. The other tagging system, the Cambridge Learner Corpus tag set, is used to tag 16 million-word corpus of learner English, which serves to develop teaching materials. CLC error-tag set and coding system can be found in Nicholls (2003). Similar to ICLE, the FreeText error tagging system was developed by the Université

Catholique de Louvain. It was designed to diagnose the errors of learners and then formulate apposite exercises as a part of a CALL program. The error tag set used in the FreeText can be found in Granger (2003). Finally, the tag set used in the National Institute of Information and Communications Technology Japanese Learner of English, which aims to build an interlanguage model of Japanese L2 English learners, was an attempt to develop an automatic error detection system. The error-tag set can be found in Izumi et al. (2004).

2.6. Corpora and Written Academic English

Corpus linguistics has achieved a significant position in the field of language teaching and learning over the past twenty years. It has transformed the way language is understood, studied, taught, and learned. Similarly, the wave of corpus linguistics and corpora has spread to written academic English. Providing rich data involving information on the characteristics of academic texts, it has brought a remarkable aspect to academic writing. Now, with the help of corpora, researchers are enabled to support their claims, make stronger interpretations, and express their ideas with greater confidence (Ädel, Dans, Campoy-Cubillo & Gea-Valor, 2010). The most significant impact of corpus linguistics in academic writing is that it has methodological power to “explain the mechanisms by which knowledge is socially constructed through language.” (Hyland, 2015a, p. 292). A careful investigation of the literature reveals that mainly native corpora are used in the field of academic writing research to carry out genre-based analysis (Yoon, 2008). However, more studies should aim at investigating learner language since second language learners suffer from a number of distinguishing problems that are likely to be discovered through a thorough corpus-based analysis (Gilquin, Granger & Paquot, 2007).

Adapting a co-textual approach, corpus-based analysis of academic texts made notable contributions to the field. First, corpus-based studies have enabled researchers to describe the distinctive features of academic text in detail. Second, it has been revealed that a highly specified phraseology is used in academic texts, and this phraseology varies across genres and disciplines (Gilquin, Granger & Paquot, 2007). Finally, corpora studies have shown that regardless of the context, which may be a scholarship application, a research article, or a piece of personal information, the aim of academic texts is to convince the readers with appropriate claims, arguments, and attitudes (Hyland, 2015a).

Corpus linguistics has been used to discover the structural regularities of a wide variety of genres: describing moves in dissertation acknowledgments (Hyland, 2004b), in Ph.D. theses

introductions (Bunton, 2014), in application statements (Ding, 2007), and in grant proposals (Connor & Upton, 2004), research article introduction section (Cortes, 2013; Ozturk, 2007), the results section (Basturkmen, 2009, 2012; Bruce, 2009) and the discussion section (Peacock, 2002). In addition to corpus-based move analysis studies, researchers have also investigated a number of features that have not been dealt with, such as this and these as pronouns (Gray & Cortes, 2011), code glosses (Hyland, 2007b; Safari, 2018), evaluative that (Hyland & Tse, 2005), attended and unattended this (Jiang & Wang, 2018; Wulff, Römer & Swales, 2012). Additionally, Biber (2006) investigated the differences between spoken and written texts and compared different types of corpora. Furthermore, Samraj (2005) examined the differences between abstracts and introductions in a research article and revealed the differences among different written genres. Likewise, Tse and Hyland (2008) studied the way male and female writers produce academic texts to make a comparison between genders. A number of studies also tried to make comparisons between expert and novice writers (Hartig & Lu, 2014; Hyland, 2006; Mansourizadeh & Ahmad, 2011). Finally, other studies investigated features such as hedges (Hu & Cao, 2011; Yang, 2013), self-mention (Walková, 2019), transition signals (Hyland, 2008a), reader engagement (Hyland & Jiang, 2016; Jiang & Ma, 2018) and bundles (Pan, Reppen & Biber, 2016; Shin, Cortes & Yoo, 2018).

One of the most extensively studied topic in corpora research is the effect of culture and first language on the written academic English. These studies (for example, Loi, 2010; Moreno & Suarez, 2008; Morton, Storch & Thompson, 2015) have bolstered the claim that the schemes used in the first language and second language writing differ, and this influence the writing of multilingual learners in English. Much research on this issue has centered on student genres and has discovered a wide variety of different aspects in first and second language writing in English (for example, Hinkel, 2002; McIntosh, Connor & Gokpinar-Shelton, 2017; Uysal, 2008). The focus of research has recently shifted to the investigation of self-representation (Dueñas, 2007; Matsuda, 2015; Sheldon, 2009) and interpersonal appropriacy (Flowerdew & Wang, 2015; Hyland, 2014) among multilingual authors.

Another contribution of corpus linguistics to written academic English is regarding the construction of academic genres (Hyland, 2015a). All individuals are members of particular social groups according to whose rules they use language. As members of academic society, scholars produce their articles and theses according to their disciplines (Hyland, 2004b). Corpus linguistics is used for identifying the features of these disciplines, enabling researchers to investigate how the disciplines are established and maintained and how the

language is used among their members (Hyland, 2015a). The seminal work on this issue belongs to Nesi and Gardner (2012), who investigated papers of students from different disciplines in different years of education. Based on British Academic Written English Corpus (BAWE), they proposed a genre classification to examine thirteen main assignment types and found that a significant variation exists among the disciplines. Another example is lexicography studies (Granger, 2017; Hyland & Tse, 2007; Vo, 2019; Yang, 2015), which investigate the variation in scholars' choices of academic vocabulary in different disciplines.

Using corpora in written academic English research has also allowed researchers to support the view that academic writing is spread throughout social interaction and intersubjectivity. It has been noticed that scholars use language to admit, build, and negotiate social relationships while producing texts. Corpora studies have been beneficial in identifying aspects that scholars utilize to build authorial self. Such studies include voice in articles (Bondi & Silver, 2004; Hyland, 2001), in student essays (Ivanič & Camps, 2001; Matsuda & Jeffery, 2012), in textbooks (Bondi, 2012; Herbel-Eisenmann, 2007) and in book reviews (Hyland & Diani, 2009) metadiscourse (Dahl, 2004a; Hyland & Tse, 2004; Kawase, 2015), stance (Aull & Lancaster, 2014; Chang & Schleppegrell, 2011; Hyland & Guinda, 2012; Jiang, 2015; Lancaster, 2016) and appraisal (Borglin & Fagerström, 2012; Hood, 2010; Hyland & Tse, 2004).

Although there is an abundance of studies investigating different features of academic writing, there are still a number of research gaps to be addressed. First, the majority of the existing corpora consist of public or easily accessible texts, which are usually published texts or student essays. Therefore, research into academic writing only deals with a limited area of the academic world. To bridge this gap, corpora-based research on texts that are unpublished, less public, or difficult to obtain should be conducted (Hyland, 2015a). Second, more disciplinary genres should be investigated and described. Third, the use of specific genres within particular contexts needs to be examined. Fourth, corpora studies should include multimodal genres, such as academic websites. Fifth, more studies focusing on non-native writers to compare the similarities and differences in their academic writing with both each other and native speakers are needed. Last, corpus studies should include more disciplines to have an understanding of the nature of and expertise in different disciplines (Hyland, 2015a).

Taking up some important points mentioned above, this study aims to collect a corpus consisting of unpublished research articles written by international scholars. Although there

are several corpora focusing on academic writing such as British Academic Written English (BAWE), The Corpus of Academic Learner English (CALE), and The Chinese Academic Written English corpus (CAWE), most of them depend on published texts or texts produced by university students as a requirement of the courses they attended. To the best knowledge of the researchers, this corpus is one of the first corpus, including unpublished and thus unedited research articles. Besides, research articles from five different disciplines, which are communication, economics, and administrative sciences, education, engineering, and medicine, are included in the corpus to gain insights about the disciplines and to compare them. The details of the corpus will be described in the methodology chapter.

2.7. The Sociolinguistic Background in Turkey

Having a population of almost 80 million and an area of 780,580 square kilometers, an overwhelming majority of Turkey lies in Asia. Whereas only a very small part of Turkey is in Europe, Turkey has always been considered as a country that turned its face more towards Europe than Asia, acting as a bridge between the two continents (Kırkgöz 2005; Uysal, 2012). Because of its geopolitical location and history, Turkey has always been caught in the dispute between the West and East, which is manifested in its state policy (Durgun, 2010). Since the foundation of the Turkish Republic in 1923, the country has been undergoing a Western-oriented modernization, which is described as a transition from being a peripheral country to a semi-periphery country adopting center-based policies (Uysal, 2012).

Turkey belongs to Kachru's expanding circle (1986), in which English does not have any colonial past and thus is used for global communication. Although the early contact of Turks with the English language dates back to the 16th century, English became apparent in state schools only in 1908 (Doğançay-Aktuna, 1998). It was not until the 1950s that English flourished in Turkey and took the place of French as the main foreign language of the state due mainly to the impact of increasing American relations. Since then, English has been taught as the main foreign language in state schools, and the schools that are offering English-medium instruction have been increasing (Selvi, 2011). Today almost 20% of programs in Turkish tertiary education use English as the medium of instruction (Arik & Arik, 2014). In this context, Doğançay-Aktuna (1998, p. 37) gives an outline of the role of English in Turkey as:

In Turkey, English carries the instrumental function of being the most studied foreign language and the most popular medium of education after Turkish. On an interpersonal level, it is used as a link language for international business and tourism while also providing a code that symbolizes modernization and elitism to the educated middle classes and those in the upper strata of the socioeconomic ladder.

English also plays an important role in the Turkish higher education system, in which 129 state universities and 73 private universities exist, in terms of both students and scholars. For students, although most of the master's programs do not require having a certain level of proficiency in a second language, they give additional points to those having a score. In addition, having a minimum score in a second language, which is mainly English, is a requirement for doctoral applications. Some universities even oblige Ph.D. candidates to publish one article at the minimum to apply for a Ph.D. program. Graduate students are also required to publish before they present their thesis in order to graduate. Although the type of publication varies according to the degree of the program, typically, students have to publish a conference proceeding for a master's degree and a research article in international indexed journals for the Ph.D. degree as a graduation requirement. As of 2019, 394,174 masters' and 96,199 Ph.D. students were enrolled in several graduate programs in the Turkish higher education system (YÖK, 2020). This means that almost 500,000 students have to publish in a relatively short period.

The role of English is more complicated for scholars since English is quite dominant in hiring, promotion, and reward processes. First, researchers need to have a minimum score of 60 in the national language proficiency exam (Foreign Language Exam) for applying for a position in the universities as well as publishing at least two articles in international journals and receiving a certain number of citations. Moreover, assistant professors have to publish internationally in order to meet the contract renewal conditions, which differ from one university to another. Second, scholars are required to publish in high-ranked journals to apply for the associate professorship examination. The candidates get 8 points for a national indexed article, 15 points for an international indexed article, and 20 points for an article indexed in SSCI, SCI, SCI-Expanded, or AHCI. Likewise, they get 15 points for books published by national publishers while they receive 20 for books published by international publishers. Third, The recent changes in the Regulation on Academic Incentive Payments emphasized the role of international activities and high-indexed journals. For example, participation in national conferences was removed, and the definition of international activities was provided. Accordingly, a conference is defined as international only when more than half of the proceedings are presented by international participants.

Similarly, the changes provided a more detailed definition for international publishers. In that sense, a publisher needs to publish at least 20 books of different scholars in the same discipline written in languages other than Turkish to be considered international. The situation is similar in terms of article publication. For articles indexed in international journals, scholars get twice to four times as many points as national activities depending on the impact factor of the journal. It is important to note that all of these international activities are carried out primarily in English.

A brief overview above clearly shows that the policies and regulations force scholars to publish in international journals. A recent report announced by University Ranking by Academic Performance (2019), a non-profit organization founded by Middle East Technical University to develop a ranking system based on academic performance meters, indicates that 9 Turkish universities ranked among the top 1000 universities and 82 Turkish universities ranked among top 2500 universities. According to the report, no Turkish university ranked in the top 100 or 500 university list. However, in 2011, 9 Turkish universities were able to be in the top 500 list, and 20 Turkish universities ranked in the top 1000 list. The reason for this decrease is the fact that Turkish scholars have recently tended to publish in low impact journals. Only 21,08% of the research articles between 2014-2018 were published in high impact journals. In addition, obtaining an impact factor of 0,77, Turkey fell beyond the world average, which was 1,00, in terms of impact factor. Therefore, it can be argued that a great majority of the publications by Turkish scholars were in low impact journals.

Although there may be several reasons for publishing in low impact journals, such as fulfilling the hiring, promotion, or graduation requirements (Demir, 2018), it is clear that such a situation is closely linked to proficiency in English. Scholars are expected to be highly proficient users of English to publish in high ranked journals. In this context, ERPP emerges as an important field of investigation in Turkey. In order to increase the academic production, writing for publication skills of scholars needs to be carefully examined. Unfortunately, there is a serious lack of studies investigating the academic writing practices of scholars in Turkey. Most of the existing studies are either limited to only one discipline, which is English Language Teaching, or deal with published texts. However, published texts do not provide sufficient information on writing practices since these texts are severely brokered during the publication process. Such a brokering activity may hinder the identification of problems or difficulties scholars experience to a great extent. In addition, the results of studies addressing

ELT scholars cannot be generalized to all Turkish scholars since ELT professionals have a high degree of proficiency in English. Therefore, studies examining the writing practices of scholars should include texts that have not been brokered and several disciplines. Attempting to diagnose the problems Turkish scholars experience in writing for publication process, this study investigates the linguistic problems faced by Turkish scholars using corpus linguistics. In that sense, 216 unedited research articles from five disciplines, communication, economics, and administrative sciences, education, engineering, and medicine, were collected, and the errors the scholars made, as well as the syntactic and lexical complexity of the texts, were analyzed. The details of the methodology of this study are presented in the following chapter.



CHAPTER III

METHODOLOGY

This chapter focuses on the methodology employed in the study. It describes the steps and procedures that were followed to carry out the present study. First, the methodological background of the study is presented. Second, the procedures for collecting the corpus, the design of the corpus, and the tagging procedure are clarified. Third, the collection procedure of the reference corpus is explained. Finally, quantitative analyses conducted in the current study are introduced.

3.1. Methodological Background

This study is exploratory study, aiming to investigate the linguistic problems in writing for publication. This study tries to identify Turkish scholars' writing practices examining their syntactic complexity and grammatical errors in research articles in English. To do this, corpus linguistics was adopted as the main research methodology. Corpus linguistics deals with large realistic and authentic language data kept in digital, usually computerized, databases to examine the language used. Corpus linguistics is considered both as a research tool and as a field of research on its own. However, "the answer to the question of whether corpus linguistics is a theory or a tool is that it can be both. It depends on how corpus linguistics is applied" (Kübler & Zinsmeister, 2015, p. 14). This study adopts corpus linguistics as a research methodology trend and uses corpus linguistics as the main methodology.

Two main research methods exist in Corpus linguistics: corpus-based and corpus-driven. Corpus-based research refers to the use of corpus data in order to discover hypotheses or theories to validate, refute, or refining the linguistic aspects and features originated in these hypotheses or theories (Tognini-Bonelli, 2001). This line of research supports corpus

linguistics as a method approach. corpus-driven research, on the other hand, refuses the construct of corpus linguistics as a method and, instead, put forward that the corpus itself must be the only basis of the hypotheses about language, making an attempt to find out new linguistic concepts by means of inductive corpora analysis (Biber, 2012). In that sense, the present study employs a corpus-based method and tries to describe linguistic features and constructs.

In general, a corpus is a digital collection of real-life language use, which may be both written and spoken, and kept in a computer-readable form (Jurafsky & Martin, 2009; Nugues, 2006; Wynne, 2005). However, there are also more detailed definitions of the corpus. According to Nesselhauf (2004b), a corpus should be designed to be used for general purposes and thus should not be built for one or a few uses. In addition, raising questions regarding the design criteria, the primary role of the corpus, and concerns of representativeness, Sinclair (2005) explains corpus as “a collection of pieces of language text in electronic form, selected according to external criteria to represent, as far as possible, a language or language variety as a source of data for linguistic research” (p.16). In sum, it can be said that a corpus should have the following features: i) it should be machine-readable, ii) it should be based on authentic text, iii) it should be representative of a specific language, or its variety (McEnery et al., 2006, p. 5).

The corpus compiled as a part of this study is a professional corpus since the data were collected from scholars aiming to publish in international journals. However, the corpus also shares two significant features with learner corpora. First, the data were collected from Turkish scholars who were non-native speakers of English, which is one of the feature of learner corpora (Granger & Leech, 2014). Second, scholars visited an academic writing center at a state university during their writing for publication process. The aim of the writing center was to offer instructional support and consulting assistance for academic members during their writing for publication process by enhancing their understanding of their writing progression and by presenting them the basic information and strategies in writing for publication in English. Therefore, it can be argued that the scholars in this corpus were also learners of academic English since the main aim of the academic writing center was not only to enhance the quality of the manuscript but to improve scholars’ awareness about English writing conventions, writing skills, mastery, and strategies by means of their own writing.

Examination of errors allows researchers and practitioners to understand the nature of the errors. These errors could occur in words, phrases, or structures, accompanied by the ways

in which they could be used accurately and inaccurately (Granger, 2003; Nesselhauf, 2004b). Therefore, error tagged corpora have been proved to be highly beneficial for making up or improving learner dictionaries by offering more details of pertinent entries to users of dictionaries as a result of identification of the most frequent errors. In addition, error tagged corpora serve as a beneficial means of measurement in that the improvement in numerous features of the target language can be identified easily with an error-tagged corpus (Buttery & Caines, 2012). Examination of errors can also be used as a foundation for developing instructional materials. An error-tagged corpus is a powerful tool for developing teaching materials that are suitable for learner needs and proficiency levels as well as their strengths and weaknesses in the target language.

3.2. The Collection of the Corpus

3.2.1. Data Collection

For this study, a new corpus, The Multidisciplinary Corpus of Writing for Publication was collected. The Corpus consists of research articles written in English by Turkish academics. The data was collected between 2015 and 2017. In order to collect the data, the academics that had visited a writing center in their publication process were contacted and informed about the study, and then they were asked to send the final drafts of their research articles written in English. The attention was paid to include the research papers that had not undergone any professional proofreading services or checked by a native speaker of English. In addition, the research articles that were translated or submitted to any journal and received any kind of feedback were excluded. Finally, 402 research articles from 18 disciplines were obtained. For this study, 5 disciplines, which are communication, economics, administrative sciences, education, engineering, and medicine, were selected. All of the received articles in the related disciplines were included in the corpus. As a result, The Multidisciplinary Corpus of Writing for Publication consisted of 216 research articles. Then, data cleaning procedure was carried out with the research articles in the corpus in line with the process used in The International Corpus of Learner English (Granger, 1993): personal data, footnotes, endnotes, statistics, tables, graphics, formulas, maps, and figures were removed, and quotations were replaced with <Q>. Finally, each research article was anonymized. The acronyms of the first language of the scholars, the genre, and disciplines were used in the anonymization process, respectively. Therefore, as the first language of all scholars was Turkish, the acronym TR

was inserted first. Next, the acronym of the discipline, for example, EDU for education, was placed in. Then, RA standing for the research article was used to indicate the genre. Finally, the number of the text was written. As a result, each of the articles had codes similar to the following: TR_ EDU _ RA _001. The data were stored in electronic format. The data collection procedure can be shown in Figure 14.

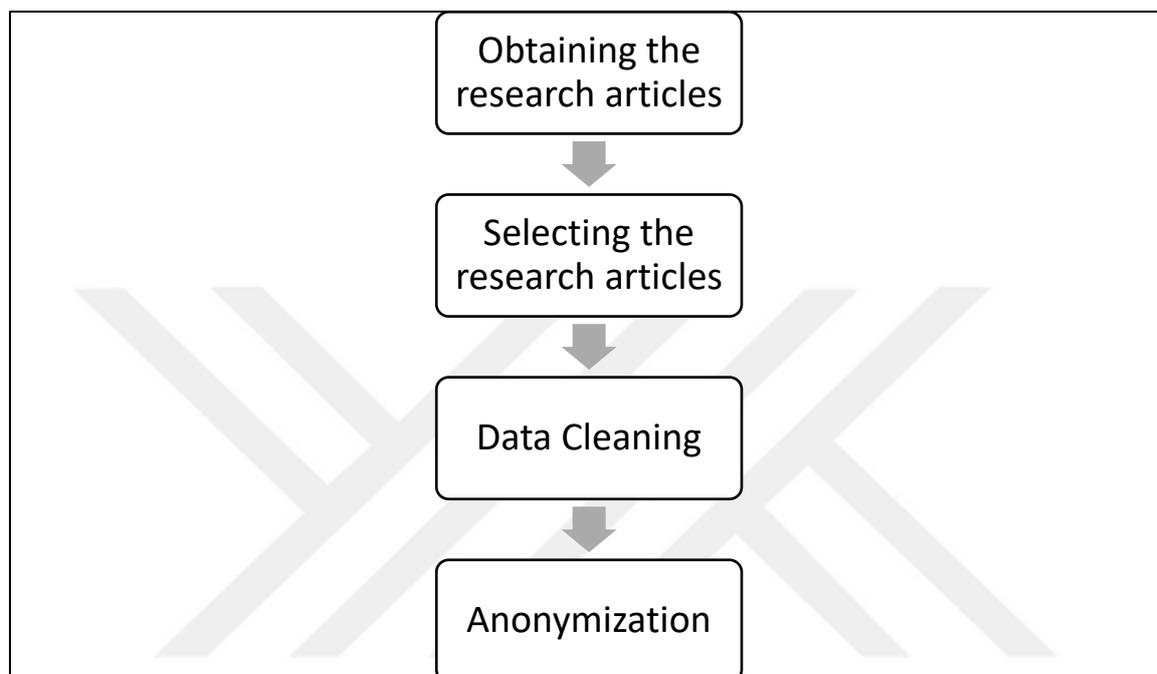


Figure 14. Data Collection Procedure

In sum, The Multidisciplinary Corpus of Writing for Publication includes five sub-corpora and consists of 216 research articles, which contain 886,482 words. Table 8 shows the disciplines, the number of research articles, and words in detail.

Table 8

The Distribution of Research Articles in The Multidisciplinary Corpus of Writing for Publication Corpus

Discipline	Number of Research Articles	Number of Words
Communication	7	25,436
Economics and Administrative Sciences	15	81,033
Education	51	210,984
Engineering	89	378,683
Medicine	54	190,346
Total	216	886,482

Learner language may vary considerably and is affected by a number of linguistic, psycholinguistic, and situational factors. Thus these variables should be controlled in order to obtain reliable and valid results (Granger, 2004). Therefore, one of the most important issues in corpus research is design criteria. In order to build a comprehensive corpus, one needs to determine its design criteria carefully and strictly. The design criteria typically cover 11 features: purpose, size, target language, availability, the learners' nativeness, proficiency level, native language, mode, genre, task type, and data annotation. In the next section, the design criteria of the corpus will be discussed.

3.2.2. Design Criteria

The purposes of corpus are usually divided into two major categories: general corpora and specific corpora. The former is likely to be used for a wide range of research features or by a large number of users, whereas the latter is used for examining particular features or by specific users. Although most of the corpora available in the literature designed for public use, a corpus can also comprise one or more purposes. For example, the purpose of The Japanese Learner of English Corpus (Izumi et al., 2004) is to support practitioners and researchers for using the corpus for “second language acquisition research, syllabus, and material design, or the development of computerized pedagogical tools, by combining it with NLP (Natural Language Processing) technology” (p.120). Similarly, the International Corpus of Learner English was designed to “make use of advances in applied linguistics and computer technology to effect a thorough investigation of the interlanguage of the foreign language learner” (Granger, 1993, p. 57). In this sense, the corpus collected as a part of the present study can be categorized as specific corpora since all of the texts used for building the corpus are research articles, and the aim of this corpus is to identify language-related difficulties of Turkish scholars in writing for publication process.

Another issue in corpus design is its size. A brief literature review reveals that sizes were relatively small at the beginning of learner corpus research. This view is supported by Granger (2003), who claims that a corpus having 200,000 words is considered as big in the field of SLA. However, “learner corpora tend to be rather large, which is a major asset in terms of representativeness of the data and generalizability of the results” (Granger, 2004, p. 125). On the other hand, Sinclair regards to size as an insignificant factor, arguing that the greatest size for a corpus does not exist and the minimum size depends on the following factors: “(a) the kind of query that is anticipated from users and (b) the methodology they

use to study the data” (p.10). In sum, what is important in the size issue is that the data in the corpus should effectively characterize the learner language. This corpus involves 216 research articles from 5 disciplines, which consist of 886,482 words and thus can be considered as both big and representative.

The target language is the language used to collect the corpus data; in other words, the language to be examined. Although the majority of the available corpora comprise data from only one language, some corpora contain more than one language. This shows that studies are likely to focus on one language instead of across languages. Similar to most of the studies in the literature, this corpus includes only one language: English.

Corpora are categorized into three major types in terms of data availability. First, there are corpora that are freely available on the Internet. Such corpora are intended to be for public use, and anyone having access to the Internet can search or download items. Corpus of Contemporary American English, Corpus of Global Web-Based English, the Michigan Corpus of Upper-level Student Papers (O’Donnell & Römer, 2009a, 2009b), and the Arabic Learners Written Corpus (Farwanah & Tamimi, 2012) are some of the examples of such corpora. The second category covers restricted corpora. Such corpora usually are limited to the use of a particular community, and other users have to pay for access. One needs to log in with a given username and password to access these kinds of corpora. For example, users willing to use The International Corpus of Learner English (Granger, 1993, 2003b; Granger et al., 2009), which is disseminated on CD-ROM, have to pay for access to the corpus. Finally, the third category covers the corpora that are under development. The Aachen Corpus of Academic Writing (ACAW) is an example of such corpora. Since they are still in progress, whether access to the corpora in this category will be available or restricted is not known at the moment. The Multidisciplinary Corpus of Writing for Publication built in this study can be considered as under development as only five disciplines are included in the corpus, and attempts to include more disciplines will be made.

Another issue in corpus design is nativeness. A corpus containing data produced by native speakers of a specific language is regarded as a *general corpus of NS*. On the other hand, a corpus is considered as a *non-native corpus* when the corpus consists of non-native data. Since the primary motivation behind building a non-native corpus is related to second language acquisition or teaching research (Granger, 2002), the data in such corpora are collected from learners of a second language. In this context, the primary focus of a learner corpus is on second language learners and their language production (Thoday, 2008). Since

the corpus built in this study consists of data from Turkish scholars regarding their unedited research articles in English, it can be classified as a non-native corpus.

The proficiency level is another issue in corpus design. Proficiency levels are typically defined according to traditional Beginner, Intermediate and Advanced classification. Nevertheless, there have been recent attempts to use the Common European Framework of Reference as the medium of classification. Although there are a number of corpora, including data from all levels, such as The SILS Learner Corpus of English, particular attention has been given to intermediate and especially advanced learners (Alfaifi, 2015). For example, The Taiwanese Corpus of Learner English (TLCE) and The LINCS Corpus include data from Intermediate and Advanced learners, while The Aachen Corpus of Academic Writing (ACAW) and The Advanced Learner English Corpus (ALEC) include only advanced-level learner data. Although a proficiency test was not administered to the participants of The Multidisciplinary Corpus of Writing for Publication, they are considered as intermediate to advanced learners for two reasons. First, all of the participants were university academic members. Researchers in Turkey need to have a minimum score of 60 in the national language proficiency exam for applying for a position in the universities. Score 60 is equivalent to the B1 level, Common European Framework of Reference for Languages (CEFR), according to YÖK. Second, a considerable amount of the participants was able to write and publish more than five articles in English. As writing and publishing several research articles requires a high level of proficiency in English, they can be regarded as advanced learners. For these reasons, the proficiency level in this study is considered intermediate to advanced.

A corpus can be classified into two categories in terms of learners' first language. Corpora in the literature have either several first languages or a single one. Generally, corpora that aim to carry out comparative studies tend to contain several first languages. This corpus contains data from only Turkish scholars, and thus it has a single first language.

The mode is another issue in corpus design. The term mode is related to the origin of the language data that is spoken or written (Sinclair, 2005). The procedure of building a corpus is, to some degree, is similar. The data need to be converted into a machine-readable textual format. Written data can easily be converted, thanks to the word processing software. Besides, there are some special programs such as Praat (Boersma & Weenink, 2014) and ELAN (Hellwig, 2019), enabling researchers to convert spoken data. Although there are spoken corpora and corpora, including both modes, exist in the literature, a great majority of

the existing corpora contain only written data (Alfaifi, 2015). Similar to most of the corpora, The Multidisciplinary Corpus of Writing for Publication includes only written data.

Genre is another part of the corpus design. A genre can be defined as any sort of communication in any mode with conventions developed through time on a social consensus. Genre selection during the compilation of a corpus can be problematic. In addition, all-inclusive genre taxonomy, which functions as a list to select, cannot be found. However, in his comprehensive study, Alfaifi (2015) lists 14 genres widely used in learner corpus research. He also states that Argumentative, Narrative, and Descriptive genres were the most preferred ones, respectively. Furthermore, a corpus may include only one genre as well as several genres. In this sense, The Multidisciplinary Corpus of Writing for Publication consists of only research articles and therefore includes one genre.

Another issue is task type; that is, how to collect the material in the corpus. The task type may differ according to the mode or genre of the material. In regard to written corpora, essays, tests, exams, and letters are the most frequent task types. With respect to spoken corpora, interviews seem to be the most common type of data collection. More than one task type may be used while building a corpus. However, the majority of the available corpora includes only one task type (Alfaifi, 2015). The reason behind including only one task type is to keep away from any misrepresentation of the results. However, this makes any task type based comparative analysis impossible. The Multidisciplinary Corpus of Writing for Publication includes a single task type. The data were obtained via e-mail.

Data annotation is the last issue in corpus design. A brief literature review reveals that although there are some corpora including only raw data, the majority of the corpora have one or more types of data annotation (Alfaifi, 2015). In terms of data annotation, researchers seem to have a high interest in particularly three kinds of data annotation: error annotation (Granger et al., 2009), Part-of-speech (POS) tagging (Boersma & Weenink, 2014), and structural tagging (Heuboeck et al., 2008). The aim of the first one is to enable researchers to carry out error analysis studies. Although several automatic tagging tools are available, the most common annotation technique is to tag the data manually according to certain criteria. The second one is frequently used in general corpora with an aim to point out the part of speech and further grammatical categories (Leech et al., 1994). As a result of noteworthy advancements in technology, there are a great number of PoS tagging tools developed for a number of different languages. The third one is used for tagging the structural aspects such as headings, sections, paragraphs, and titles. Such a tagging supports

researchers to analyze different functions such as certain parts or styles (Hammarberg, 2010). In terms of data annotation, The Multidisciplinary Corpus of Writing for Publication adopts error annotation and POS Tagging. Louvain Error Tagging Taxonomy was used for the error annotation. Stanford POS Tagger was used for the POS Tagging. The details of the Error Taxonomy and tagging process will be discussed in the following sections. Table 9 shows the overall design of The Multidisciplinary Corpus of Writing for Publication.

Table 9

Design Criteria of The Multidisciplinary Corpus of Writing for Publication

Purpose	Specific Use
Size	886,482 words
Target language	English
Availability	Under development
Learners' nativeness	Non-native
Learners' proficiency level	Advanced
Learners' first language	Turkish
Material mode	Written
Material genre	Research Article
Task type	Single
Data annotation	Error Tagged

3.2.3. Error Taxonomy

In order to increase the validity and reliability of the study, an error taxonomy that offers a compact and language-specific error classification was used. Although there are several error taxonomies in the literature (see Section 2.5.2.3.3), the taxonomy used in The International Corpus of Learner English (ICLE) was employed to increase the comparability of the present study. The International Corpus of Learner English error-tagging system contains eight main error categories, which are mostly broken down into subdomains. The main error categories are concisely described in Table 10. In addition, the full taxonomy is presented in Figure 15.

Table 10

Main Error Categories

Error Tags	Definition	Description
F	Formal Errors	Spelling or morphological errors that result in a non-existent English word
G	Grammatical Errors	Errors that break the general rules of English grammar
L	Lexical Errors	Errors involving the lexicosemantic properties of words or phrases (conceptual, collocational, or connotative)
X	Lexicogrammatical Errors	Errors that violate the lexicogrammatical properties of words; i.e., erroneous dependent prepositions, erroneous complementation patterns, or countable/uncountable noun confusion
Q	Punctuation Errors	Errors that target punctuation problems; e.g., confusion between punctuation markers, missing or redundant markers
W	Word Redundant/ Missing/Order Errors	Unnecessary use of words, missing necessary words, or misordered words
S	Style Errors	Sentence fragments and incomprehensible sentences
R	Register	Improper use of words, phrases, or grammatical rules

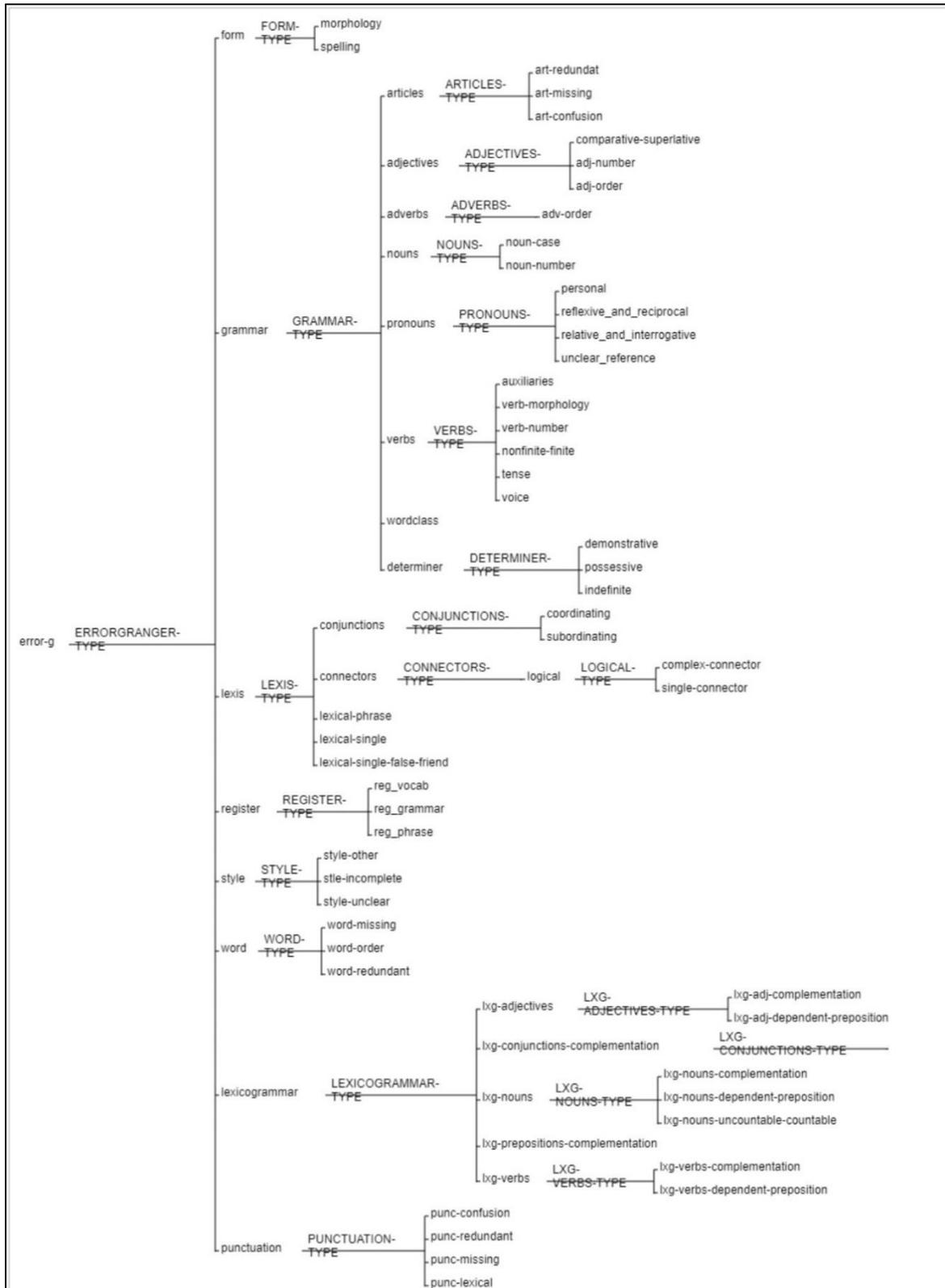


Figure 15. The Complete Error Tagging Taxonomy

3.2.4. Error Tagging Procedure

Error tagging is considered a tiresome, difficult, time, and energy-consuming task (Andersen, 2011). The tagging process requires spending an extensive amount of time, a number of human-based judgments, and managing a variety of problems. In computer-aided error analysis, texts are typically annotated manually using a software that enables error tagging. The software for error tagging is basically an editor for inserting the tags in the annotation procedure. Tag types are organized categorically on a menu-driven interface. Thus, users can move from the main categories to subcategories to achieve an appropriate definition of the erroneous element.

All of the text in the present corpus were manually and comprehensively annotated for errors using UAM CorpusTool. It is a free software and enables the multi-layer annotation of text files in a corpus. The tool was developed by Michael O'Donnell in 2008. It can be downloaded from <http://www.corpustool.com/download.html>. It allows researchers to manually tag the errors according to the error taxonomy that is built in its system or is designed for a particular study by the researchers. The tool also consists of a gloss section (coding criteria) related to each feature. This section enables researchers to select the correct category and thus contributes to inter-coder reliability.

The error tagging procedure followed three steps:

1. First, the erroneous element was selected. To ensure the reliability between the coders, only the erroneous element was coded instead of the entire grammatical unit.
2. Error codes were selected. The assigning procedure started with selecting the main error category appropriate for the error and then selecting more detailed subcategories. Considering the excessive number of errors in the taxonomy, this process enabled the researcher to easily focus on the suitable category without searching all of the categories.
3. A correction was provided.

Figure 16 shows an example of the UAM CorpusTool error tagging screen. The textual data appears in the upper part. Straightaway below, there are three boxes, the first two of which are related to error taxonomy and the last of which functions as a gloss, indicating the tagging criterion for the selected category. The function of the gloss section is important in that coders become familiarized with the error categories at the beginning of the tagging process. Finally, there are two boxes for writing down corrections and noting comments that coders may wish to discuss with each other at the bottom of the screen. The note section was

extremely useful in ensuring the inter-coder reliability. The tool enabled the coders to discuss the cases where coders did not agree and thus need to arrive at a consensus.

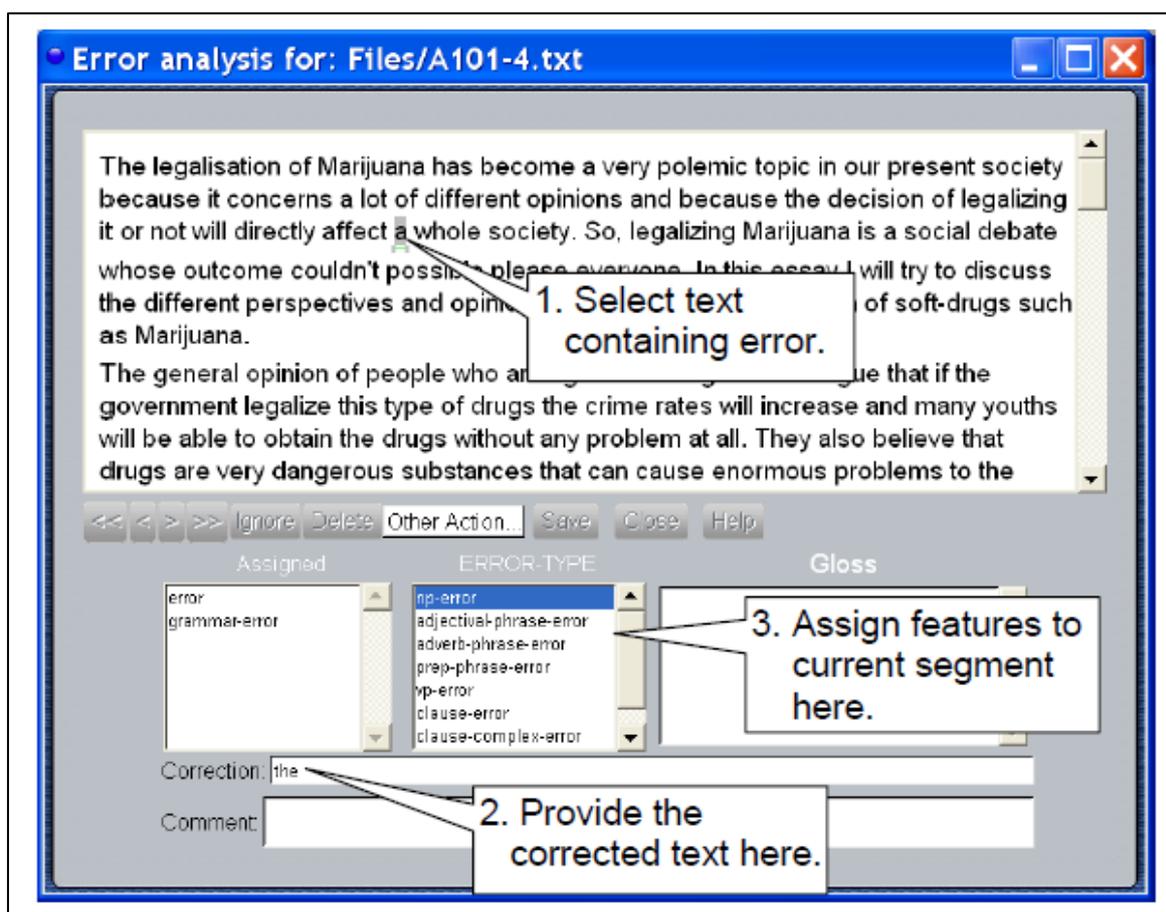


Figure 16. An example of the UAM corpustool error tagging screen

3.2.5. Piloting

The error taxonomy procedure was piloted before the tagging procedure began. Two coders, the researcher, and an experienced university lecturer, independently detected and tagged the errors. To do this, eight research articles that were not included in the corpus were selected randomly. The piloting was carried out in 3 stages. First, the researcher trained the second-rater on error taxonomy and how to use the error tagging software for tagging errors. After both of the coders were familiarized with the error taxonomy, error tagging, and the software, the coders independently detected and tagged the errors in the same articles selected for the pilot study using the software. The coders benefitted from the gloss and comment section described above during the tagging. Last, the two raters compared the errors and tags and discussed the cases where a disagreement occurred in order to reach a consensus. To calculate reliability, the inter-rater reliability was also conducted. Cohen's

Kappa coefficient was used for this purpose. Cohen's kappa (k) statistic is a widely used method for calculating inter-rater reliability (Stemler, 2004). Cohen's Kappa is "a chance-corrected measure, which takes into account the likelihood that the agreement between coders has occurred by chance" (Moreno & Swales, 2018, p. 55). The findings show that the level of agreement was 98.20%. Furthermore, the minimum Kappa coefficient was 0.77, while the maximum Kappa coefficient was 0.97. The Kappa statistic ranges from 0 to 1, where values below 0.40 are considered 'poor agreement'; 0.40 - 0.75 is 'fair to good agreement'; over 0.75 is 'excellent agreement.' The findings of the pilot study revealed that Cohen's Kappa coefficients ranged between 0.77 and 0.97, which demonstrates that there was excellent agreement in the study.

3.2.6. Inter-Rater Coding Reliability in the Study

The inter-rater coding reliability analysis was carried out in four stages. As done in the pilot study, the second coder was trained on error taxonomy and how to use the error tagging software for tagging errors. Then, the second coder independently detected and tagged the errors in randomly selected articles in the corpus. Next, the inter-rater coding reliability between the coders was measured using Cohen's Kappa coefficient. Finally, the coders evaluated the results of the inter-rater analysis and discussed the errors and tags to arrive at an agreement. 21 randomly selected texts, consisting of 52,725 words, were used to be tagged and to calculate interrater reliability as the use of 10% of the total data is considered sufficient for such an analysis (Lombard, Synder-Duch & Bracken, 2004, p. 4). Similar to the pilot study, Cohen's kappa (k) statistic was used to determine the degree of agreement. Cohen's Kappa was found to be 0.93 in the present study. In addition, the findings show that the level of agreement was 98%

3.3. Reference Corpus

A reference corpus is designed or collected to present comprehensive information about a language. A reference corpus should be large enough to characterize all the pertinent features of the language and the typical vocabulary in order to be used as a reliable basis. A reference corpus was collected in this study to compare the results of the MCWP. Therefore, published research articles written by native speakers of English were obtained. Since the data collection of the study took part between 2015 and 2017, attention was paid to obtain research articles published within the same period. In addition, there are a number of fields

within each discipline and these fields may have their own discourse communities and thus their own research patterns. Therefore, in order to make a reliable comparison, attention was also paid to include the research articles of the same fields within the disciplines. To do this, first, the keywords of the articles in the MCWP were identified to include research articles with the same fields in the reference corpus. Next, the keywords were searched in SSCI/SCI-Expanded / AHCI database to identify related journals. Then, all of the issues published between 2015 and 2017 were searched, and the research articles written by native speakers were downloaded. The list of the articles included in the reference corpus is presented in Appendices 1-5. Although research articles with one author were preferred, this was not possible sometimes as joint articles are quite widespread in some disciplines such as engineering and medicine. In such situations, research articles with a corresponding author who is a native speaker of English were included in the reference corpus. Attention was also paid to collect a reference corpus with the equal size of the MCWP. In sum, the reference Corpus includes 5 sub-corpora and consists of 163 research articles, which contain 885,791 words. Table 11 shows the disciplines, the number of research articles, and words in detail.

Table 11

The Distribution of Research Articles in the Reference Corpus

Discipline	Number of Research Articles	Number of Words
Communication	5	22.798
Economics and Administrative Sciences	12	80.147
Education	33	210.517
Engineering	71	379.959
Medicine	42	192.370
Total	163	885.791

3.4. Quantitative Analyses

Being one of the fastest-growing methodological fields, corpus studies are scientific ways of analyzing language, and thus, they need to provide empirical evidence to support any claims put forward about the data (Brezina, 2018). To do this, corpus studies basically adopt quantitative methodology in order to deal with numbers that indicate the frequencies of words and phrases (McEnery & Hardie, 2011). Besides, quantitative studies frequently present higher degrees of comparability, replicability, objectivity, and precision (Gries, 2013). Therefore, statistics are quite important in corpus studies since it aids the researchers to deal efficiently with the quantitative data (Brezina, 2018). Consequently, corpus studies

should include a number of different statistical tests to understand the phenomena under question. This section presents quantitative analyses used in this study.

3.4.1. Type / Token Ratio

Type-token ratio (TTR) is one of the simplest statistical measurements for lexical diversity. It gives the number of separate word forms (types) in relation to the number of running words (tokens). Nevertheless, TTR is quite responsive to text length; that is, it decreases in longer texts since more words are used repetitively in such texts. Therefore, it can be performed when texts having the same length are compared (Brezina, 2018). Although there are several ways to correct TTR for the length of the text, they were proved to be insufficient (Tweedie & Baayen, 1998). Hence, two different measures were put forward. The first one is a standardized type/token ratio (STRR), which computes the mean TTR for every N-words (for example, 1000 words) (Scott, 2004). The second one is the moving average type/ token ratio (MATTR) (Covington & McFall, 2010). Similar to STTR, MATTR uses a window of a fixed size which moves through the text and estimates the TTR for each window position. For these considerations, all measures of TTR were employed in this study.

3.4.2. Syntactic Complexity

Each article in both the MCWP and the reference corpus was analyzed using the L2 Syntactic Complexity Analyzer (L2SCA) developed by Lu (2010). L2SCA is a free software tool designed to examine the English writing syntactic complexity computing 14 indices of syntactic complexity. It uses the Stanford Parser (Klein & Manning, 2003) for Part-of-Speech Tagging, the Penn Treebank guidelines (Marcus, Santorini & Marcinkiewics, 1993) for sentence parsing, and Stanford Tregex (Levy & Andrew, 2006) for specified syntactic patterns. Lu (2010) presented a strong correlation, from .834 to 1.000, between the indices in the tool and human raters in a corpus of Chinese EFL learners. Similarly, Yoon and Polio (2017) reported high-reliability scores between L2SCA and manually coded essays. To calculate the indices, the tool produces frequency measures for the following units for each text file: word, sentence, verb phrase, clause, T-unit, dependent clause, complex T-Unit, coordinate phrase, and complex nominal. Table 12 shows the measures of syntactic complexity and their definitions. Once the values for all 14 indices were obtained, the data were exported to MS Excel for further statistical analysis.

Table 12

Definitions of Syntactic Complexity Measures

Category	Measure	Code	Definition
Length of the production unit	Mean length of clause	MLC	# of words/# of clauses
	Mean length of sentence	MLS	# of words/# of sentences
	Mean length of T-unit	MLT	# of words/# of T-units
Amount of subordination	Clauses per T-unit	C/T	# of clauses/# of T-unit
	Complex T-units per T-unit	CT/T	# of complex T-units/# of T-units
	Dependent clauses per clause	DC/C	# of dependent clauses/# of clauses
Amount of coordination	Dependent clauses per T-unit	DC/T	# of dependent clauses/# of T-units
	Coordinate phrases per clause	CP/C	# of coordinate phrases/# of clauses
	Coordinate phrases per T-unit	CP/T	# of coordinate phrases/# of T-units
Degree of phrasal sophistication	T-units per sentence	T/S	# of T-units/# of sentences
	Complex nominals per clause	CN/C	# of complex nominals/# of clauses
	Complex nominals per T-unit	CN/T	# of complex nominals/# of T-units
Overall sentence complexity	Verb phrases per T-unit	VP/T	# of verb phrases/# of T-units
	Clauses per sentence	C/S	# of clauses/# of sentences

3.4.3. Statistical Analysis

Statistical tests were employed in order to calculate descriptive statistics and compare groups using SPSS 23.0. To determine whether the data were distributed normally or not, Kolmogorov–Smirnov, and Shapiro–Wilk tests of Normality was carried out. It was found that the data did not follow a normal distribution ($p < 0.05$). However, since these tests are quite sensitive to sample size, that is even a small deviation from normality lead to a failure in normality test in a large sample size despite the fact that such a small deviation does not intervene to the results of parametric tests (Field, 2009; Oztuna et al., 2006), and the dataset used in the present study is relatively large, skewness and kurtosis values were also calculated to investigate the normality of the data. Skewness is used for measuring the symmetry of distribution, and kurtosis refers to the height and sharpness of the central peak. A skewness and kurtosis value between +1 and -1 is regarded as excellent, and the data that meet these criteria considered as normally distributed (George & Mallery, 2012; Hair et al., 2013). On the basis of this information, the data in the present study are treated as having a normal distribution. Therefore, parametric tests were used to compare the groups. Nevertheless, two groups, namely Communication and Economics corpora, had relatively small sample sizes, and thus, non-parametric tests were used for the comparison of these groups in order to obtain more reliable results though they had normal distributions. In addition, The Bonferroni correction was employed to eliminate the family-wise error rate in

the datasets. Therefore, the alpha value for each comparison was adjusted to $.05/14$, or $.0036$, where $.05$ is the significance level for the complete set of tests, and 14 is the number of tests.

3.4.3.1. Mann Whitney U

The Mann-Whitney U test, which is the nonparametric version of the parametric t-test, was used for comparing two groups on a single ordinal variable with no specific distribution (Mann & Whitney, 1947). Contrary to the t-test, it compares mean ranks rather than medians to determine the difference. It requires two independently sampled groups and assesses whether two groups differ on a single, continuous variable (Hollander & Wolfe, 1999). Therefore, the Mann-Whitney U test was used for comparing Communication and Economics corpora for each of the 14 syntactic complexity measures. $p < .0036$ significance level was set as the basis to interpret the statistical tests.

3.4.3.2. Independent Samples T-test

An independent samples t-test is a parametric test to examine whether a statistically significant difference exists between two independent groups. It compares the means of these groups to determine the difference (Levine, 2014). Therefore, a set of independent samples t-tests were used for comparing the MCWP and the reference corpus as well as the subcorpora for each of the 14 syntactic complexity measures. $p < .0036$ significance level was set as the basis to interpret the statistical tests.

3.4.3.3. ANOVA

One-way analysis of variance (ANOVA) is a statistical test to find out whether any statistically significant differences exist among three or more independent groups. Similar to the independent samples t-test, ANOVA uses means of the groups for comparison. The means are considered as equal in ANOVA, and thus statistically significant results indicate that there is a difference in the means of the group (Fisher, 1992). Therefore, ANOVA was used to compare the Turkish corpora with regard to their syntactic complexity measures. Nevertheless, as an omnibus test, ANOVA does not indicate which groups significantly differed from one another. Therefore, post-hoc tests should be performed to determine the differences between particular groups (Hilton & Armstrong, 2006). As the sizes of the groups in the present dataset were not equal, Games-Howell test was used as post-hoc since

equal sample sizes and variances are not an assumption of it (Malhotra & Krishina, 2018). $p < .0036$ significance level was set as the basis to interpret both ANOVA and Games-Howell tests.





CHAPTER IV

RESULTS

This chapter reports the results of the type/token ratio, syntactic complexity measures, and computer-aided error analysis as well as the comparison of subcorpora. First, the results of the MCWP treated as a separate body are presented and then compared with the reference corpus. Second, results regarding each sub-corpus in the MCWP are presented and compared with its counterpart in the reference corpus. Lastly, the result of the comparison of the sub-corpora in the MCWP with each other is presented.

4.1. The Multidisciplinary Corpus of Writing for Publication

In the present study, a corpus of unedited research articles written by Turkish scholars was compiled. The results regarding the type/token ratio, syntactic complexity, and computer-aided error analysis of the MCWP are presented in this section. In addition, the results of the MCWP are compared with a reference corpus to identify the differences.

4.1.1. Type/Token Ratio

Table 13 shows the type/token ratio of the MCWP and the reference corpus. The MCWP consists of 889,648 words, 51,793 types, and 851,643 tokens. In that sense, the type-token ratio of the MCWP was found to be 0.06. For the considerations explained in the method section, additional STTR and MATTR were also calculated. Both STTR and MATTR of the MCWP were calculated as 0.68. On the other hand, the reference corpus consists of 885,791 words, 45,608 types, and 879,552 tokens. The type-token ratio of the reference corpus was 0.052, and both STTR and MATTR were found to be 0.70.

Table 13

Type/Token Ratio of the MCWP and Reference Corpus

	MCWP	Reference Corpus
Total Words	889,648	885,791
Types	51,793	45,608
Tokens	851,643	879,552
TTR	0.06	0.052
STTR	0.68	0.70
MATTR	0.68	0.70

On the basis of the fact that a higher ratio indicates a higher lexical diversity (Biber, Conrad & Leech, 2002), the findings reveal that Turkish scholars used a wide range of vocabulary in their research articles. However, when compared to the reference corpus, the vocabulary they used is slightly less varied.

4.1.2. Syntactic Complexity

Syntactic complexity was measured automatically using 14 indices in 5 broad categories, which are the length of the production unit, amount of subordination, amount of coordination, degree of phrasal sophistication, and overall sentence complexity.

To measure the length of the production unit, three indices were used: mean length of clause (MLC), mean length of T-Unit (MLT), and mean length of sentence (MLS). Figure 17 indicates the means of these measures. Accordingly, the MLC, MLT, and MLS scores of the MCWP were found to be 14.6, 21.65, and 23.21, respectively. In contrast, the scores of the reference corpus in the same indices were 16.15, 26.27, and 28.61, respectively. The results show that the reference corpus produced longer sentences, T-units, and clauses than the MCWP.

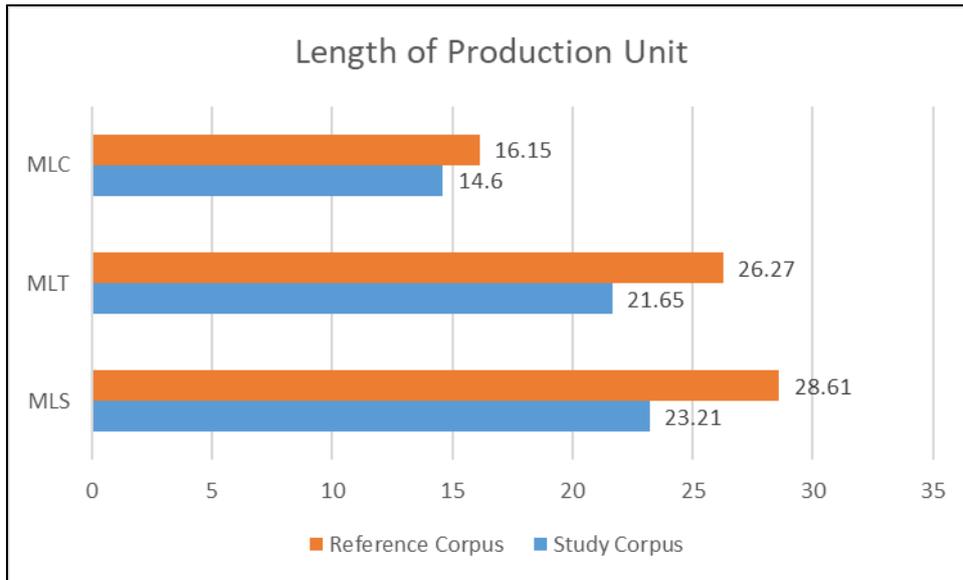


Figure 17. The comparison of mean values of the length of production unit measures

The amount of subordination was measured using the following four indices: clauses per T-unit (C/T), complex T-units per T-unit (CT/T), dependent clauses per clause (DC/C), and dependent clauses per T-unit (DC/T). Figure 18 presents the mean scores of these indices. The results show that the scores the MCWP had in these measures were 0.45, 0.29, 0.34, and 1.48, respectively. On the contrary, the reference corpus had means of 0.60, 0.36, .043, and 1.63, respectively. These findings suggest that the MCWP had lower mean values than the reference corpus in all of the indices.

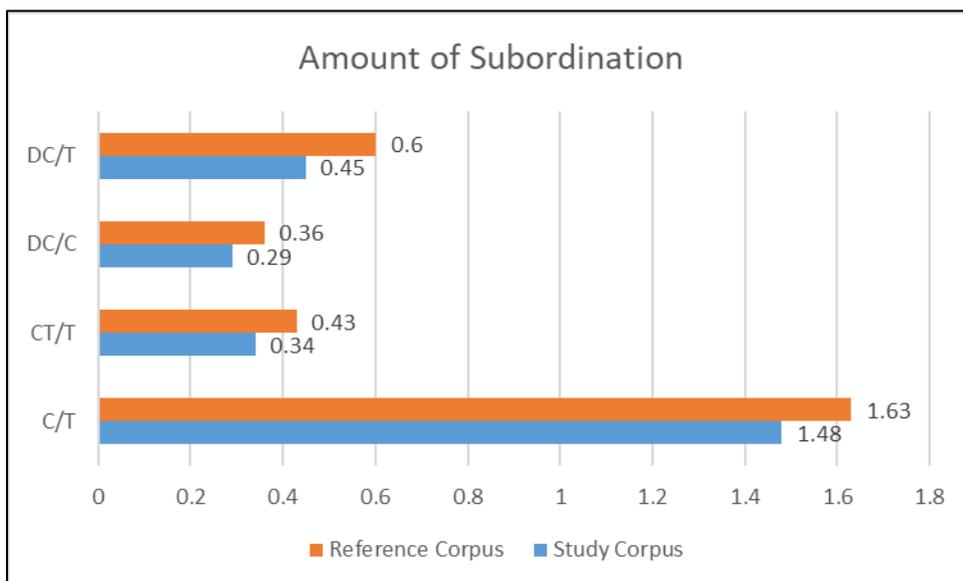


Figure 18. The comparison of mean values of the the amount of subordination measures

The amount of coordination was measured using three indices: coordinate phrases per clause (CP/C), coordinate phrases per T-unit (CP/T), and T-units per sentence (T/S). As presented in Figure 19, the mean scores of the MCWP were 1.05, 0.67, and 0.45, respectively. On the other hand, the reference corpus had 1.09, 0.70, and 0.43. The findings show that the reference corpus had slightly higher means in T/S and CP/T, while the MCWP had a higher mean value in CP/C.

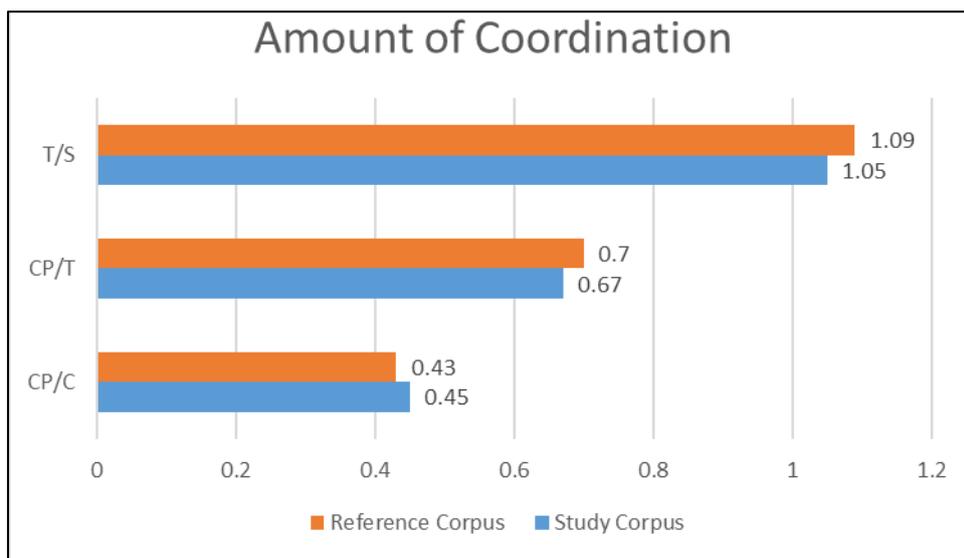


Figure 19. The comparison of mean values of the amount of coordination measures

The degree of sophistication was measured using the following indices: complex nominals per clause (CN/C), complex nominal per T-unit (CN/T), and verb phrases per T-unit (VP/T). Figure 20 shows the mean values of the MCWP and the reference corpus. Accordingly, the CN/C, CN/T, and VP/T means of the MCWP were 1.93, 2.80, and 1.89, respectively. In contrast, the reference corpus had means of 2.21, 3.44, and 2.12, respectively. Based on these findings, the reference group used more sophisticated phrases than the MCWP in all of the indices in this category.

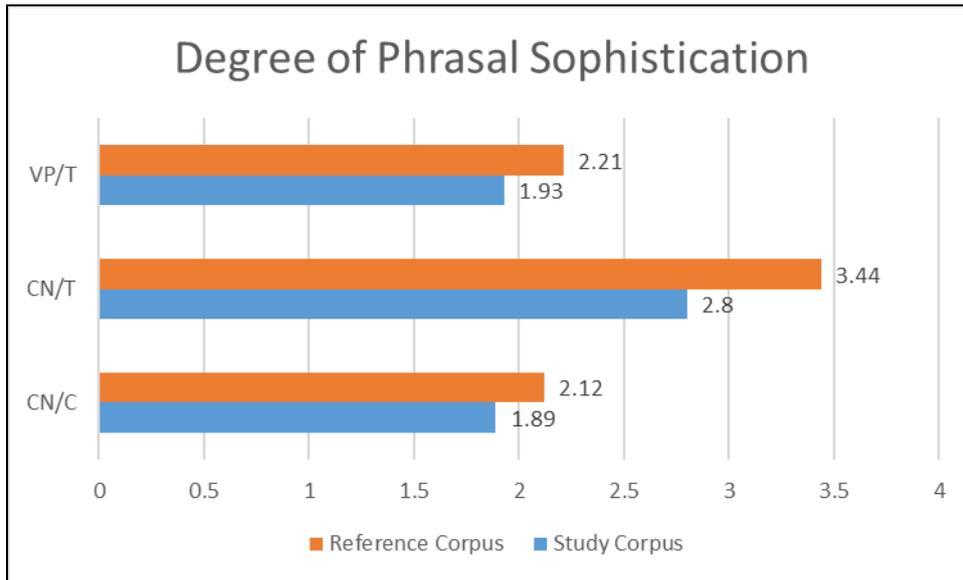


Figure 20. The comparison of mean values of the degree of phrasal sophistication measures

Finally, overall sentence complexity was measured using the clauses per sentence (C/S) index. Figure 21 shows the means of the MCWP and the reference corpus. The MCWP had a mean of 1.58, while the reference corpus had 1.79. This finding shows that the reference corpus had a higher level of overall sentence complexity.

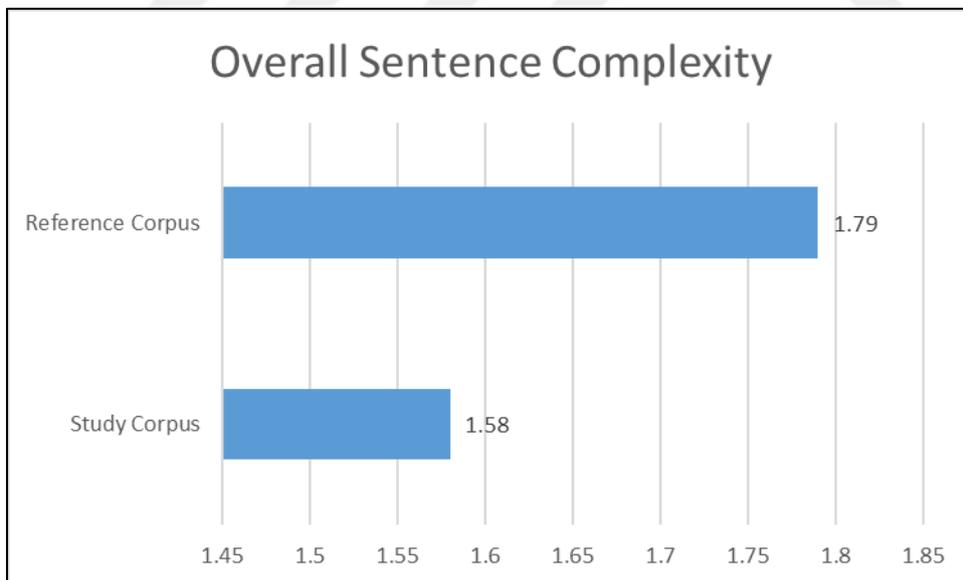


Figure 21. The comparison of mean values of the overall sentence complexity measure

A set of independent-samples t-tests were run to investigate whether there are any statistically significant differences between the MCWP and the reference corpus for each of the syntactic complexity measures. In addition, The Bonferroni correction was employed to eliminate the family-wise error rate in the dataset. Therefore, the alpha value for each

comparison was adjusted to .05/14, or.0036, where .05 is the significance level for the complete set of tests, and 14 is the number of tests performed. Table 14 presents the t-test results. The results suggest that the MCWP and the reference corpus differed significantly in 12 out of 14 syntactic complexity measures. The two measures that were not statistically significant were coordinate phrases per clause ($p=.225$) and coordinate phrases per T-unit ($p=.157$). The only category the MCWP and the reference corpus did not differ significantly was the amount of coordination. These findings demonstrate that Turkish scholars used less syntactically complex sentences than scholars in the reference corpus.

Table 14

The Comparison of the T-Test Results of the Syntactic Complexity Measures

			N	Mean	Std. Deviation	t	p
Length of the production unit	MLS	Study	212	23,21	4,99	10,263	<.001*
		Reference	163	28,61	5,13		
	MLT	Study	212	21.65	4.15	10,179	<.001*
		Reference	163	26.27	4.60		
	MLC	Study	212	14.60	2.30	-6,101	<.001*
		Reference	163	16.15	2.60		
Amount of subordination	C/T	Study	212	1.48	0.19	-7,233	<.001*
		Reference	163	1.63	0.22		
	CT/T	Study	212	0.34	0.10	-7,604	<.001*
		Reference	163	0.43	0.10		
	DC/C	Study	212	0.29	0.07	-7,947	<.001*
		Reference	163	0.36	0.08		
	DC/T	Study	212	0.45	0.17	-7,799	<.001*
		Reference	163	0.60	0.19		
Amount of coordination	CP/C	Study	212	0.45	0.15	1,216	.225
		Reference	163	0.43	0.15		
	CP/T	Study	212	0.67	0.22	-1,417	.157
		Reference	163	0.70	0.25		
	T/S	Study	212	1.05	0.09	-4,266	<.001*
		Reference	163	1.09	0.06		
Degree of phrasal sophistication	CN/C	Study	212	1.89	0.30	-6,844	<.001*
		Reference	163	2.12	0.35		
	CNT	Study	212	2.80	0.56	11,108	<.001*
		Reference	163	3.44	0.54		
	VP/T	Study	212	1.93	0.31	-8,550	<.001*
		Reference	163	2.21	0.32		
Overall sentence complexity	C/S	Study	212	1.58	0.25	-7,656	<.001*
		Reference	163	1.79	0.27		

Note: * $p<.0036$; see Table 12 for the definitions.

4.1.3. Computer-Aided Error Analysis

A computer-aided error analysis was carried out to provide a general picture of the Turkish scholars' academic writing problems. Louvain Error Tagging Taxonomy was used to investigate errors Turkish scholars made in their manuscripts. The Louvain taxonomy examines the errors in eight broad categories, which are further divided into subdomains. Therefore, first, the results of the main categories and then the results of subcategories are presented.

Table 15 shows the number of errors. A total of 23,253 errors were identified in the corpus. Turkish scholars made the greatest number of errors in the Grammar category with 11,331 errors (48.73%) and the fewest number of errors in the Style category with 524 errors (2.25%). Punctuation was the second most frequent error category with 2,789 errors which equals to 11.99% of all errors, followed by Register (N=2,153), Lexis (N=2,011), Word (N=1,793), Lexicogrammar (N=1,540) and Form (N=1,112) errors, respectively.

Table 15

The Computer-Aided Error Analysis Results of the MCWP

	N	%	Per 1000 tokens
Form	1,112	4.78%	0.96
Grammar	11,331	48.73%	9.78
Lexis	2,011	8.65%	1.74
Register	2,153	9.26%	1.86
Style	524	2.25%	0.45
Word	1,793	7.71%	1.55
Lexicogrammar	1,540	6.62%	1.33
Punctuation	2,789	11.99%	2.41
Total	23,253		

A total of 1,112 errors were identified in the Form category. Table 16 shows the subcategories and distribution of the errors. A great majority of the errors in the Form category belong to spelling (N=1,103), which is equal to 4.74% of all errors, while nine errors were made in the morphology subcategory.

Table 16

The Distribution of Form Errors in the MCWP

Form	N	%	Per 1000 tokens
Morphology	9	0.04%	0.01
Spelling	1103	4.74%	0.95
Total	1112	4.78%	0.96

Table 17 demonstrates the subcategories and the distribution of errors in the Grammar category. It is seen that the most frequent error type was articles, which equals to 23.33% of all errors. Of 5,425 errors made in the Articles subcategory, 4,552 were missing, 747 were redundant, and 126 were confusion errors. For the Adjectives (N=160), 78 comparative-superlative errors and 82 order errors were tagged. Further, 129 adverb errors were made by Turkish scholars. In Nouns subcategory (N=1,325), 1132 number and 193 case errors were made. A total of 237 errors were made in the Pronouns category, which is divided into personal (N=82), reflexive and reciprocal (N=43), relative and interrogative (N=42), and unclear reference (N=70) errors. In the Verb subcategory (N=3523), 2112 tense errors, 569 voice errors, 511 verb-number errors, 173 finite-nonfinite errors, 130 auxiliaries errors, and 28 verb-morphology errors were made. In addition, 337 errors were made in the Word Class category. Finally, 92 demonstratives, 35 possessives, and 68 indefinite errors were detected in the Determiner subcategory (N=195).

Table 17

The Distribution of Grammar Errors in the MCWP

	Grammar	N	%	Per 1000 tokens
Articles	Redundant	747	3.21%	0.64
	Missing	4552	19.58%	3.93
	Confusion	126	0.54%	0.11
	Total	5425	23.33%	4.68
Adjectives	Comparative-superlative	78	0.34%	0.07
	Order	82	0.36%	0.08
	Total	160	0.69%	0.14
Adverbs		129	0.55%	0.11
Nouns	Case	193	0.83%	0.17
	Number	1132	4.87%	0.98
	Total	1325	5.70%	1.14
Pronouns	personal	82	0.35%	0.07
	Reflexive and reciprocal	43	0.18%	0.04
	Relative and interrogative	42	0.18%	0.04
	Unclear reference	70	0.30%	0.06
	Total	237	1.02%	0.20
Verbs	Auxiliaries	130	0.56%	0.11
	Verb-morphology	28	0.12%	0.02
	Verb-number	511	2.20%	0.44
	Nonfinite-finite	173	0.74%	0.15
	Tense	2112	9.08%	1.82
	Voice	569	2.45%	0.49
	Total	3523	15.15%	3.04
Word Class		337	1.45%	0.29
Determiner	Demonstrative	92	0.40%	0.08
	Possessive	35	0.15%	0.03
	Indefinite	68	0.29%	0.06
	Total	195	0.84%	0.17
	Total	11331	48.73%	9.78

Table 18 shows the subcategories and the distribution of errors in the Lexis category. It is seen that most of the errors were made in the lexical-single (N=1428) subcategory, which equals to 6.14% of all errors. With a total of 308 errors, 244 coordinating, and 64 subordinating errors, the conjunction category was found to be the second category. In addition, Turkish scholars made 140 lexical phrases. The fewest number of errors were annotated in connectors with 135 errors. These findings demonstrate that the main lexical problem of Turkish scholars is using the proper words in their manuscripts.

Table 18

The Distribution of Lexis Errors in the MCWP

Lexis	N	%	Per 1000 tokens	
Conjunctions	Coordinating	244	1.05%	0.21
	Subordinating	64	0.28%	0.06
	Total	308	1.32%	0.27
Connectors	135	0.58%	0.12	
Lexical-phrase	140	0.60%	0.12	
Lexical-single	1428	6.14%	1.23	
Total	2011	8.65%	1.74	

Another error tagging category was Register. As presented in Table 19, the Register category is divided into three subcategories: vocabulary, grammar, and phrase. The greatest number of errors in this category was within the vocabulary subcategory (N=1222). The second-largest group was grammar (N=564), which was followed by phrase errors (N=367).

Table 19

The Distribution of Register Errors in the MCWP

Register	N	%	Per 1000 tokens
Vocabulary	1222	5.26%	1.05
Grammar	564	2.43%	0.49
Phrase	367	1.58%	0.32
Total	2153	9.26%	1.86

Compilation and categorization of errors show that errors in the Style category account for 2.25 % of the total number of errors in the corpus. Table 20 shows the distribution of style errors. Accordingly, Unclear (N=391) is the most frequent error type in this category. Besides, Turkish scholars made 81 errors in the incomplete subcategory and 52 errors in other subcategories.

Table 20

The Distribution of Style Errors in the MCWP

Style	N	%	Per 100 tokens
Other	52	0.22%	0.04
Incomplete	81	0.35%	0.07
Unclear	391	1.68%	0.34
Total	524	2.25%	0.45

A total of 1,1793 errors were annotated in the Word category. Table 21 shows the distribution of errors in this category. The major subcategory was Redundant (N=1106),

followed by Order (N=348) and Missing (N=339). The findings demonstrate that the main problem Turkish scholars experienced in this category was using unnecessary words in their manuscripts.

Table 21

The Distribution of Word Errors in the MCWP

Word	N	%	Per 1000 tokens
Missing	339	1.46%	0.29
Order	348	1.50%	0.30
Redundant	1106	4.76%	0.95
Total	1793	7.71%	1.55

Table 22 shows the subcategories and the distribution of errors in the Lexicogrammar category. The errors in this category make up 6.62 % of the total number of errors. The greatest number of errors in this category was within the verbs subcategory (N=597). The second-largest group was prepositions (N=531), which was followed by nouns (N=328) adjectives (N=46), and conjunctions (N=38).

Table 22

The Distribution of Lexicogrammar Errors in the MCWP

Lexicogrammar	N	%	Per 1000 tokens	
Adjectives	Complementation	9	0.04%	0.01
	Dependent preposition	37	0.16%	0.03
	Total	46	0.20%	0.04
Conjunctions	38	0.16%	0.03	
Nouns	Complementation	178	0.77%	0.15
	Dependent preposition	144	0.62%	0.12
	Countable / Uncountable	6	0.03%	0.01
	Total	328	1.41%	0.28
Prepositions	531	2.28%	0.46	
Verbs	Complementation	207	0.89%	0.18
	Dependent preposition	390	1.68%	0.34
	Total	597	2.57%	0.52
Total	1540	6.62%	1.33	

The last category is Punctuation. Compilation and categorization of errors show that errors in the Punctuation category (N=2789) account for 11.99 % of the total number of errors in the corpus. It is seen in Table 23 that the largest subcategory was Missing with 1,602 errors and the fewest subcategory was lexical with 38 errors. In addition, 790 Redundant errors and

359 confusion errors were annotated in this category. The findings suggest that Turkish scholars had some problems utilizing the correct punctuation.

Table 23

The Distribution of Punctuation Errors in the MCWP

Punctuation	N	%	Per 1000 tokens
Confusion	359	1.54%	0.31
Redundant	790	3.40%	0.68
Missing	1602	6.89%	1.38
Lexical	38	0.16%	0.03
Total	2789	11.99%	2.41

4.2. Communication Subcorpora

4.2.1. Type/Token Ratio

Table 24 shows the type/token ratio of the Turkish Communication (TR_COM) corpus and its counterpart (REF_COM) in the reference corpus. TR_COM corpus consists of 25,833 words, 4,438 types and 25,706 tokens. In that sense, the type-token ratio of the MCWP was found to be 0.17. For the considerations explained in the method section, additional STTR and MATTR were also calculated. Both STTR and MATTR of the TR_COM corpus were calculated as 0.70. On the other hand, the REF_COM corpus consists of 22,798 words, 4,606 types, and 22,668 tokens. The type-token ratio of the reference corpus was calculated as 0.20, and both STTR and MATTR were 0.74.

Table 24

Type/Token Ratio Results of Communication Subcorpora

	TR_COM	REF_COM
Total Words	25,833	22,798
Types	4,438	4,606
Tokens	25,706	22,668
TTR	0.17	0.20
STTR	0.70	0.74
MATTR	0.70	0.74

The findings demonstrate that the TTR values of both Communication subcorpora were above those of the main corpora. In addition, it can be said that participants in these corpora corpus used a wide range of vocabulary in their research articles. However, a comparison of

the TR_COM and REF_COM corpora reveals that Turkish scholars used less varied vocabulary than their counterparts.

4.2.2. Syntactic Complexity

Table 25 shows the mean values of 14 syntactic complexity measures for Communication Corpora. It can be seen in the table that the REF_COM corpus had higher mean values 13 out of 14 syntactic complexity measures. The only measure that TR_COM corpus had a higher mean is the clauses per sentence index (C/S). This finding demonstrates that participants in REF_COM subcorpus made use of more syntactically complex sentences in their manuscripts.

Table 25

Means of Syntactic Complexity Measures for Communication Corpora

			N	Mean	Std. Deviation
Length of production unit	MLS	Study	7	25,58	6,26
		Reference	5	29,53	1,46
	MLT	Study	7	22.75	4.32
		Reference	5	28.08	1.67
	MLC	Study	7	13.15	3.26
		Reference	5	15.95	0.65
Amount of subordination	C/T	Study	7	1.63	0.14
		Reference	5	1.78	0.09
	CT/T	Study	7	0.43	0.06
		Reference	5	0.49	0.05
	DC/C	Study	7	0.39	0.08
		Reference	5	0.41	0.04
	DC/T	Study	7	0.61	0.14
		Reference	5	0.74	0.10
Amount of coordination	CP/C	Study	7	0.36	0.10
		Reference	5	0.56	0.12
	CP/T	Study	7	0.63	0.12
		Reference	5	0.89	0.09
	T/S	Study	7	0.96	0.37
		Reference	5	1.06	0.06
Degree of phrasal sophistication	CN/C	Study	7	1.87	0.70
		Reference	5	2.09	0.16
	CN/T	Study	7	3.19	0.90
		Reference	5	3.72	0.44
	VP/T	Study	7	2.09	0.26
		Reference	5	2.48	0.14
Overall sentence complexity	C/S	Study	7	1.95	0.33
		Reference	5	1.86	0.13

In order to see whether the two communication corpora significantly differed or not, Mann-Whitney U tests were applied for each of the syntactic complexity measures. Although data sets had a normal distribution, the sample sizes did not allow us to run an independent samples t-test, and thus we had to use non-parametric statistical analysis. In addition, The Bonferroni correction was employed to eliminate the family-wise error rate in the dataset. Therefore, the alpha value for each comparison was adjusted to $.05/14$, or $.0036$, where $.05$ is the significance level for the complete set of tests, and 14 is the number of tests. Table 26 shows the results of the Mann-Whitney U test. No significant differences were found for any measures. Although REF_COM had higher means in 13 out of 14 syntactic complexity measures, they were not found to be statistically significant.



Table 26

Mann Whitney U Results of Communication Corpora

			N	Mean Rank	Sum of Ranks	U	P
Length of production unit	MLS	Study	7	5,43	38,00	10,000	.223
		Reference	5	8,00	40,00		
	MLT	Study	7	4,86	34,00	6,000	.062
		Reference	5	8,80	44,00		
	MLC	Study	7	4,71	33,00	5,000	.042
		Reference	5	9,00	45,00		
Amount of subordination	C/T	Study	7	6,57	46,00	17,000	.935
		Reference	5	6,40	32,00		
	CT/T	Study	7	6,21	43,50	15,500	.745
		Reference	5	6,90	34,50		
	DC/C	Study	7	6,29	44,00	16,000	.808
		Reference	5	6,80	34,00		
	DC/T	Study	7	6,57	46,00	17,000	.935
		Reference	5	6,40	32,00		
Amount of coordination	CP/C	Study	7	4,71	33,00	5,000	.042
		Reference	5	9,00	45,00		
	CP/T	Study	7	4,14	29,00	1,000	.007
		Reference	5	9,80	49,00		
	T/S	Study	7	7,86	55,00	8,000	.123
		Reference	5	4,60	23,00		
Degree of phrasal sophistication	CN/C	Study	7	5,00	35,00	7,000	.088
		Reference	5	8,60	43,00		
	CN/T	Study	7	5,29	37,00	9,000	.167
		Reference	5	8,20	41,00		
	VP/T	Study	7	5,71	40,00	12,000	.372
		Reference	5	7,60	38,00		
Overall sentence complexity	C/S	Study	7	6,43	45,00	17,000	.935
		Reference	5	6,60	33,00		

Note: * $p < .0036$; see Table 12 for the definitions.

4.2.3. Computer-Aided Error Analysis

Table 27 shows the error categories, the total number of errors, and the number of errors per 1000 tokens in the TR_COM subcorpus. A total of 971 errors were annotated in the TR_COM subcorpus. It is seen that the greatest number of errors was detected in Grammar (N= 355), which is followed by Punctuation (N=168) and Lexis (N=131). On the other hand, the fewest number of errors were identified in the Style category with 29 errors. A normalized number of errors shows that 31.48 errors were tagged per 1000 tokens. On the basis of the fact that the mean length of sentence in TR_COM subcorpus was found to be

25.58, it can be argued that 1000 tokens account for 39.09 sentences. Therefore 31.48 errors were annotated in 39.09 sentences. In addition, considering that seven research articles were included in this subcorpus, each research article averaged 138.7 errors.

Table 27

The Computer-Aided Error Analysis Results of TR_COM Subcorpus

	N	%	Per 1000 tokens
Form	30	3.09%	0.97
Grammar	355	36.56%	11.51
Lexis	131	13.49%	4.25
Register	95	9.78%	3.08
Style	29	2.99%	0.94
Word	71	7.31%	2.30
Lexicogrammar	92	9.47%	2.98
Punctuation	168	17.30%	5.45
Total	971	100%	31.48

An in-depth examination of Form errors showed that all of the errors tagged in this category belong to the spelling subcategory. The normalized number of Form errors was calculated as 0.97. In other words, there is a 0.97 error per 1000 tokens in the TR_COM subcorpus.

Table 28

The Distribution of Form Errors in TR_COM Subcorpus

Form	N	%	Per 1000 tokens
Morphology	0	0.00%	0.00
Spelling	30	3.09%	0.97
Total	30	3.09%	0.97

Table 29 shows the distribution of Grammar errors in the TR_COM subcorpus. Categorization and annotation of errors show that errors in the Grammar category (N=355) account for 36.56% of the total number of errors in the TR_COM subcorpus. The major subcategory was found to be Articles (N=171), followed by Verbs (N=94) and Nouns (N=57). An investigation of articles errors showed that 144 out of 171 were missing article errors. The distribution of Verb errors demonstrates that the 37 tense, 27 number, 21 voice, five nonfinite-finite, and four auxiliaries errors were tagged. In the Noun category, 42 number and 15 case errors were detected.

Table 29

The Distribution of Grammar Errors in TR_COM Subcorpus

Grammar	N	%	Per 1000 tokens	
Articles	Redundant	25	2.57%	0.81
	Missing	144	14.83%	4.67
	Confusion	2	0.21%	0.06
	Total	171	17.61%	5.55
Adjectives	Comparative-superlative	3	0.31%	0.09
	Order	1	0.10%	0.03
	Total	4	0.41%	0.13
Adverbs	1	0.10%	0.03	
Nouns	Case	15	1.54%	0.49
	Number	42	4.33%	1.36
	Total	57	5.87%	1.85
Pronouns	Personal	4	0.41%	0.13
	Reflexive and reciprocal	2	0.21%	0.06
	Relative and interrogative	0	0.00%	0.00
	Unclear reference	0	0.00%	0.00
	Total	6	0.62%	0.19
Verbs	Auxiliaries	4	0.41%	0.13
	Verb-morphology	0	0.00%	0.00
	Verb-number	27	2.78%	0.88
	Nonfinite-finite	5	0.51%	0.16
	Tense	37	3.81%	1.20
	Voice	21	2.16%	0.68
	Total	94	9.68%	3.05
Word Class	9	0.93%	0.29	
Determiner	Demonstrative	5	0.51%	0.16
	Possessive	0	0.00%	0.00
	Indefinite	8	0.82%	0.26
	Total	13	1.34%	0.42
Total	355	36.56%	11.51	

A further investigation of Lexis errors, presented in Table 30, shows that the greatest number of the errors in this category was made in the lexical-single subcategory (N=103). Furthermore, 20 conjunctions, four connectors, and four lexical phrase errors were identified in this category. The errors in this category equal to 13.49% of all of the errors in this subcorpus and 4.25 errors per 1000 tokens.

Table 30

The Distribution of Lexis Errors in TR_COM Subcorpus

Lexis	N	%	Per 1000 tokens	
Conjunctions	Coordinating	17	1.75%	0.55
	Subordinating	3	0.31%	0.10
	Total	20	2.06%	0.65
Connectors	4	0.41%	0.13	
lexical-phrase	4	0.41%	0.13	
lexical-single	103	10.61%	3.34	
Total	131	13.49%	4.25	

The register is another error category. The number of errors identified in this corpus is presented in Table 31. The errors in this category make up 9.78% of the total number of errors in this subcategory. It is seen that 46 vocabulary errors, 30 grammar errors, and 19 phrase errors were tagged. In sum, register errors were made 3.08 times per 1000 tokens in the TR_COM subcorpus.

Table 31

The Distribution of Register Errors in TR_COM Subcorpus

Register	N	%	Per 1000 tokens
Vocabulary	46	4.74%	1.49
Grammar	30	3.09%	0.97
Phrase	19	1.96%	0.62
Total	95	9.78%	3.08

Table 32 shows the subcategories and the distribution of errors in the Style category. It is seen that most of the errors were made in the unclear (N=24) subcategory, which equals to 2.47% of all errors in this subcorpus. In addition, one error in another category and four errors in the incomplete category were tagged. Overall, Turkish scholars in this subcorpus made 0.94 style errors per 1000 tokens.

Table 32

The Distribution of Style Errors in TR_COM Subcorpus

Style	N	%	Per 1000 tokens
Other	1	0.10%	0.03
Incomplete	4	0.41%	0.13
Unclear	24	2.47%	0.78
Total	29	2.99%	0.94

A total of 71 errors were identified in the Word category. Table 33 shows the subcategories and distribution of the errors. Thirty-six redundant, 20 missing, and 15 order errors were annotated in this subcorpus. It can be seen that there were 0.94-word errors per 1000 tokens.

Table 33

The Distribution of Word Errors in TR_COM Subcorpus

Word	N	%	Per 1000 tokens
Missing	20	2.06%	0.65
Order	15	1.54%	0.49
Redundant	36	3.71%	1.17
Total	71	7.31%	2.30

Compilation and categorization of errors show that errors in the Lexicogrammar (N=92) category account for 9.47 % of the total number of errors in the corpus. As shown in Table 34, the greatest number of errors was tagged in the verbs (N=50) subcategory, which was followed by prepositions (N=20) and Nouns (N=18), respectively. In addition, one adjective and three conjunctions errors were detected in this category. Examination of normalized occurrences of errors demonstrates that Lexicogrammar errors were annotated 2.98 times per 1000 tokens.

Table 34

The Distribution of Lexicogrammar Errors in TR_COM Subcorpus

Lexicogrammar	N	%	Per 1000 tokens	
Adjectives	Complementation	1	0.10%	0.03
	Dependent preposition	0	0.00%	0.00
	Total	1	0.10%	0.03
Conjunctions	3	0.31%	0.10	
Nouns	Complementation	12	1.24%	0.39
	Dependent preposition	6	0.62%	0.19
	Countable / Uncountable	0	0.00%	0.00
	Total	18	1.85%	0.58
Prepositions	20	2.06%	0.65	
Verbs	Complementation	24	2.47%	0.78
	Dependent preposition	26	2.68%	0.84
	Total	50	5.15%	1.62
Total	92	9.47%	2.98	

Punctuation was found to be the second category; the greatest number of errors were tagged in the TR_COM subcorpus. An investigation of the distribution of Punctuation errors reveals that the subcategory Turkish scholar made the most errors was the missing punctuations

(N=118), which equals to 12.15% of all errors in this category. In addition, 39 errors in the redundant subcategory and 11 errors in the confusion subcategory were tagged. The normalized occurrences of Punctuation errors were found to be 5.45 per 1000 tokens.

Table 35

The Distribution of Punctuation Errors in TR_COM Subcorpus

Punctuation	N	%	Per 1000 tokens
Confusion	11	1.13%	0.36
Redundant	39	4.02%	1.26
Missing	118	12.15%	3.83
Lexical	0	0.00%	0.00
Total	168	17.30%	5.45

4.3. Economics Subcorpora

4.3.1. Type/Token Ratio

Table 36 shows the type/token ratio of the Turkish Economics (TR_EAS) subcorpus and its counterpart (REF_EAS) in the reference corpus. TR_EAS subcorpus consists of 85,282 words, 9,158 types, and 82,853 tokens. Accordingly, the type-token ratio of the TR_EAS corpus was found to be 0.11. For the considerations explained in the method section, additional STTR and MATTR were also calculated. Both STTR and MATTR of the TR_COM corpus were calculated as 0.68. On the other hand, the REF_COM corpus consists of 80.147 words, 7,948 types, and 79,524 tokens. The type-token ratio of the reference corpus was calculated as 0.09, and both STTR and MATTR were 0.70. These findings suggest that Turkish scholars in the Economics subcorpus had a wide range of lexical diversity in their manuscripts, though their performance was below their counterparts in the reference corpus.

Table 36

Type/Token Ratio Results of Economics Subcorpora

	TR_EAS	REF_EAS
Total Words	85,282	80.147
Types	9,158	7948
Tokens	82,853	79524
TTR	0.11	0.09
STTR	0.68	0.70
MATTR	0.68	0.70

4.3.2. Syntactic Complexity

Table 37 shows the mean values of 14 syntactic complexity measures for Economics Corpora. It was found that the REF_EAS corpus had higher mean values 9 out of 14 syntactic complexity measures. In the length of production unit category, REF_EAS had higher means in MLS and MLC measures, while TR_EAS performed better in the MLC measure. Besides, in the REF_EAS subcorpus, the mean values of all measures related to the amount of subordination unit were higher than those in the TR_EAS subcorpus. On the other hand, TR_EAS had higher means in all of the amounts of coordination measures. In the Degree of phrasal sophistication unit, TR_EAS had a higher mean value in CN/C and lower mean values in CN/T and VP/T measures. Finally, the REF_EAS subcorpus outperformed the TR_EAS subcorpus in the overall sentence complexity measure. Therefore, it can be said that the REF_EAS subcorpus outperformed the TR_EAS subcorpus in all of the dimensions except for the amount of coordination. These findings show that participants in REF_EAS subcorpus employed relatively longer sentences, T-units, and clauses, and more subordinate phrases in their manuscripts than their Turkish counterparts.

Table 37

Means of Syntactic Complexity Measures for Economics Corpora

			N	Mean	Std. Deviation	
Length of the production unit	MLS	Study	14	25.18	4.44	
		Reference	12	28.16	4.70	
	MLT	Study	14	23.02	3.75	
		Reference	12	26.01	4.04	
	MLC	Study	14	15.04	1.95	
		Reference	12	14.35	1.34	
Amount of subordination	C/T	Study	14	1.53	0.12	
		Reference	12	1.78	0.16	
	CT/T	Study	14	0.37	0.64	
		Reference	12	0.50	0.66	
	DC/C	Study	14	0.31	0.39	
		Reference	12	0.42	0.50	
	DC/T	Study	14	0.49	0.96	
		Reference	12	0.77	0.15	
	Amount of coordination	CP/C	Study	14	0.50	0.15
			Reference	12	0.39	0.12
CP/T		Study	14	0.75	0.22	
		Reference	12	0.69	0.20	
T/S		Study	14	1.09	0.89	
		Reference	12	1.08	0.48	
Degree of phrasal sophistication	CN/C	Study	14	2.01	0.29	
		Reference	12	1.85	0.25	
	CN/T	Study	14	3.06	0.51	
		Reference	12	3.33	0.52	
	VP/T	Study	14	2.08	0.23	
		Reference	12	2.44	0.26	
Overall sentence complexity	C/S	Study	14	1.67	0.18	
		Reference	12	1.96	0.24	

A Mann-Whitney U test is run to determine whether there is a significant difference in the syntactic complexity units among the TR_EAS and REF_EAS data. We choose the non-parametric Mann-Whitney U test rather than a t-test because the sample sizes did not allow us to run a parametric independent samples t-test, although values of the syntactic measures in both Subcorpora were a normal distribution. In addition, The Bonferroni correction was employed to eliminate the family-wise error rate in the dataset. Therefore, the alpha value for each comparison was adjusted to $.05/14$, or $.0036$, where $.05$ is the significance level for the complete set of tests, and 14 is the number of tests performed. The Mann Whitney U test results are summarized in Table 38. The results reveal that the two corpora differed in 5 out of 14 measures. It was found that there were statistically significant differences in all

four measures of the amount of subordination and overall sentence complexity (measured using C/S index). No significant differences were found for other measures. These findings suggest that TR_EAS corpora produced shorter sentences, used fewer subordinates, showed less phrasal sophistication, and formed less complex sentences than their counterparts in the REF_EAS subcorpus. On the other hand, they used more coordination at both phrasal and sentential levels than participants in the REF_EAS subcorpus.

Table 38

Mann Whitney U Results of Economics Corpora

			N	Mean Rank	Sum of Ranks	U	P
Length of production unit	MLS	Study	14	11,36	159,00	54,000	.123
		Reference	12	16,00	192,00		
	MLT	Study	14	10,93	153,00	48,000	.064
		Reference	12	16,50	198,00		
	MLC	Study	14	14,86	208,00	65,000	.328
		Reference	12	11,92	143,00		
Amount of subordination	C/T	Study	14	8,86	124,00	19,000	.001*
		Reference	12	18,92	227,00		
	CT/T	Study	14	8,29	116,00	11,000	<0.001*
		Reference	12	19,58	235,00		
	DC/C	Study	14	7,86	110,00	5,000	<0.001*
		Reference	12	20,08	241,00		
	DC/T	Study	14	8,14	114,00	9,000	<0001*
		Reference	12	19,75	237,00		
Amount of coordination	CP/C	Study	14	15,86	222,00	51,000	.090
		Reference	12	10,75	129,00		
	CP/T	Study	14	14,57	204,00	69,000	.440
		Reference	12	12,25	147,00		
	T/S	Study	14	13,57	190,00	83,000	.959
		Reference	12	13,42	161,00		
Degree of phrasal sophistication	CN/C	Study	14	15,14	212,00	61,000	.237
		Reference	12	11,58	139,00		
	CN/T	Study	14	12,07	169,00	64,000	.304
		Reference	12	15,17	182,00		
	VP/T	Study	14	9,50	133,00	28,000	.004
		Reference	12	18,17	218,00		
Overall sentence complexity	C/S	Study	14	9,43	132,00	27,000	.003*
		Reference	12	18,25	219,00		

Note: * $p < .0036$; see Table 12 for the definitions.

4.3.3. Computer-Aided Error Analysis

A total of 1,513 errors were identified in this TR_EAS subcorpus (14.35 errors per 1,000 words). Based on the fact that there were 14 research articles in the TR_EAS subcorpus, each research article averaged 108 errors. Table 39 shows the number of errors in general categories. It is seen that grammatical errors represent 48.12% of all errors (6.90 errors per 1000 tokens). The number of the errors tagged in other categories was as follows: Punctuation (N=229), Lexis (N=138), Register (N=123), Lexicogrammar (N=122), Form (N=80), Word (N=70), and Style (N=23), respectively. Errors at each level will be further analyzed and examined in detail.

Table 39

The Computer-Aided Error Analysis Results of TR_EAS Subcorpus

	N	%	Per 1000 tokens
Form	80	5.29%	0.76
Grammar	728	48.12%	6.90
Lexis	138	9.12%	1.31
Register	123	8.12%	1.17
Style	23	1.52%	0.22
Word	70	4.63%	0.66
Lexicogrammar	122	8.06%	1.16
Punctuation	229	15.14%	2.17
Total	1,513	100%	14.35

Table 40 shows the distribution of Form errors in the TR_EAS subcorpus. It can be seen that all of the errors in this category belong to the spelling category. A total of 80 errors were tagged, and the normalized number of the errors were found to be 0.76 per 1000 words.

Table 40

The Distribution of Form Errors in TR_EAS Subcorpus

Form	N	f	Per 1000 tokens
Morphology	0	0.00%	0.00
Spelling	80	5.29%	0.76
Total	80	5.29%	0.76

A further investigation of Grammar errors, presented in Table 41, shows that the three most problematic grammatical features were articles (N=309, 2.93 per 1000 tokens), verbs (N=238, 2.26 per 1000 tokens), and nouns (N=124, 1.18 per 1000 toke), respectively. A close look at article errors reveals that the participants made the greatest number of errors in the missing category with 241 errors (2.28 per 1000 tokens). In addition, 57 redundant and

11 confusion errors were tagged in the Article category. In the adjectives category, 2 comparative-superlative and four order errors were detected. A total of 124 (20 cases and 104 number) errors were identified as Noun errors. The pronoun errors included five personal and one unclear reference error. The analysis of verb errors shows that the most problematic verb aspect was tense (N=131, 1.24 per 1000 tokens), followed by verb number (N=49), voice (N=34), auxiliaries (N=15), and nonfinite-finite (N=9), respectively. There were also 28-word-class errors and eight determiners (8 demonstratives and 4 indefinites) in this category.

Table 41

The Distribution of Grammar Errors in TR_EAS Subcorpus

Grammar	N	f	Per 1000 tokens
Articles	Redundant	57	3.77%
	Missing	241	15.93%
	Confusion	11	0.73%
	Total	309	20.42%
Adjectives	Comparative-Superlative	2	0.13%
	Order	4	0.26%
	Total	6	0.40%
Adverbs	5	0.33%	0.05
Nouns	Case	20	1.32%
	Number	104	6.87%
	Total	124	8.20%
Pronouns	Personal	5	0.33%
	Reflexive and reciprocal	0	0.00%
	Relative and interrogative	0	0.00%
	Unclear reference	1	0.07%
	Total	6	0.40%
Verbs	Auxiliaries	15	0.99%
	Verb-morphology	0	0.00%
	Verb-number	49	3.24%
	Nonfinite-finite	9	0.59%
	Tense	131	8.66%
	Voice	34	2.25%
	Total	238	15.73%
Word Class	28	1.85%	0.27
Determiner	Demonstrative	8	0.53%
	Possessive	0	0.00%
	Indefinite	4	0.26%
	Total	12	0.79%
Total	728	48.12%	6.90

Compilation and categorization of errors show that errors in the Lexis (N=138, 1.31 per 1000 tokens) category account for 9.12 % of the total number of errors in this subcorpus. As shown in Table 42, the greatest number of errors was tagged in the lexical single (N=112, 1.06 per 1000 tokens) subcategory, which was followed by conjunctions (N=15), lexical-phrase, and connectors (N=6), respectively. It is seen that the most problematic lexical feature is the lexical-single category. In other words, Turkish scholars in this subcorpus misused a single word 1.06 times per 1000 tokens.

Table 42

The Distribution of Lexis Errors in TR_EAS Subcorpus

Lexis	N	%	Per 1000 tokens
Conjunctions	Coordinating	11	0.73%
	Subordinating	4	0.26%
	Total	15	0.99%
Connectors	5	0.33%	0.05
Lexical-phrase	6	0.40%	0.06
Lexical-single	112	7.40%	1.06
Total	138	9.12%	1.31

Another error category is Register. Table 43 shows that a total of 123 errors (1.17 per 1000 tokens) were detected in this category. It is seen that the greatest number of errors were tagged in the vocabulary category (N=76). There were also 44 grammar and 3 phrase errors in the Register category.

Table 43

The Distribution of Register Errors in TR_EAS Subcorpus

Register	N	f	Per 1000 tokens
Vocabulary	76	5.02%	0.72
Grammar	44	2.91%	0.42
Phrase	3	0.20%	0.03
Total	123	8.13%	1.17

Table 44 shows the subcategories and the distribution of errors in the Style category. It is seen that most of the errors were made in the unclear (N=12) subcategory. In addition, 8 errors in the incomplete category and 3 errors in the other category were tagged. Overall, Turkish scholars in this subcorpus made 0.22 style errors per 1000 tokens.

Table 44

The Distribution of Style Errors in TR_EAS Subcorpus

Style	N	f	Per 1000 tokens
Other	3	0.20%	0.03
Incomplete	8	0.53%	0.08
Unclear	12	0.79%	0.11
Total	23	1.52%	0.22

Table 45 shows the distribution of Word errors in the TR_EAS subcorpus. It is seen that the most frequent error type was redundant (N=36, 0.34 per 1000 tokens). Besides, 18 order and 16 missing errors were detected in this category.

Table 45

The Distribution of Word Errors in TR_EAS Subcorpus

Word	N	f	Per 1000 tokens
Missing	16	1.06%	0.15
Order	18	1.19%	0.17
Redundant	36	2.38%	0.34
Total	70	4.63%	0.66

An in-depth examination of Lexicogrammar errors showed that the greatest number of the errors tagged in this category belongs to the prepositions subcategory (N=55, 0.53 per 1000 tokens). Furthermore, 36 verbs, 29 nouns, and 1 conjunction and 1 adjectives errors were identified in this category. The errors in this category equal to 8.06% of all of the errors in this subcorpus and 1.16 errors per 1000 tokens. In other words, lexicogrammatical errors were made 1.16 times in TR_EAS subcorpus.

Table 46

The Distribution of Lexicogrammar Errors in TR_EAS Subcorpus

Lexicogrammar	N	f	Per 1000 tokens	
Adjectives	Complementation	0	0.00%	0.00
	Dependent preposition	1	0.07%	0.01
	Total	1	0.07%	0.01
Conjunctions	1	0.07%	0.01	
Nouns	Complementation	16	1.06%	0.15
	Dependent preposition	13	0.86%	0.12
	Countable / Uncountable	0	0.00%	0.00
	Total	29	1.92%	0.27
Prepositions	55	3.64%	0.52	
Verbs	Complementation	12	0.79%	0.11
	Dependent preposition	24	1.59%	0.23
	Total	36	2.38%	0.34
Total	122	8.06%	1.16	

A further investigation of Punctuation errors, presented in Table 47, shows that the most problematic feature in this category was missing (N=113, 1.07 per 1000 tokens). In addition, 56 redundant, 47 confusion, and 13 lexical errors were tagged. The normalized value of this category shows that 2.17 punctuation errors were committed per 1000 tokens.

Table 47

The Distribution of Lexicogrammar Errors in TR_EAS Subcorpus

Punctuation	N	f	Per 1000 tokens
Confusion	47	3.11%	0.45
Redundant	56	3.70%	0.53
Missing	113	7.47%	1.07
Lexical	13	0.86%	0.12
Total	229	15.14%	2.17

4.4. Education Subcorpora

4.4.1. Type/Token Ratio

Table 48 shows the type/token ratio of the Turkish Education (TR_EDU) subcorpus and its counterpart (REF_EDU) in the reference corpus. TR_COM corpus consists of 209,973 words, 16,701 types and 199,780 tokens. The type-token ratio of the MCWP was found to be 0.08. However, for the considerations explained in the method section, additional STTR and MATTR were also calculated. Both STTR and MATTR of the TR_COM corpus were calculated as 0.67. On the other hand, the REF_COM corpus consists of 210,517 words,

14,609 types, and 209,289 tokens. The type-token ratio of the reference corpus was calculated as 0.07, and both STTR and MATTR were 0.71.

Table 48

Type/Token Ratio of Economics Subcorpora

	TR_EDU	REF_EDU
Total Words	209,973	210,517
Types	16,701	14609
Tokens	199,780	209289
TTR	0.08	0.07
STTR	0.67	0.71
MATTR	0.67	0.71

Based on the fact that the TTR value is quite sensible to text length and tends to decrease with the length of the texts (Brezina, 2018), SATTR and MATTR calculations were taken into consideration. These values show that Turkish scholars in Education subcorpus used an average range of vocabulary in their manuscripts. However, when compared to the reference corpus, the vocabulary they used is less varied.

4.4.2. Syntactic Complexity

Table 49 shows the mean values and t-test results of 14 syntactic complexity measures for Education Corpora. It was found that the REF_EDU subcorpus had higher mean values 12 out of 14 syntactic complexity measures. The TR_EDU subcorpus had higher means in one length of production unit measure (MLC) and one amount of coordination measure (CP/C). These findings show that the REF_EAS subcorpus outperformed the TR_EAS subcorpus in almost all of the syntactic complexity measures.

Table 49

Mean Values and T-test Results of Education Subcorpora

			N	Mean	Std. Deviation	t	p	Cohen's d
Length of the production unit	MLS	Study	49	26,31	4,05	-3,685	<.001*	0.81
		Reference	33	30,06	5,13			
	MLT	Study	49	24.52	3.67	-3,573	.001*	0.78
		Reference	33	27.84	4.71			
	MLC	Study	49	15.57	2.20	.985	.327	0.22
		Reference	33	15.07	2.24			
Amount of subordination	C/T	Study	49	1.59	0.22	-5,098	<.001*	1.18
		Reference	33	1.84	0.20			
	CT/T	Study	49	0.40	0.12	-5,487	<.001*	1.27
		Reference	33	0.53	0.08			
	DC/C	Study	49	0.34	0.09	-5,697	<.001*	1.30
		Reference	33	0.44	0.06			
	DC/T	Study	49	0.55	0.21	-5,444	<.001*	1.28
		Reference	33	0.78	0.14			
Amount of coordination	CP/C	Study	49	0.50	0.17	.413	.681	0.06
		Reference	33	0.49	0.14			
	CP/T	Study	49	0.78	0.24	-2,199	.031	0.51
		Reference	33	0.90	0.23			
	T/S	Study	49	1.07	0.05	1,308	.194	0.4.
		Reference	33	1.09	0.05			
Degree of phrasal sophistication	CN/C	Study	49	1.96	0.24	-.955	.342	0.19
		Reference	33	2.01	0.27			
	CN/T	Study	49	3.11	0.59	-4,868	<.001*	1.08
		Reference	33	3.75	0.59			
	VP/T	Study	49	2.17	0.31	-4,987	<.001*	1.13
		Reference	33	2.49	0.25			
Overall sentence complexity	C/S	Study	49	1.72	0.28	-4,942	<.001*	1.12
		Reference	33	2.03	0.27			

Note: * $p < .0036$; see Table 12 for the definitions.

A set of independent-samples t-tests were run to examine whether there are any statistically significant differences between TR_EDU subcorpus and REF_EDU subcorpus for each of the syntactic complexity measures since the data had a normal distribution. In addition, The Bonferroni correction was employed to eliminate the family-wise error rate in the dataset. Therefore, the alpha value for each comparison was adjusted to $.05/14$, or $.0036$, where $.05$ is the significance level for the complete set of tests, and 14 is the number of tests performed. The results suggest that TR_EDU and REF_EDU differed significantly in 9 out of 14 syntactic complexity measures. No statistically significant differences were found for one

length of production measure (MLC $t=.985$, $p=.327$) three measures of amount of coordination measures (CP/C $t=413$, $p=.681$, CPI/T $t=-2,199$, $p=0.31$ and T/S $t=1.308$, $p=.194$) and one degree of phrasal sophistication measure (CN/C $t=-.955$, $p=.342$). In addition, in order to investigate the effect size, Cohen's d value for each of the measures was calculated. Considering that a typical interpretation is to refer to effect sizes as small ($d=0.2$), medium ($d=0.5$), and large ($d=0.8$) (Cohen, 1988), 8 out of 10 measures that were found to be significantly different had a large effect size, and one (MLT $d=0.78$) had a medium effect. These findings demonstrate that the TR_EDU subcorpus used significantly less syntactically complex features than the REF_EDU corpus.

4.4.3. Computer-Aided Error Analysis

The results of computer-aided error analysis demonstrate that a total of 5,708 errors were annotated in the TR_EDU subcorpus. Considering that the TR_EDU corpus consisted of 49 manuscripts, an average of 116.48 errors were tagged in each text. In addition, the normalized value of the number of errors accounts for 20.50 per 1000 tokens. On the basis of the fact that the mean length of the sentence was calculated as 26,31, it can be said that 20.50 errors were detected in 38 sentences in the TR_EDU subcorpus. As can be seen in Table 50, more than half of the errors in this subcorpus were Grammar errors (N=2936, 10.54 per 1000 tokens). The greatest number of the errors was annotated in Register (N=626, 2.25 per 1000 tokens), Lexis (N=563, 2.02 per 1000 tokens), Punctuation (N=417, 1.50 per 1000 tokens), Lexicogrammar (N=387, 1.39 per 1000 tokens), Word (N=366, 1.31 per 1000 tokens), Style (N=222, 0.80 per 1000 tokens) and Form (N=191, 0.69 per 1000 tokens), respectively.

Table 50

The Computer-Aided Error Analysis Results of TR_EDU Subcorpus

	N	%	Per 1000 tokens
Form	191	3.35%	0.69
Grammar	2936	51.44%	10.54
Lexis	563	9.86%	2.02
Register	626	10.97%	2.25
Style	222	3.89%	0.80
Word	366	6.41%	1.31
Lexicogrammar	387	6.78%	1.39
Punctuation	417	7.31%	1.50
Total	5,708	100%	20.50

A detailed investigation of Form errors shows that almost all of the errors in this category were spelling errors. Table 51 demonstrates that only 1 error in morphology was tagged while 190 errors (0.68 per 1000 tokens) were annotated in the spelling category.

Table 51

The Distribution of Form Errors in TR_EDU Subcorpus

Form	N	f	Per 1000 tokens
Morphology	1	0.02%	0.01
Spelling	190	3.33%	0.68
Total	191	3.35%	0.69

Table 52 shows the subcategories and the distribution of errors in the Grammar category. It is seen that the greatest number of the was made in articles (N=1505, 5.40 per 1000 tokens), followed by verbs (N=824, 2.96 per 1000 tokens) and nouns (N=355, 1.27 per tokens). An in-depth investigation of article errors demonstrates that a great majority of the errors were detected in the missing category (N=1246, 4.47 per 1000 tokens) while there were 241 redundant and 18 confusion errors in this category. In adjectives category, 32 comparative-superlative (0.11 per 1000 tokens), 30 order (0.11 per 1000 tokens) were detected. In addition, a total of 39 errors were identified as in adverbs (0.14 per 1000 tokens). There were 306 numbers (1.10 per 1000 tokens), and 49 cases (0.18 per 1000 tokens) were identified in Nouns. The pronoun errors included 19 personal, 5 reflexive and reciprocal, 1 relative and interrogative, and 27 unclear reference errors. The analysis of verb errors shows that the most problematic verb aspect was tense (N=539, 1.94 per 1000 tokens), followed by verb number (N=115), voice (N=83), nonfinite-finite (N=51), auxiliaries (N=35) and verb morphology (N=1), respectively. There were also 67-word-class errors and 32 determiner errors in this category

Table 52

The Distribution of Grammar Errors in TR_EDU Subcorpus

Grammar	N	f	Per 1000 tokens	
Articles	Redundant	241	4.22%	
	Missing	1246	21.83%	
	Confusion	18	0.32%	
	Total	1505	26.37%	
Adjectives	Comparative-Superlative	32	0.55%	
	Order	30	0.53%	
	Total	62	1.09%	
Adverbs	39	0.68%	0.14	
Nouns	Case	49	0.86%	
	Number	306	5.36%	
	Total	355	6.22%	
Pronouns	Personal	19	0.33%	
	Reflexive and reciprocal	5	0.09%	
	Relative and interrogative	1	0.02%	
	Unclear reference	27	0.47%	
	Total	52	0.91%	
Verbs	Auxiliaries	35	0.61%	
	Verb-morphology	1	0.02%	
	Verb-number	115	2.01%	
	Nonfinite-finite	51	0.89%	
	Tense	539	9.44%	
	Voice	83	1.45%	
	Total	824	14.44%	
Word Class	67	1.17%	0.24	
Determiner	Demonstrative	15	0.26%	
	Possessive	7	0.12%	
	Indefinite	10	0.18%	
	Total	32	0.56%	
Total		2936	51.44%	10.54

An in-depth examination of Lexis errors showed that the greatest number of the errors in this category was tagged in the lexical-single subcategory (N=396, 1.42 per 1000 tokens). Furthermore, 57 connector errors, as well as 57 conjunctions errors (49 coordinating and 8 subordinating), were tagged. Besides, 52 lexical phrases were identified. Errors in the Lexis category represent 9.86% of all errors in this subcategory. The normalized value of the number of the errors accounts for 2.02 per 1000 tokens.

Table 53

The Distribution of Lexis Errors in TR_EDU Subcorpus

Lexis	N	f	Per 1000 tokens
Conjunctions	Coordinating	49	0.86%
	Subordinating	8	0.14%
	Total	57	1.00%
Connectors	57	1.00%	0.20
Lexical-phrase	53	0.93%	0.19
Lexical-single	396	6.94%	1.42
Total	563	9.86%	2.02

The register was the second most problematic category in TR_EDU Subcorpora. Table 54 shows that a total of 626 errors (2.25 per 1000 tokens) were annotated in this category. It is seen that the greatest number of errors was tagged in the vocabulary category (N=344, 1.24 per 1000 tokens). There were also 196 grammar and 86 phrase errors in the Register category.

Table 54

The Distribution of Lexis Errors in TR_EDU Subcorpus

Register	N	f	Per 1000 tokens
Vocabulary	344	6.03%	1.24
Grammar	196	3.43%	0.70
Phrase	86	1.51%	0.31
Total	626	10.97%	2.25

Table 55 shows the distribution of Style errors in the TR_EDU subcorpus. An annotation of errors shows that participants in this category had the most difficulty in the unclear category (N=161, 0.58 per 1000 tokens). In addition, 41 errors in incomplete and 20 in the other category were detected.

Table 55

The Distribution of Style Errors in TR_EDU Subcorpus

Style	N	f	Per 1000 tokens
Other	20	0.35%	0.07
Incomplete	41	0.72%	0.15
Unclear	161	2.82%	0.58
Total	222	3.89%	0.80

Table 56 shows the distribution of Word errors in the TR_EDU subcorpus. It is seen that the most frequent error type was redundant (N=151, 0.90 per 1000 tokens). Besides, 70 missing and 45 order errors were detected in this category.

Table 56

The Distribution of Style Errors in TR_EDU Subcorpus

Word	N	f	Per 1000 tokens
Missing	70	1.23%	0.25
Order	45	0.79%	0.16
Redundant	251	4.40%	0.90
Total	366	6.41%	1.31

A further investigation of Lexicogrammar errors, presented in Table 57, shows that the most problematic lexicogrammatical feature was verbs (N=151), followed by prepositions (N=136). A close look at verb errors reveals that the participants made the greatest number of errors in the dependent preposition category with 112 errors. In addition, there were 82 noun errors, 38 conjunctions errors, and 8 adjective errors in this category.

Table 57

The Distribution of Lexicogrammar Errors in TR_EDU Subcorpus

Lexicogrammar	N	f	Per 1000 tokens	
Adjectives	Complementation	2	0.04%	0.01
	Dependent preposition	6	0.11%	0.02
	Total	8	0.14%	0.03
Conjunctions	38	0.16%	0.04	
Nouns	Complementation	39	0.68%	0.14
	Dependent preposition	37	0.65%	0.13
	Countable / Uncountable	6	0.11%	0.02
	Total	82	1.44%	0.29
Prepositions	136	2.38%	0.49	
Verbs	Complementation	39	0.68%	0.14
	Dependent preposition	112	1.96%	0.40
	Total	151	2.65%	0.54
Total	387	6.78%	1.39	

The last category was Punctuation. Table 58 shows the distribution of errors in this category. The greatest number of the errors was annotated in the missing subcategory (N=238, 0.85 per 1000 tokens). Other punctuation errors included 103 redundant, 63 confusion, and 13 lexical errors.

Table 58

The Distribution of Punctuation Errors in TR_EDU Subcorpus

Punctuation	N	f	Per 1000 tokens
Confusion	63	1.10%	0.23
Redundant	103	1.80%	0.37
Missing	238	4.17%	0.85
Lexical	13	0.23%	0.05
Total	417	7.31%	1.50

4.5. Engineering Subcorpora**4.5.1. Type/Token Ratio**

Table 59 shows the type/token ratio of the Turkish Engineering (TR_ENG) subcorpus and its counterpart (REF_ENG) in the reference corpus. TR_ENG corpus consists of 378,439 words, 26,152 types, and 371,465 tokens. In that sense, the type-token ratio of the MCWP was found to be 0.07. For the considerations explained in the method section, additional STTR and MATTR were also calculated. Both STTR and MATTR of the TR_ENG corpus were calculated as 0.68. In contrast, REF_ENG subcorpus consists of 309,398 words, 22,883 types, and 376,873 tokens. The type-token ratio of REF_ENG was calculated as 0.06, and both STTR and MATTR were 0.68. These findings demonstrate that scholars in both subcorpora used a wide range of vocabulary in their manuscripts, and they performed nearly the same in lexical diversity.

Table 59

Type/Token Ratio of Engineering Subcorpora

	TR_ENG	REF_ENG
Total Words	378,439	309,398
Types	26,152	22883
Tokens	371,465	376873
TTR	0.07	0.06
STTR	0.68	0.68
MATTR	0.68	0.68

4.5.2. Syntactic Complexity

Table 60 shows the mean values and t-test results of 14 syntactic complexity measures for Education Corpora. It was found that REF_ENG subcorpus had higher mean values 12 out of 14 syntactic complexity measures. TR_ENG subcorpus had higher means in two amounts

of coordination measures, which are CP/C and CP/T. These results show that REF_ENG subcorpus outperformed TR_ENG subcorpus in all categories except for amount coordination.

Table 60

Mean Values and T-test Results of Engineering Corpora

			N	Mean	Std. Deviation	t	p	Cohen
Length of production unit	MLS	Study	88	22.02	4.91	-5,958	<.001*	0.96
		Reference	71	26.35	4.07			
	MLT	Study	88	20.84	3.95	-5,613	<.001*	0.90
		Reference	71	24.13	3.30			
	MLC	Study	88	14.57	2.36	-3,212	.002*	0.51
		Reference	71	15.72	2.06			
Amount of subordination	C/T	Study	88	1.43	0.14	-5,019	<.001*	0.77
		Reference	71	1.55	0.17			
	CT/T	Study	88	0.32	0.09	-4,664	<.001*	0.77
		Reference	71	0.39	0.09			
	DC/C	Study	88	0.27	0.05	-5,383	<.001*	1.08
		Reference	71	0.33	0.06			
	DC/T	Study	88	0.40	0.13	-5,226	<.001*	0.82
		Reference	71	0.52	0.16			
Amount of coordination	CP/C	Study	88	0.43	0.15	2,493	.014	0.44
		Reference	71	0.37	0.12			
	CP/T	Study	88	0.60	0.21	1,237	.218	0.20
		Reference	71	0.56	0.18			
	T/S	Study	88	1.05	0.06	-3,795	<.001*	0.46
		Reference	71	1.08	0.07			
Degree of phrasal sophistication	CN/C	Study	88	1.84	0.28	-4,602	<.001*	0.75
		Reference	71	2.07	0.33			
	CN/T	Study	88	2.62	0.45	-8,026	<.001*	1.29
		Reference	71	3.17	0.40			
	VP/T	Study	88	1.87	0.25	-6,059	<.001*	0.94
		Reference	71	2.12	0.28			
Overall sentence complexity	C/S	Study	88	1.50	0.17	-5,979	<.001*	0.99
		Reference	71	1.69	0.21			

Note: * $p < .0036$; see Table 12 for the definitions.

A set of independent-samples t-tests were to determine whether there is a significant difference in the syntactic complexity units among the TR_ENG and REF_ENG data since the data had a normal distribution. In addition, The Bonferroni correction was employed to eliminate the family-wise error rate in the dataset. Therefore, the alpha value for each comparison was adjusted to $.05/14$, or $.0036$, where $.05$ is the significance level for the

complete set of tests, and 14 is the number of tests performed. The results suggest that the Turkish corpus and the reference corpus differed significantly in 12 out of 14 syntactic complexity measures. No statistically significant differences were found for two amount of coordination measures (CP/T, $t=1,237$, $p=.218$; CP/C, $t=2,493$, $p=.014$). In addition, in order to investigate the effect size, Cohen's d value for each of the measures was calculated. Cohen's d results showed that 7 out of 13 measures that were found to be significantly different had a large effect size; four measures (MLC $d=0.51$, C/T $d=0.77$, CT/T $d=0.77$, and CN/C $d=0.75$) had a medium effect, and two measures (CP/C $d=0.44$ and CP/T $d=0.46$) had a small effect. These findings demonstrate that the TR_ENG subcorpus used significantly less syntactically complex features than the REF_ENG corpus except for three measures in the number of coordination measures.

4.5.3. Computer-Aided Error Analysis

Table 61 shows the results of the computer-aided error analysis. A total of 10,098 errors were annotated in the TR_ENG subcorpus. More than half of the errors in this subcorpus were annotated in the Grammar category ($N=5307$, 2.28 per 1000 tokens). Punctuation ($N=1134$, 1.50 per 1000 tokens), Word ($N=854$, 1.72 per 1000 tokens), Register ($N=814$, 1.64 per 1000 tokens), Lexis ($N=747$, 1.50 per 1000 tokens), Lexicogrammar ($N=659$, 1.32 per 1000 tokens), Form ($N=440$, 0.88 per 1000 tokens) and Style ($N=143$, 0.29 per 1000 tokens) was found to be the most problematic categories, respectively. Considering that the TR_EDU corpus consisted of 88 manuscripts, an average of 114.75 errors were tagged in each text. In addition, the normalized value of the number of errors was 20.29 per 1000 tokens. On the basis of the fact that the mean length of the sentence was calculated as 22.02 in this subcorpus, 20.29 errors were detected in each of the 45.41 sentences in the TR_ENG subcorpus.

Table 61

The Computer-Aided Error Analysis Results of TR_ENG Subcorpus

	N	%	Per 1000 tokens
Form	440	4.36%	0.88
Grammar	5307	52.55%	10.66
Lexis	747	7.40%	1.50
Register	814	8.06%	1.64
Style	143	1.42%	0.29
Word	854	8.46%	1.72
Lexicogrammar	659	6.53%	1.32
Punctuation	1134	11.23%	2.28
Total	10,098	100%	20.29

An investigation of Form errors shows that a great majority of the errors in this category were spelling errors. Table 62 demonstrates that 436 spelling (0.88 per 1000 tokens) and 4 morphology errors were annotated in TR_ENG subcorpus.

Table 62

The Distribution of Form Errors in TR_ENG Subcorpus

Form	N	f	Per 1000 tokens
Morphology	4	0.04%	0.01
Spelling	436	4.32%	0.88
Total	440	4.36%	0.88

A detailed examination of Grammar errors, presented in Table 63, shows that the three most problematic grammatical features were articles (N=2392, 4.81 per 1000 tokens), verbs (N=1900, 3.82 per 1000 tokens), and nouns (N=543, 1.09 per 1000 tokens), respectively. A close look at article errors shows that the greatest number of errors was tagged in the missing category with 1984 errors (3.99 per 1000 tokens). Furthermore, 328 redundant and 80 confusion errors were tagged in the Article category. In the adjectives category, 27 comparative-superlative and 43 order adjectives were detected. In addition, there were 56 Adverb errors in the TR_ENG subcorpus. A total of 543 (476 numbers and 67 cases) errors were identified in the Noun category. The pronoun errors included 31 reflexive and reciprocal, 26 personal, 25 unclear references, and 6 relative and interrogative errors. The analysis of verb errors shows that the most problematic verb aspect was tense (N=1196, 2.40 per 1000 tokens), followed by voice (N=325), verb number (N=227), nonfinite-finite (N=78) auxiliaries (N=51) and morphology (N=23), respectively. There were also 174-word-class

errors and 8 determiners (46 demonstratives, 24 indefinites, and 14 possessives) in this category.

Table 63

The Distribution of Grammar Errors in TR_ENG Subcorpus

Grammar		N	f	Per 1000 tokens
Articles	Redundant	328	3.25%	0.66
	Missing	1984	19.65%	3.99
	Confusion	80	0.79%	0.16
	Total	2392	23.69%	4.81
Adjectives	Comparative-superlative	27	0.27%	0.05
	Order	43	0.43%	0.09
	Total	70	0.69%	0.14
Adverbs		56	0.55%	0.11
Nouns	Case	67	0.66%	0.13
	Number	476	4.71%	0.96
	Total	543	5.38%	1.09
Pronouns	Personal	26	0.26%	0.05
	Reflexive and reciprocal	31	0.31%	0.06
	Relative and interrogative	6	0.06%	0.01
	Unclear reference	25	0.25%	0.05
	Total	88	0.87%	0.18
Verbs	Auxiliaries	51	0.51%	0.10
	Verb-morphology	23	0.23%	0.05
	Verb-number	227	2.25%	0.46
	Nonfinite-finite	78	0.77%	0.16
	Tense	1196	11.84%	2.40
	Voice	325	3.22%	0.65
	Total	1900	18.82%	3.82
Word Class		174	1.72%	0.35
Determiner	Demonstrative	46	0.46%	0.09
	Possessive	14	0.14%	0.03
	Indefinite	24	0.24%	0.05
	Total	84	0.83%	0.17
Total		5307	52.55%	10.66

Compilation and categorization of errors show that errors in the Lexis (N=747, 1.50 per 1000 tokens) category account for 7.40 % of the total number of errors in this subcorpus. As shown in Table 64, the greatest number of the errors tagged in this category was in the lexical-single subcategory (N=560, 1.13 per 1000 tokens). Furthermore, 95 conjunctions (75 coordinating and 20 subordinating), 58 lexical phrases, and 34 lexical phrase errors were identified, respectively.

Table 64

The Distribution of Lexis Errors in TR_ENG Subcorpus

Lexis	N	f	Per 1000 tokens
Conjunctions	Coordinating	75	0.74%
	Subordinating	20	0.20%
	Total	95	0.94%
Connectors	34	0.34%	0.07
Lexical-phrase	58	0.57%	0.12
Lexical-single	560	5.55%	1.13
Total	747	7.40%	1.50

The register was found to be the second most problematic category in the TR_ENG subcorpus. A total of 814 errors were identified in the Register category. Table 65 shows the subcategories and distribution of the errors in this category. Accordingly, 499 vocabulary (1.00 per 1000 tokens), 164 grammar (0.33 per 1000 tokens), and 151 order (0.30 per 1000 tokens) errors were tagged. It can be seen that there were 1.64 register errors per 1000 tokens.

Table 65

The Distribution of Register Errors in TR_ENG Subcorpus

Register	N	f	Per 1000 tokens
Vocabulary	499	4.94%	1.00
Grammar	164	1.62%	0.33
Phrase	151	1.50%	0.30
Total	814	8.06%	1.64

Table 66 shows the subcategories and the distribution of errors in the Style category. It is seen that most of the errors were made in the unclear (N=103) subcategory. In addition, 21 errors in the incomplete category and 19 errors in the other category were annotated. Overall, scholars in TR_ENG made 0.22 style errors per 1000 tokens.

Table 66

The Distribution of Style Errors in TR_ENG Subcorpus

Style	N	f	Per 1000 tokens
Other	19	0.19%	0.04
Incomplete	21	0.21%	0.04
Unclear	103	1.02%	0.21
Total	143	1.42%	0.29

Table 67 shows the subcategories and the distribution of errors in the Word category. A total of 854 errors were identified in this category. The number of errors in each subcategory was as follows: 526 redundant, 186 order, and 142 missing errors. It can be seen that there were 1.72-word errors per 1000 tokens.

Table 67

The Distribution of Word Errors in TR_ENG Subcorpus

Word	N	f	Per 1000 tokens
Missing	142	1.41%	0.29
Order	186	1.84%	0.37
Redundant	526	5.21%	1.06
Total	854	8.46%	1.72

An in-depth investigation of Lexicogrammar errors showed that the greatest number of the errors was tagged in the verbs subcategory (N=240, 0.48 per 1000 tokens). Furthermore, 215 prepositions, 164 nouns, and 22 adjectives, and 18 conjunctions errors were identified in this category. The errors in this category equal to 6.53% of all of the errors in this subcorpus and the normalized value was found to be 1.32 errors per 1000 tokens. In sum, lexicogrammatical errors were made 1.32 times in the TR_ENG subcorpus.

Table 68

The Distribution of Lexicogrammar Errors in TR_ENG Subcorpus

Lexicogrammar	N	f	Per 1000 tokens	
Complementation	3	0.03%	0.01	
Adjectives	Dependent preposition	19	0.19%	0.04
	Total	22	0.22%	0.04
Conjunctions	18	0.18%	0.04	
Complementation	99	0.98%	0.20	
Nouns	Dependent preposition	65	0.64%	0.13
	Countable / Uncountable	0	0.00%	0.00
	Total	164	1.62%	0.33
Prepositions	215	2.13%	0.43	
Complementation	113	1.12%	0.23	
Verbs	Dependent preposition	127	1.26%	0.26
	Total	240	2.38%	0.48
Total	659	6.53%	1.32	

Punctuation was the third category; the greatest number of errors was tagged in this subcorpus. An examination of the distribution of Punctuation errors shows that the most problematic feature in this category was missing (N=661, 1.33 per 1000 tokens). In addition,

331 redundant, 132 confusion, and 10 lexical errors were tagged. The normalized value of the errors in this category reveals that 2.28 punctuation errors were committed per 1000 tokens.

Table 69

The Distribution of Punctuation Errors in TR_ENG Subcorpus

Punctuation	N	f	Per 1000 tokens
Confusion	132	1.31%	0.27
Redundant	331	3.28%	0.66
Missing	661	6.55%	1.33
Lexical	10	0.10%	0.02
Total	1134	11.23%	2.28

4.6. Medicine Subcorpora

4.6.1. Type/Token Ratio

Table 70 shows the type/token ratio of the Turkish Medicine (TR_MED) subcorpus and its counterpart (REF_MED) in the reference corpus. TR_MED corpus consists of 190,121 words, 18,132 types and 171,839 tokens. In that sense, the type-token ratio of TR_MED subcorpus was found to be 0.10. For the considerations explained in the method section, additional STTR and MATTR were also calculated. Both STTR and MATTR of the TR_ENG corpus were calculated as 0.68. In contrast, the REF_MED subcorpus consists of 192,370 words, 18,044 types, and 191,198 tokens. The type-token ratio of REF_MED was calculated as 0.09, and both STTR and MATTR were 0.70. The findings demonstrate that the TTR values of both Medicine subcorpora were the same as those of the main corpora. In addition, it can be said that participants in these corpora used a wide range of vocabulary in their research articles. However, a comparison of TR_MED and REF_MED corpora reveals that Turkish scholars used less varied vocabulary than their counterparts.

Table 70

Type/Token Ratio of Medicine Subcorpora

	TR_MED	REF_MED
Total Words	190,121	192,370
Types	18,132	18,044
Tokens	171,839	191,198
TTR	0.10	0.09
STTR	0.68	0.70
MATTR	0.68	0.70

4.6.2. Syntactic Complexity

Table 71 shows the mean values and t-test results of 14 syntactic complexity measures for Medicine Corpora. It was found that the REF_MED subcorpus had higher mean values in all of the measures. In other words, the REF_MED subcorpus outperformed the TR_MED subcorpus in all of the syntactic complexity measures.

In order to investigate whether there were any statistically significant differences between TR_MED subcorpus and REF_MED subcorpus, a set of independent-samples t-tests were run to examine for each of the syntactic complexity measures since the data had a normal distribution. In addition, The Bonferroni correction was employed to eliminate the family-wise error rate in the dataset. Therefore, the alpha value for each comparison was adjusted to $.05/14$, or $.0036$, where $.05$ is the significance level for the complete set of tests, and 14 is the number of tests performed. The results show that the TR_MED corpus and the reference corpus differed significantly in 11 out of 14 syntactic complexity measures. No statistically significant differences were found for one amount of subordination measure (C/T $t=-2,939$, $p=0.004$) and two amount of coordination measures (CP/C $t=-1,585$, $p=.116$; CP/T $t=-2,401$, $p=0.018$). Furthermore, Cohen's d value for each of the measures was calculated in order to investigate the effect size. Based on the fact that that a typical interpretation of Cohen's d value is to refer to effect sizes as small ($d =0.2$), medium ($d=0.5$), and large ($d=0.8$) (Cohen, 1988), 9 out of 13 measures found to be significantly different had a large effect size, and 2 (DC/T $d=0.64$ and CT/T $d=0.73$) had a medium effect. These findings demonstrate that the TR_MED subcorpus used significantly less syntactically complex features than the REF_MED corpus.

Table 71

Mean Values and T-test Results of Medicine Subcorpora

		N	Mean	Std. Deviation	t	p	Cohen		
Length of the production unit	MLS	Study	54	21,52	4,43	-9,637	<.001*	1.95	
		Reference	42	31,33	5,54				
	MLT	Study	54	19.86	3.53	-9,596	<.001*	1.92	
		Reference	42	28.49	5.26				
	MLC	Study	54	13.83	1.87	-9,067	<.001*	1.81	
		Reference	42	18.24	2.88				
Amount of subordination	C/T	Study	54	1.42	0.18	-2,939	.004	0.62	
		Reference	42	1.53	0.17				
	CT/T	Study	54	0.31	0.10	-3,501	.001*	0.73	
		Reference	42	0.38	0.09				
	DC/C	Study	54	0.26	0.06	-3,797	<.001*	0.92	
		Reference	42	0.32	0.07				
	DC/T	Study	54	0.40	0.15	-3,381	.001*	0.64	
		Reference	42	0.50	0.16				
	Amount of coordination	CP/C	Study	54	0.45	0.13	-1,585	.116	0.34
			Reference	42	0.50	0.16			
CP/T		Study	54	0.65	0.20	-2,401	.018	0.48	
		Reference	42	0.76	0.25				
T/S		Study	54	1.05	0.06	-3,855	<.001*	0.90	
		Reference	42	1.10	0.05				
Degree of phrasal sophistication	CN/C	Study	54	1.87	0.29	-7,536	<.001*	1.53	
		Reference	42	2.38	0.37				
	CN/T	Study	54	2.69	0.52	-9,068	<.001*	1.88	
		Reference	42	3.65	0.50				
	VP/T	Study	54	1.79	0.28	-4,809	<.001*	1.00	
		Reference	42	2.07	0.28				
Overall sentence complexity	C/S	Study	54	1.51	0.22	-4,230	<.001*	0.86	
		Reference	42	1.71	0.24				

Note: * $p < .0036$; see for the definitions.

4.6.3. Computer-Aided Error Analysis

A total of 4,963 errors were identified (20.19 errors per 1,000 words) in the TR_MED subcorpus. Based on the fact that there were 54 research articles in this subcorpus, each research article averaged 91.9 errors. In addition, since the mean length of the sentence was calculated as 21,52, it can be said that there were 20.19 errors in 40.46 sentences. Table 72 shows the number of errors in general categories. It is seen that grammatical errors represent 40.40% of all errors (N= 2005, 8.15 errors per 1000 tokens). The number of the errors tagged in other categories was as follows: Punctuation (N=841), Register (N=495), Lexis (N=432),

Word (N=432), Form (N=371), Lexicogrammar (N=280) and Style (N=107), respectively. Errors at each level will be further analyzed and examined in detail.

Table 72

The Computer Aided Error Analysis Results of TR_MED Subcorpus

	N	%	per 1000 tokens
Form	371	7.48%	1.51
Grammar	2005	40.40%	8.15
Lexis	432	8.70%	1.76
Register	495	9.97%	2.01
Style	107	2.16%	0.44
Word	432	8.70%	1.76
Lexicogrammar	280	5.64%	1.14
Punctuation	841	16.95%	3.42
Total	4,963	100%	20.19

An investigation of Form errors demonstrates that a great majority of the errors in this category were spelling errors. Table 73 shows that 367 spelling (1.49 per 1000 tokens) and 4 morphology errors were annotated in TR_MED subcorpus.

Table 73

The Distribution of Form Errors in TR_MED Subcorpus

Form	N	f	Per 1000 tokens
Morphology	4	0.08%	0.02
Spelling	367	7.39%	1.49
Total	371	7.48%	1.51

Table 74 shows the distribution of Grammar errors in the TR_COM subcorpus. Categorization and annotation of errors show that errors in the Grammar category (N=2005) account for 40.40% of the total number of errors in the TR_MED subcorpus. It is seen that the greatest number of errors was made in articles (N=1048, 4.26 per 1000 tokens), followed by verbs (N=467, 1.90 per 1000 tokens) and nouns (N=246, 0.90 per tokens). An in-depth investigation of article errors demonstrates that a great majority of the errors were detected in the missing category (N=937, 3.81 per 1000 tokens) while there were 96 redundant and 15 confusion errors in this category. In the adjectives category, 14 comparative-superlative and 27 order errors were detected. In addition, 28 errors were identified in the adverbs subcategory. Furthermore, 204 numbers (0.83 per 1000 tokens) and 42 cases (0.17 per 1000 tokens) were identified in Nouns. The pronoun errors included 35 relative and interrogative, 28 personal, 27 unclear, and 5 reflexive and reciprocal reference errors. The analysis of verb

errors shows that the most problematic verb aspect was tense (N=209, 0.85 per 1000 tokens), followed by voice (N=106), verb number (N=93), nonfinite-finite (N=30), auxiliaries (N=25), and verb morphology (N=4), respectively. There were also 59 word-class errors and 54 determiner errors in this category.

Table 74

The Distribution of Grammar Errors in TR_MED Subcorpus

Grammar		N	f	Per 1000 tokens
Articles	Redundant	96	1.93%	0.39
	Missing	937	18.88%	3.81
	Confusion	15	0.30%	0.06
	Total	1048	21.12%	4.26
Adjectives	Comparative-superlative	14	0.28%	0.06
	Order	4	0.08%	0.02
	Total	18	0.36%	0.07
Adverbs		28	0.56%	0.11
Nouns	Case	42	0.85%	0.17
	Number	204	4.11%	0.83
	Total	246	4.96%	1.00
Pronouns	Personal	28	0.56%	0.11
	Reflexive and reciprocal	5	0.10%	0.02
	Relative and interrogative	35	0.71%	0.14
	Unclear reference	17	0.34%	0.07
	Total	85	1.71%	0.35
Verbs	Auxiliaries	25	0.50%	0.10
	Verb-morphology	4	0.08%	0.02
	Verb-number	93	1.87%	0.38
	Nonfinite-finite	30	0.60%	0.12
	Tense	209	4.21%	0.85
	Voice	106	2.14%	0.43
	Total	467	9.41%	1.90
Word Class		59	1.19%	0.24
Determiner	Demonstrative	18	0.36%	0.07
	Possessive	14	0.28%	0.06
	Indefinite	22	0.44%	0.09
	Total	54	1.09%	0.22
	Total	2005	40.40%	8.15

A detailed investigation of Lexis errors, presented in Table 75, shows that the greatest number of errors was made in lexical-single (N=103 1.04 per 1000 tokens). Furthermore, 121 conjunctions, 35 connectors, and 18 lexical phrases were identified in this category. The errors in this category equal to 13.49% of all of the errors in the TR_MED subcorpus and 1.76 errors per 1000 tokens.

Table 75

The Distribution of Lexis Errors in TR_MED Subcorpus

Lexis	N	f	Per 1000 tokens	
Conjunctions	Coordinating	92	1.85%	0.37
	Subordinating	29	0.58%	0.12
	Total	121	2.44%	0.49
Connectors	35	0.71%	0.14	
Lexical-phrase	19	0.38%	0.07	
Lexical-single	257	5.18%	1.04	
Total	432	8.70%	1.76	

The register was the third most problematic category in TR_EDU Subcorpora. Table 76 shows that a total of 495 errors (2.01 per 1000 tokens) were annotated in this category. It is seen that the greatest number of errors was tagged in the vocabulary category (N=257, 1.04 per 1000 tokens). There were also 130 grammar and 108 phrase errors in the Register category.

Table 76

The Distribution of Register Errors in TR_MED Subcorpus

Register	N	f	Per 1000 tokens
Vocabulary	257	5.18%	1.04
Grammar	130	2.62%	0.53
Phrase	108	2.18%	0.44
Total	495	9.97%	2.01

Table 77 shows the subcategories and the distribution of errors in the Style category. It is seen that most of the errors were made in the unclear (N=91) subcategory. In addition, 9 errors in the other category and 7 errors in the incomplete category were tagged. Overall, Turkish scholars in this subcorpus made 0.44 style errors per 1000 tokens.

Table 77

The Distribution of Style Errors in TR_MED Subcorpus

Style	N	f	Per 1000 tokens
Other	9	0.18%	0.04
Incomplete	7	0.14%	0.03
Unclear	91	1.83%	0.37
Total	107	2.16%	0.44

Table 78 shows the distribution of Word errors in the TR_MED subcorpus. It is seen that the most frequent error type was redundant (N=257, 1.04 per 1000 tokens). Besides, 91 missing and 84 order errors were detected in this category.

Table 78

The Distribution of Word Errors in TR_MED Subcorpus

Word	N	f	Per 1000 tokens
Missing	91	1.83%	0.37
Order	84	1.69%	0.34
Redundant	257	5.18%	1.04
Total	432	8.70%	1.76

Annotation and categorization of errors show that Lexicogrammar errors (N=280 1.14 per 1000 tokens) account for 5.64 % of the total number of errors in this subcorpus. As shown in Table 79, the greatest number of errors was tagged in the verbs (N=120) subcategory, which was followed by prepositions (N=105) and Nouns (N=35), respectively. In addition, 14 adjectives and 6 conjunctions errors were detected in this category. Examination of the normalized value of the errors demonstrates that Lexicogrammar errors were annotated 1.14 times per 1000 words.

Table 79

The Distribution of Lexicogrammar Errors in TR_MED Subcorpus

Lexicogrammar	N	f	Per 1000 tokens	
Complementation	3	0.06%	0.01	
Adjectives	Dependent preposition	11	0.22%	0.04
	Total	14	0.28%	0.06
Conjunctions	6	0.12%	0.02	
	Complementation	12	0.24%	0.05
Nouns	Dependent preposition	23	0.46%	0.09
	Countable / Uncountable	0	0.00%	0.00
	Total	35	0.71%	0.14
Prepositions	105	2.12%	0.43	
	Complementation	19	0.38%	0.08
Verbs	Dependent preposition	101	2.04%	0.41
	Total	120	2.42%	0.49
Total	280	5.64%	1.14	

Punctuation was the second most problematic category in the TR_MED subcorpus. Table 80 shows the subcategories and the distribution of errors in the Punctuation category. An investigation of the distribution of Punctuation errors reveals that the subcategory in which

the greatest number of errors was annotated in missing (N=472 1.92 per 1000 token), which equals to 9.51% of all errors in this category. In addition, 276 errors (1.06 per 1000 tokens) in the redundant subcategory, 106 errors in the confusion subcategory, and 2 errors in the lexical subcategory were tagged. The normalized occurrences of Punctuation errors were found to be 3.42 per 1000 tokens.

Table 80

The Distribution of Punctuation Errors in TR_MED Subcorpus

Punctuation	N	f	Per 1000 tokens
Confusion	106	2.14%	0.43
Redundant	261	5.26%	1.06
Missing	472	9.51%	1.92
Lexical	2	0.04%	0.01
Total	841	16.95%	3.42

4.7. Comparison of Turkish Subcorpora

4.7.1. Type/Token Ratio

Table 81 shows the total words and type/token ratios of Turkish subcorpora. Based on the fact that the TTR value decreases when the length of the texts increases, SATTR and MATTR values should be taken into consideration while comparing several corpora (Brezina, 2018). A comparison of these values shows that TR_COM subcorpus had the highest STTR (0.70) and MATTR (0.70) values, whereas TR_EDU subcorpus had the lowest values (STTR=0.67, MATTR=0.67). The other subcorpora had the same STTR and MATTR values (0.68). These findings demonstrate that the Turkish scholars used a wide variety of vocabulary in their manuscripts. However, scholars in the Communication discipline had slightly more lexical diversity than the scholars in other disciplines. On the other hand, scholars in the field of education exhibited slightly less lexical diversity when compared to other disciplines.

Table 81

Type/Token Ratios of Turkish Subcorpora

	TR_COM	TR_EAS	TR_EDU	TR_ENG	TR_MED
Total Words	25,833	85,282	209,973	378,439	190,121
Types	4,438	9,158	16,701	26,152	18,132
Tokens	25,706	82,853	199,780	371,465	171,839
TTR	0.17	0.11	0.08	0.07	0.10
STTR	0.70	0.68	0.67	0.68	0.68
MATTR	0.70	0.68	0.67	0.68	0.68

4.7.2. Syntactic Complexity

A one-way between-groups analysis of variance (ANOVA) was conducted to compare five subcorpora with regard to 14 syntactic complexity measures in 5 dimensions. In addition, a post hoc test was applied to understand the source of this difference where appropriate. Based on the fact that the sample size was not equal in the dataset, the Games-Howell posthoc test was used since equal variances and sample sizes are not required as assumptions in this test. Furthermore, the Bonferroni correction was employed to eliminate the family-wise error rate in the dataset. Therefore, the alpha value for each comparison was adjusted to $.05/14$, or $.0036$, where $.05$ is the significance level for the complete set of tests, and 14 is the number of tests performed.

Table 82 gives information about the means and standard deviations of the three-length of unit measures as well as the significant differences and effect sizes. The results demonstrate that the TR_EDU subcorpus had the highest mean values in all three measures, whereas TR_MED had the lowest values.

Table 82

ANOVA Results of Length of Production Unit Category

	Group	<i>M</i>	<i>SD</i>	<i>F</i>	<i>p</i>	η^2	Post Hoc	
Length of the production unit	MLS	TR_COM	25,58	6,26	9.873	<.001*	.160	3>4 3>5
		TR_EAS	25,18	4,44				
		TR_EDU	26,31	4,05				
		TR_ENG	22,02	4,91				
		TR_MED	21,52	4,43				
	MLT	TR_COM	22,75	4,32	11.709	<.001*	.185	3>4 3>5
		TR_EAS	23,02	3,75				
		TR_EDU	24,52	3,67				
		TR_ENG	20,84	3,95				
		TR_MED	19,86	3,53				
	MLC	TR_COM	13,15	3,26	4.833	.001*	.085	3>5
		TR_EAS	15,04	1,95				
		TR_EDU	15,57	2,20				
		TR_ENG	14,57	2,36				
		TR_MED	13,83	1,87				

Note: * $p < .0036$; see Table 12 for the definitions; Criteria: 1:TR_COM, 2:TR_EAS, 3:TR_EDU, 4:TR_ENG, 5:TR_MED

ANOVA results suggest that Turkish subcorpora differed significantly in all of the indices in the length of production unit measures. In order to understand the source of this difference, the Games Howell post hoc test was performed. The results showed that both MLS and MLT values were found to be statistically higher in TR_EDU subcorpus than TR_ENG and TR_MED subcorpora, and MLC values were statistically higher in TR_EDU subcorpus than TR_MED subcorpus ($p < .0036$).

Table 83 summarizes the means and standard deviations of the number of subordination indices as well as the significant differences and effect sizes. It is seen that the TR_COM subcorpus had the highest mean scores in all four amounts of subordination. On the other hand, TR_MED subcorpus had the lowest means in three measures, and TR_MED and TR_ENG had the lowest means in one measure (DC/T). In order to understand the source of this difference, the Games Howell post hoc test was performed. The results demonstrated that TR_COM and TR_EDU subcorpora produced significantly more subordination than TR_ENG and TR_MED subcorpora

Table 83

ANOVA Results of Amount of Subordination Category

	Group	<i>M</i>	<i>SD</i>	<i>F</i>	<i>p</i>	η^2	Post Hoc	
Amount of subordination	C/T	TR_COM	1.63	0.14	9.648	<.001*	.157	1>4
		TR_EAS	1.53	0.12				1>5
		TR_EDU	1.59	0.22				3>4
		TR_ENG	1.43	0.14				3>5
		TR_MED	1.42	0.18				
	CT/T	TR_COM	0.43	0.06	7.901	<.001*	.132	1>4
		TR_EAS	0.37	0.06				1>5
		TR_EDU	0.40	0.12				2>5
		TR_ENG	0.32	0.09				3>4
		TR_MED	0.31	0.10				3>5
	DC/C	TR_COM	0.39	0.08	13.078	<.001*	.202	1>4
		TR_EAS	0.31	0.04				1>5
		TR_EDU	0.34	0.09				2>4
		TR_ENG	0.27	0.05				2>5
		TR_MED	0.26	0.06				3>4
	DC/T	TR_COM	0.61	0.14	11.074	<.001*	.176	1>5
		TR_EAS	0.49	0.10				3>4
		TR_EDU	0.55	0.21				3>5
		TR_ENG	0.40	0.13				
		TR_MED	0.40	0.15				

Note: * $p < .0036$; see Table 12 for the definitions; Criteria: 1:TR_COM, 2:TR_EAS, 3:TR_EDU, 4:TR_ENG, 5:TR_MED

Table 84 shows the means and standard deviations of the number of coordination measures as well as the significant differences and effect sizes. The results show that TR_EDU and TR_EAS had the highest mean in CP/C, TR_EDU had the highest mean in CP/T, and TR_EAS had the highest mean in T/S measure. One-way ANOVA test showed that statistically significant differences were found only in CP/T measure. Games Howell posthoc test revealed that it was statistically higher in TR_EDU than TR_ENG and TR_MED.

Table 84

ANOVA Results of Amount of Coordination Category

	Group	<i>M</i>	<i>SD</i>	<i>F</i>	<i>p</i>	η^2	Post Hoc
Amount of coordination	TR_COM	0.36	0.10				
	TR_EAS	0.50	0.15				
	CP/C TR_EDU	0.50	0.17	2.882	.024	.053	
	TR_ENG	0.43	0.15				
	TR_MED	0.45	0.13				
	TR_COM	0.63	0.12				
	TR_EAS	0.76	0.22				
	CP/T TR_EDU	0.78	0.24	5.887	<.001*	.102	3>4 3>5
	TR_ENG	0.60	0.21				
	TR_MED	0.65	0.20				
	TR_COM	0.96	0.37				
	TR_EAS	1.09	0.09				
	T/S TR_EDU	1.07	0.05	3.241	.013	.059	
	TR_ENG	1.05	0.06				
	TR_MED	1.05	0.06				

Note: * $p < .0036$; see Table 12 for the definitions; Criteria: 1:TR_COM, 2:TR_EAS, 3:TR_EDU, 4:TR_ENG, 5:TR_MED

Table 85 gives information about the means and standard deviations of the three degrees of phrasal sophistication measures as well as the significant differences and effect sizes. The results demonstrate that TR_EAS had the highest mean value in CN/C, TR_COM had the highest value in CN/T, and TR_EDU had the highest mean in VP/T measure. ANOVA results suggest that Turkish subcorpora differed significantly in two of the measures (CN/T $p < .001^*$; VP/T $p < .001^*$) of the degree of phrasal sophistication. In order to understand the source of this difference, the Games Howell post hoc test was performed. The results revealed that TR_EDU subcorpus produced significantly more complex nominal per T-unit (CP/T) than TR_ENG and TR_MED subcorpora. Furthermore, it was also found that TR_EAS and TR_EDU produced significantly more verb phrases per T-unit than TR_ENG and TR_MED subcorpora. Nevertheless, no significant difference was found for complex nominals per clause (CN/C) measure.

Table 85

ANOVA Results of Degree of Phrasal Sophistication Category

	Group	<i>M</i>	<i>SD</i>	<i>F</i>	<i>p</i>	η^2	Post Hoc	
Degree of phrasal sophistication	CN/C	TR_COM	1.87	0.70	1.718	.147	.032	
		TR_EAS	2.01	0.29				
		TR_EDU	1.96	0.24				
		TR_ENG	1.84	0.28				
		TR_MED	1.87	0.29				
	CN/T	TR_COM	3.19	0.90	9.224	<.001*	.151	3>4
		TR_EAS	3.06	0.51				3>5
		TR_EDU	3.11	0.59				
		TR_ENG	2.62	0.45				
		TR_MED	2.69	0.52				
	VP/T	TR_COM	2.09	0.26	15.556	<.001*	.231	2>4
		TR_EAS	2.08	0.23				2>5
		TR_EDU	2.17	0.31				3>4
		TR_ENG	1.87	0.25				3>5
		TR_MED	1.79	0.28				

Note: * $p < .0036$; see Table 12 for the definitions; Criteria: 1:TR_COM, 2:TR_EAS, 3:TR_EDU, 4:TR_ENG, 5:TR_MED

Table 86 summarizes the means and standard deviations of overall sentence complexity as well as the significant differences and effect sizes. The results demonstrate that TR_COM and TR_MED had the highest mean, whereas TR_ENG had the lowest mean. The results of ANOVA demonstrated that there were statistically significant differences in overall sentence complexity. In order to understand the source of this difference, the Games Howell post hoc test was performed. It was found that TR_EAS and TR_EDU showed significantly more syntactic complexity than TR_ENG. In addition, the TR_MED subcorpus showed significantly more syntactic complexity than the TR_EDU subcorpus.

Table 86

ANOVA Results of Overall Sentence Complexity Category

	Group	<i>M</i>	<i>SD</i>	<i>F</i>	<i>p</i>	η^2	Post Hoc	
Overall sentence complexity	C/S	TR_COM	1.95	0.33	14.212	<.001*	.215	
		TR_EAS	1.67	0.18				2>4
		TR_EDU	1.72	0.28				3>4
		TR_ENG	1.50	0.17				5>3
		TR_MED	1.95	0.33				

Note: * $p < .0036$; see Table 12 for the definitions; Criteria: 1:TR_COM, 2:TR_EAS, 3:TR_EDU, 4:TR_ENG, 5:TR_MED

4.7.2. Computer-Aided Error Analysis

Since the subcorpora in this corpus were of different sizes, the use of normalized frequency was necessary in order to accurately compare five subcorpora. Therefore, the occurrences of errors were normalized to a value of per 1000 tokens. Table 87 shows the comparison of computer-aided error analysis results. The results demonstrate that the greatest number of errors were detected in the TR_COM subcorpus, while the fewest number of errors were annotated in the TR_EAS subcorpus. It can be seen that the most problematic category was grammar in each subcorpus, whereas the least problematic category was formed. Errors at each level will be further analyzed and examined in detail.

Table 87

The Comparison of Computer-Aided Error Analysis Results

	TR_COM	TR_EAS	TR_EDU	TR_ENG	TR_MED
Form	0.97	0.76	0.69	0.88	1.51
Grammar	11.51	6.90	10.54	10.66	8.15
Lexis	4.25	1.31	2.02	1.50	1.76
Register	3.08	1.17	2.25	1.64	2.01
Style	0.94	0.22	0.80	0.29	0.44
Word	2.30	0.66	1.31	1.72	1.76
Lexicogrammar	2.98	1.16	1.39	1.32	1.14
Punctuation	5.45	2.17	1.50	2.28	3.42
Total	31.48	14.35	20.50	20.29	20.19

Table 88 shows the normalized value of form errors in each corpus. The results show that spelling was the most problematic aspect of Form. In addition, TR_MED subcorpus made the greatest number of errors (1.51 per 1000 tokens), while TR_EDU made the fewest number of errors (0.69 per 1000 tokens) in this category.

Table 88

The Comparison of Form Errors

Form	TR_COM	TR_EAS	TR_EDU	TR_ENG	TR_MED
Morphology	0.00	0.00	0.00	0.01	0.02
Spelling	0.97	0.76	0.68	0.88	1.49
Total	0.97	0.76	0.69	0.88	1.51

Table 89 shows the subcategories and the normalized occurrences of grammar errors in each corpus. The results show that TR_COM made the greatest number of errors (11.51 per 1000 tokens) whereas TR_EAS made the fewest number of errors in the grammar category (6.90 per 1000 tokens). This situation was found to be the same in articles and determiners

subcategories. A detailed investigation of adjective errors demonstrates that TR_EDU subcorpus made the greatest number of errors (0.14 per 1000 tokens), and TR_COM made the fewest number of errors. In addition, TR_COM made noun errors 1.85 times per 1000 tokens, while TR_MED made them 1.00 times per 1000 tokens. For the pronouns, TR_MED subcorpus were found to be the subcorpus in which the greatest number of pronoun errors (0.35 per 1000 tokens) were tagged while TR_EAS made the fewest number of the errors (0.06 per 1000 tokens). A further investigation of verb errors shows that TR_ENG subcorpus made the greatest number of errors (3.82 per 1000 tokens), and TR_MED subcorpus made verb errors 1.90 times per 1000 tokens. This was the same for word-class errors.

Table 89

The Comparison of Grammar Errors

Grammar		TR_COM	TR_EAS	TR_EDU	TR_ENG	TR_MED
Articles	Redundant	0.81	0.54	0.87	0.66	0.39
	Missing	4.67	2.28	4.47	3.99	3.81
	Confusion	0.06	0.10	0.06	0.16	0.06
	Total	5.55	2.93	5.40	4.81	4.26
Adjectives	Comparative-superlative	0.09	0.02	0.11	0.05	0.06
	Order	0.03	0.04	0.11	0.09	0.02
	Total	0.13	0.06	0.22	0.14	0.07
Adverbs		0.03	0.05	0.14	0.11	0.11
Nouns	Case	0.49	0.19	0.18	0.13	0.17
	Number	1.36	0.99	1.10	0.96	0.83
	Total	1.85	1.18	1.27	1.09	1.00
Pronouns	Personal	0.13	0.05	0.07	0.05	0.11
	Reflexive and reciprocal	0.06	0.00	0.02	0.06	0.02
	Relative and interrogative	0.00	0.00	0.00	0.01	0.14
	Unclear reference	0.00	0.01	0.10	0.05	0.07
	Total	0.19	0.06	0.19	0.18	0.35
Verbs	Auxiliaries	0.13	0.14	0.13	0.10	0.10
	Verb-morphology	0.00	0.00	0.00	0.05	0.02
	Verb-number	0.88	0.46	0.41	0.46	0.38
	Nonfinite-finite	0.16	0.09	0.18	0.16	0.12
	Tense	1.20	1.24	1.94	2.40	0.85
	voice	0.68	0.32	0.30	0.65	0.43
	total	3.05	2.26	2.96	3.82	1.90
Word Class		0.29	0.27	0.24	0.35	0.24
Determiner	Demonstrative	0.16	0.08	0.05	0.09	0.07
	Possessive	0.00	0.00	0.03	0.03	0.06
	Indefinite	0.26	0.04	0.04	0.05	0.09
	Total	0.42	0.11	0.11	0.17	0.22
	Total	11.51	6.90	10.54	10.66	8.15

Table 90 shows the normalized values of lexis errors in each corpus. The results show that the greatest number of errors was detected in the TR_COM subcorpus (4.25 per 1000 tokens) while the fewest number of errors was identified in the TR_EAS subcorpus (1.31 per 1000 tokens). This situation was the same for all subcategories except for the lexical-single category, in which the fewest number of errors were tagged in TR_MED subcorpus (1.04 per 1000 tokens).

Table 90

The Comparison of Lexis Errors

Lexis	TR_COM	TR_EAS	TR_EDU	TR_ENG	TR_MED
	0.55	0.10	0.18	0.15	0.37
Conjunctions	0.10	0.04	0.03	0.04	0.12
	0.65	0.14	0.20	0.19	0.49
Connectors	0.13	0.05	0.20	0.07	0.14
Lexical-phrase	0.13	0.06	0.19	0.12	0.07
Lexical-single	3.34	1.06	1.42	1.13	1.04
Total	4.25	1.31	2.02	1.50	1.76

Table 91 shows the subcategories and the normalized occurrences of register errors in each corpus. The results show TR_COM subcorpus made the most number of errors (3.08 per 1000 tokens), and TR_EAS subcorpus made the least number of errors (1.17 per 1000 tokens). This situation was the same for all subcategories except for grammar, in which the fewest number of errors were annotated in TR_MED subcorpus.

Table 91

The Comparison of Register Errors

Register	TR_COM	TR_EAS	TR_EDU	TR_ENG	TR_MED
Vocabulary	1.49	0.72	1.24	1.00	1.04
Grammar	0.97	0.42	0.70	0.33	0.53
Phrase	0.62	0.03	0.31	0.30	0.44
Total	3.08	1.17	2.25	1.64	2.01

Table 92 shows the comparison of style errors. The results reveal that the subcorpora that experienced the most difficulty in the Style category was TR_COM (0.94 errors per 1000 tokens) while the TR_EAS subcorpus made the fewest number of errors (0.22 per 1000 tokens). This situation was also valid for the unclear subcategory. However, in the Other subcategory, TR_EDU made the greatest number of errors, and TR_COM and TR_EAS made the fewest number of errors. Furthermore, the TR_MED subcorpus was found to be

the subcorpus in which the fewest number of errors were detected in the incomplete subcategory.

Table 92

The Comparison of Style Errors

Style	TR_COM	TR_EAS	TR_EDU	TR_ENG	TR_MED
Other	0.03	0.03	0.07	0.04	0.04
Incomplete	0.13	0.08	0.15	0.04	0.03
Unclear	0.78	0.11	0.58	0.21	0.37
Total	0.94	0.22	0.80	0.29	0.44

Table 93 shows the subcategories and the normalized occurrences of grammar errors in each corpus. The results show that, in line with other categories, the greatest number of errors was annotated in the TR_COM subcorpus (2.30 per 1000 tokens) and the fewest number in the TR_EAS subcorpus (0.66 per 1000 tokens). This finding was the same for all of the subcategories in Word errors.

Table 93

The Comparison of Word Errors

Word	TR_COM	TR_EAS	TR_EDU	TR_ENG	TR_MED
Missing	0.65	0.15	0.25	0.29	0.37
Order	0.49	0.17	0.16	0.37	0.34
Redundant	1.17	0.34	0.90	1.06	1.04
Total	2.30	0.66	1.31	1.72	1.76

Table 94 shows the subcategories and the normalized occurrences of grammar errors in each corpus. The results show that TR_COM made the greatest number of errors (2.98 per 1000 tokens) though the fewest number of Lexicogrammar errors were annotated in TR_MED subcorpus (1.14 per 1000 tokens). This finding is similar in nouns and prepositions subcategories. A detailed investigation of adjectives errors reveals that the greatest number of the errors were detected in the TR_MED subcorpus and the fewest number in TR_EAS. The normalized occurrences of conjunctions errors show that TR_COM (0.10 per 1000 tokens) was the subcorpus in which the greatest number of errors was tagged while the fewest number of errors was detected in TR_EAS (0.01 per 1000 tokens) subcorpus. Finally, the investigation of verb errors showed that participants in TR_COM (1.62 per 1000 tokens) made the greatest number of errors though TR_EAS (0.34 per 1000 tokens) committed the least number of errors in this subcategory.

Table 94

The Comparison of Lexicogrammar Errors

Lexicogrammar		TR_COM	TR_EAS	TR_EDU	TR_ENG	TR_MED
Adjectives	Complementation	0.03	0.00	0.01	0.01	0.01
	Dependent preposition	0.00	0.01	0.02	0.04	0.04
	Total	0.03	0.01	0.03	0.04	0.06
Conjunctions		0.10	0.01	0.04	0.04	0.02
Nouns	Complementation	0.39	0.15	0.14	0.20	0.05
	Dependent preposition	0.19	0.12	0.13	0.13	0.09
	Countable / Uncountable	0.00	0.00	0.02	0.00	0.00
	Total	0.58	0.27	0.29	0.33	0.14
Prepositions		0.65	0.52	0.49	0.43	0.43
Verbs	Complementation	0.78	0.11	0.14	0.23	0.08
	Dependent preposition	0.84	0.23	0.40	0.26	0.41
	Total	1.62	0.34	0.54	0.48	0.49
	Total	2.98	1.16	1.39	1.32	1.14

Table 95 shows the normalized values of punctuation errors in each subcorpora. The results show that the greatest number of errors was annotated in the TR_COM subcorpora, while the fewest number of errors was detected in the TR_EDU subcorpus. This finding is the same in all subcategories except for the lexical subcategory in which the greatest number of the errors was made by TR_EAS, and the fewest number of errors was identified in the TR_COM subcorpus.

Table 95

The Comparison of Punctuation Errors

Punctuation	TR_COM	TR_EAS	TR_EDU	TR_ENG	TR_MED
Confusion	0.36	0.45	0.23	0.27	0.43
Redundant	1.26	0.53	0.37	0.66	1.06
Missing	3.83	1.07	0.85	1.33	1.92
Lexical	0.00	0.12	0.05	0.02	0.01
Total	5.45	2.17	1.50	2.28	3.42

CHAPTER V

DISCUSSION

This chapter focuses on the interpretation of the findings with respect to the research questions put forward in this study. First, the lexical variation of Turkish scholars, both in general and discipline-specific, is discussed using the results of type-token ratio. Second, syntactic complexity results are discussed. Third, computer-aided error analysis results are described and contrasted.

5.1. Lexical Diversity

The aim of this research was to investigate the linguistic problems faced by Turkish scholars in writing for publication. In this context, three linguistic dimensions were examined: lexical diversity, syntactic complexity, and errors they made in their manuscripts. Lexical diversity is a useful indicator to determine whether the writers have a high degree of proficiency in using a wide range of vocabulary (Laufer, 1991). The present study used type-token ratio (TTR) as the measure of lexical diversity since it is one of the most frequently used measures in language research (Richards, 1987). TTR is the proportion of the number of diverse words to the total number of words and indicates the lexical richness of a text (Nation, 2013). However, TTR is quite sensitive to text length; that is, it tends to decrease when the length of the text increases. Therefore, two additional measures of TTR were employed to compare the results: standardized type/token ratio (STRR) and moving average type/ token ratio (MATTR). The results of the TTR range between 0 and 1. A greater TTR indicates a larger lexical diversity in that when it approximates to 1; it conveys that every one of the words is used once in a text. On the other hand, when it approximates to 0, it indicates a lower TTR, pointing out a lesser lexical diversity in the text since this shows a frequent repetition of the words in the text (Brezina, 2018; Thomas, 2005).

Type/token ratio results of both the MCWP (TTR=0.06, STTR=0.68, MATTR=0.68) and the reference corpus (TTR=0.052, STTR=0.70, MATTR=0.70) showed that scholars in both corpora used a wide variety of vocabulary in their manuscripts. Considering that the participants in the corpus were advanced learners or native speakers of English, it can be said that the use of varied vocabulary is expected. In addition, since academic writing is characterized by a formal tone of writing, scholars are likely to avoid using the same words since one of the ways of achieving formality on a lexical basis is to use a more varied vocabulary (Bui, 2018). These findings are consistent with other studies in the literature. For example, Biber and Gray (2013) found that participants having higher TOEFL IBT scores showed a wide variety of lexical diversity. Similarly, Kang and Wang (2014) proved that participants who obtained high C1-C2 level Cambridge English Exam scores used a more varied vocabulary.

On the other hand, it seems that the MCWP used slightly less varied vocabulary than the reference corpus. The situation is similar in subcorpora as well. In all subcorpora except for engineering, participants in the reference subcorpora had slightly higher TTR values. However, in the Engineering subcorpora, the participants had the same values. This may be attributed to the fact that the texts in the MCWP were unedited manuscripts while the texts in the reference corpus were published in leading high-ranked international journals, which means that each text in the reference corpus had undergone several processes such as editorial and peer-review and proof-reading (Lee & Chen, 2009).

The comparison of Turkish corpora among themselves showed that TR_COM (TTR=0.17, STTR=0.70, MATTR=0.70) subcorpus had the highest TTR values while TR_EDU subcorpus showed the lowest values (TTR=0.08, STTR=0.67, MATTR=0.67). This finding indicates that the participants in the communication discipline demonstrated more lexical diversity in their manuscripts compared to other disciplines. However, these findings should be interpreted carefully since scholars in different disciplines are required to follow their own research patterns and use their own specialized vocabulary, which may differ across the disciplines as the nature of science requires (Chen, Deng, Zhong & Zhang, 2020). In addition, there are certain differences in the academic expression and style among the disciplines (Swales & Feak, 2011). Therefore, these requirements may demand them to use some terms recurrently, which then paves the way for lower TTR values.

5.2. Syntactic Complexity

The second linguistic dimension examined as a part of the present study was syntactic complexity. In second language acquisition writing research, syntactic complexity is regarded as a typical indicator of the linguistic performance and proficiency of the writers (Norris & Ortega, 2009). It is hypothesized that when they develop their skills and improve their abilities to control the language, they produce more syntactically complex sentences (Larsen-Freeman, 1978). Syntactic complexity was measured automatically using 5 broad categories: length of the production unit, amount of subordination, amount of coordination, degree of phrasal sophistication, and overall sentence complexity.

The length of the production unit was measured using three indices: MLC, MLT, and MLS. The results revealed that the reference corpus had higher means in all of the measures than the MCWP. The situation is similar in all of the subcorpora as well. This finding shows that native writers produced longer sentences, clauses, and T-units than Turkish scholars. Statistically significant differences were found between the MCWP and the reference corpus for all of the measures as well as between the subcorpora except for all measures in Communication and Economics and for MLC in Education.

The findings are consistent with those of Ai and Lu (2013), who reported that native students outperformed nonnative students in all of the length of production unit measures. However, they run contrary to Wu et al.'s (2020) study in which they compared the unedited research articles of ELF scholars and published articles of native speakers of English. They collected a corpus of ELF papers, including 789,671 words, and compared it with the research papers in the COCA corpus. They reported significant differences for MLS and MLC in favor of the ELF group, indicating that ELF scholars produced longer clauses and sentences. Furthermore, Lu and Ai (2015) compared the second language writers with 8 different L1 backgrounds with the native speakers of English. They first treated the second language writers as one group and reported significant differences for MLC measure in favor of the native group. Then, they compared each nonnative group with the native group and reported that native corpus failed to use longer sentences, clauses, and T-units than all nonnative groups. In that sense, they reported that Japanese and Chinese groups produced significantly shorter sentences; Japanese, Tswana, and Russian groups employed significantly shorter clauses, and Japanese, Russian and Chinese groups used significantly shorter T-units than the native group. The findings of Lu and Ai (2015) indicates the effect of the first language on syntactic complexity. In that vein, Uysal (2008) clearly demonstrated the effect of Turkish

on L2 English writing, concluding that Turkish writers share a number of similarities in L2 English writing with Japanese and Chinese learners. In that sense, together with the findings in Lu and Ai (2015), the finding of the present study supports her claim.

The studies in the literature have demonstrated that writers producing longer units are likely to be proficient learners (Ortega, 2003; Wolfe-Quintero et al., 1998). Therefore, international scholars having higher proficiency levels are expected to produce units that are more or less equal to native speakers (Lu, 2011). In that sense, although the Turkish scholars produced shorter production units than native writers and differed significantly from them, their performance can be considered satisfactory considering C1 and C2 level French EFL learners' MLC performance (C1=11.16, C2=11.50) in Paquot's study (2017). Similarly, they also produced longer units than both the advanced NNS and NS students in Mancilla et al.'s study (2017) in which they compared class discussions of 243 graduate-level students in the education discipline across 6 years. As a result, it can be concluded that Turkish scholars produced a level of the length of production units between advanced nonnative writers and native writers. The use of longer production units is consistent with the previous studies arguing that academic writing is characterized by longer production units (Brown & Yule, 1983; O' Donnell, 1974)

In addition, Turkish scholars were compared among themselves in order to see whether there were any disciplinary changes. Accordingly, TR_EDU subcorpus differed significantly from TR_ENG and TR_MED subcorpora in MLS and MLT and from TR_MED subcorpus in MLC ($p < .0036$). These findings point out the variation among the disciplines as the rate of recurrence, and a variety of complexity properties may vary from one discipline to another (Ansarifar et al., 2018). The typical characteristics of disciplines are shaped by different research concerns and traditions in that, though both of them are technical disciplines, engineering deals with generalizations in science in order to seek a solution to real-world issues, whereas science concentrates on the natural world to find out and explain how it works (Biber & Conrad, 2009). In addition, the studies in the literature have revealed numerous textual features that vary in the sciences and humanities (Charles, 2003; Hyland, 2000; MacDonald, 1994). In that sense, the findings of the present study support the claim that variation exists among the disciplines.

The second syntactic complexity category was the amount of subordination which was measured using four indices: C/T, CT/T, DC/C, and DC/T. The results demonstrated that the MCWP and the reference corpus differed significantly in all of the measures in favor of the

reference corpus. In other words, the native writers used more subordination than Turkish scholars. Statistically significant differences were also found for the subcorpora except for all of the measures in TR_COM and one measure (C/T) in TR_MED. These findings contrast with Wu et al.'s (2020) study in which statistically significant difference was found only for C/T measure though NS group employed a greater amount of subordination. Similarly, Lu and Ai (2015) did not find any significant differences for any of the subordination measures between the NNS and NS corpora. However, when they took L1 into consideration, it was found that Chinese, Japanese and Russian learners used significantly less amount of subordination compared to native writers. Similarly, making a proficiency-based comparison, Ai and Lu (2013) reported that native writers employed a greater amount of dependent clauses in clauses and T-units than Chinese learners. In addition, Mancilla et al. (2017) found statistical differences between the NNS group and the NS group in two subordination measures: DC/T and DC/C, arguing that the NNS group used a lesser amount of dependent clauses. A potential explanation of the use of a lesser amount of subordination by Turkish scholars may be that they did not perceive the importance of subordination in academic writing (Hyland, 2002b) as they produce significantly lesser subordination in their manuscripts.

In addition, Turkish scholars were compared among themselves in order to see whether there were any disciplinary changes. It was found that TR_COM and TR_EDU subcorpora used a greater amount of subordination than TR_ENG and TR_MED subcorpora. This finding is in line with Hundt et al. (2016) and Seoane and Hundt (2018), who pointed out the differences between hard and soft sciences. The language is used as a way to interpret, organize, and clarify a number of aspects of the subject in social studies in which writers are required to make analysis and synthesis (Lent, 2015). Another potential explanation is that Turkish scholars might follow different developmental patterns as to syntactic complexity regardless of proficiency (Lu & Ai, 2015; Polat et al., 2019). However, additional research is required to completely understand the nature of this difference.

The third category was the amount of coordination, which was measured using three indices: CP/C, CP/T, and T/S. The results revealed that the only measure the MCWP had a higher mean value was CP/C. However, a statistically significant difference was found only for the T/S measure in favor of the reference corpus. These results suggest that Turkish scholars used significantly fewer T-units per sentence than native scholars in their research papers as well as more coordinate phrases per clause and less coordinate phrases per T-unit though

they were not significant. These findings are, in some part, contrary to those in the literature. For example, Ai and Lu (2013) reported that nonnative students and native students differed only in the CP/T index and that nonnative students had higher means only in T/S index though it was not significant. They also added that no difference for any measures was found between the proficiency levels. Besides, the nonnative students in Mancilla et al.'s study (2017) was found to employ a greater amount of coordinate phrases; however, they significantly differ only in CP/C measure from their native counterparts. Furthermore, Lu and Ai (2015) stated that the nonnative group and native group did not differ in terms of the number of coordination measures. However, when they considered L1 of the participants, they reported significant variations. It was found that similar to Turkish scholars, the Chinese group produced significantly fewer T-units per clause than native speakers. In addition, Wu et al. (2020) stated that ELF scholars in their study had significantly higher mean ranks in all of the number of coordination measures. The literature on coordination argues that beginner writers or those having lower L2 proficiency use greater amounts of coordinate sentences in their writing (Bardovi-Harlig, 1992; Norris & Ortega, 2009). Therefore, the frequent use of coordination in writing can be regarded as a characteristic of a nonnative usage of English (Ai & Lu, 2013; Mancilla et al., 2017). However, such an argumentation should be put forward carefully since coordination plays distinctive roles in academic writing in that it helps increase clarity by connecting the same sort of patterns (Halliday & Hasan, 1976). Overall, it can be said that Turkish scholars used comparable amounts of coordination with native scholars.

Turkish scholars were compared among themselves in order to see whether there were any disciplinary changes. The only significant difference was found for CP/T measure, revealing that TR_EDU used a greater amount of coordinate phrases per T-unit than TR_ENG and TR_MED. There are conflicting findings as to the variation among coordination in the literature. For example, investigating the changes in syntactic complexity through the school year, Lu (2011) concluded that CP/C and CP/T measures are indicatives of proficiency and differentiate school grades, but no significant differences exist for T/S measure. However, Neary-Sundquist (2017) did not report significant differences in the coordination amount between intermediate, advanced, and highly advanced second language learners. The reason for this finding may be attributed to the fact that research in the field of education may be anecdotal and aims to find out the best classroom practices and thus requires more explaining

and clarification than engineering and medicine disciplines, which are relatively more technical and procedural (Biber & Conrad, 2009).

The fourth syntactic complexity category was the degree of phrasal sophistication which was measured using four indices: CN/C, CN/T, and VP/T. The results demonstrated that native speakers used significantly more complex phrases than Turkish scholars. There are conflicting results regarding the use of complex phrases in the literature in that some studies presented results claiming that nonnative learners use more complex phrases while others reported results in favor of native speakers. For example, the nonnative students in Mancilla et al.'s study (2017) were found to use significantly more complex nominal per clause. Similarly, Wu et al. (2020) reported that ELF researchers employed a greater amount of complex nominals but the fewer amount of verb phrases than the native researchers. On the other hand, Ai and Lu (2013) found that native and nonnative writers significantly differed in CN/C and CN/T measures, which were also proved to discriminate the proficiency levels. They also added that more proficient second language learners employed a significantly greater amount of complex nominal per T-unit and that they used more complex nominals per clause though it was not significant. Likewise, Lu and Ai (2015) found that native speakers used significantly more complex nominals in their writing than nonnative students. In addition, taking the influence of L1 on the syntactic complexity, they reported variation among different L1 groups. Similar to Turkish scholars, Japanese and Tswana groups were found to use significantly less complex nominals, and Japanese, Russian, and Chinese groups employed significantly fewer verb phrases per T-unit than native speakers.

Contrary to the main corpora, the results as to the subcorpora are conflicting. No statistical significance was found for any of the measures in Communication and Economics corpora, whereas Engineering and Medicine corpora differed in all three measures. Finally, it was found that the Education subcorpora differed in two measures: CN/T and VP/T. The degree of phrasal sophistication has been reported to be influenced by topic (Yang et al., 2015), grammatical instruction (Vyatkina, 2013), genre (Lu, 2011; Yoon & Polio, 2016), discipline (Karakaya, 2017), proficiency in a second language (Lu & Ai, 2015) as well as being a strong predictor of quality of the papers and changing over time (Bulté & Housen, 2014; Yang et al., 2015). Although the effects of topic and genre were eliminated by including only research articles of the same topic in the corpus, more research is needed to understand the nature of this variation in the subcorpora.

Turkish scholars were compared among themselves in order to see whether there was any variation among them in terms of the degree of phrasal sophistication. The findings demonstrated that scholars in the Education discipline used significantly more complex nominal per T-unit than scholars in Engineering and Medicine. It was also revealed that scholars in Economics and Education discipline used significantly more verb phrases per T-unit than scholars in Engineering and Medicine. No statistically significant difference was found for the CN/C index. These findings indicate disciplinary variation in the degree of phrasal sophistication among Turkish scholars. Similar to the present study, based on a corpus study in which three disciplines were compared, Karakaya (2017) reported disciplinary changes in the use of complex phrases concluding that scholars in agronomy used more complex nominals compared to scholars in applied linguistics and industrial and manufacturing systems engineering. However, as mentioned before, there are a number of variables that have an effect on the degree of phrasal sophistication measurements, and thus more research is needed to fully understand the nature of this variation.

Finally, overall sentence complexity results showed that the Turkish scholars employed significantly fewer clauses per sentence. The situation is also similar in the subcorpora except for Communication in which no statistical difference was found though the native writers were found to have a higher mean than Turkish scholars. The results run contrary to those in the literature, which reveals no significance as to this measure. For example, both Wu et al. (2020) and Lu and Ai (2015) did not report any statistically significant differences in their studies, although it was found that nonnative groups used more clauses per sentence in both of the studies. However, taking the L1 effect into consideration, Lu and Ai (2015) found that Chinese and Japanese groups used significantly fewer clauses per sentence than native speakers.

Turkish scholars were compared among themselves in order to see whether any disciplinary changes exist or not. The results revealed that Turkish scholars in Economics and Education disciplines showed significantly more syntactic complexity than those in Engineering. In addition, Turkish scholars in Medicine showed significantly more syntactic complexity than scholars in Education. Similar to other syntactic complexity categories, disciplinary variation is evident for overall sentence complexity as well.

In sum, native scholars showed significantly higher syntactic complexity in all of the measures except for two amount of coordination indices (CP/C and CP/T) in which Turkish scholars were found to use comparable amounts of coordinate phrases. The most similar

study to the present one was Wu et al.'s study (2020) who compared ELF scholars and native scholars in terms of syntactic complexity. The results of the present study are in contrast with Wu et al. (2020) in which ELF scholars were found to use more syntactically complex structures than native scholars. In this context, it seems that the case of Turkish scholars does not fit the conventions of ELF scholars described in Wu et al. (2020). However, as clearly proved by Lu and Ai (2015) there exists a significant relationship between the L1 of the writers and syntactic complexity measures. Considering that the corpus used in Wu et al. (2020) included 10 first languages, 9 of which belongs to European languages, and 5 disciplines that were different from the ones in the present study, it can be put forward that Turkish scholars do not have to conform to the norms described in Wu et al. (2020). In this sense, the present study gives support to Lu and Ai's (2015) proposal that the L1 of the writers should be taken into consideration while investigating the syntactic complexity of the second language learners. Turkish scholars were found to have remarkable similarities with upper intermediate Japanese and Chinese learners in Ai and Lu's study (2013). A similar finding was stressed by Uysal (2008). Investigating the effect of culture on writing, she concluded that Turkish writers shared a number of similarities with Japanese and Chinese writers. However, attributing the reasons for the findings of the present study only to L1 may be misleading since, as mentioned before, there are several variables influencing syntactic complexity. Therefore, more research is required to completely uncover the reasons for the difference underlined in the present study.

The present study revealed variation among Turkish scholars in terms of discipline as well. It was found that scholars in education and economics used significantly more syntactically complex features than those in engineering and medicine. In this regard, this study supports the claim that disciplinary variation exists in scholarly writing (Biber & Conrad, 2009; Bulte & Housen, 2012; Karakaya, 2017; Lu et al., 2020). However, this finding should also be dealt with carefully since learners do not necessarily develop their syntactic complexity skills alike in all categories and areas (Lu & Ai, 2015; Polat et al., 2019). For example, Polat et al. (2009) investigated the syntactic development of Turkish EFL learners in a period of three semesters and found that they did not develop in a linear and incremental way. Consequently, there is a need for further studies to explain the reasons for this disciplinary variation.

5.3. Computer-Aided Error Analysis

The last linguistic dimension investigated in this study was the errors Turkish scholars committed by using computer-aided error analysis. Louvain Error Tagging Taxonomy was used to investigate the errors in the manuscripts. In this context, the errors were investigated in eight broad categories, which are further divided into subdomains.

The overall results showed that the most problematic three categories were Grammar, Punctuation, and Register, respectively. These findings are consistent with other studies in the literature (Aziz, Kashif & Aijaz, 2016; Kasperavičienė & Motiejūnienė, 2013; McDowell, 2016; Mestre, 2011). In addition, they made 20.08 errors per 1000 tokens with an average of 110.15 errors per text. The average number of errors shows that the Turkish scholar made a considerable amount of errors in their manuscripts. However, McDowell and Liardét (2020) reported that the Japanese scholars in their study made 31 errors per 1000 words. In this regard, it can be said that Turkish scholars made fewer errors than Japanese scholars. Nevertheless, the manuscripts submitted to scientific journals are expected to be accurate, clear, concise, and appropriate (Andrews & Blicke, 1982; Zall, 1963). Otherwise, they are most probably requested to be revised or at least checked by native speakers, if they are not rejected by the journal gatekeepers (Bordage, 2001; Duszak & Lewkowicz, 2008; Li, 2005; Römer, 2001). Therefore, it can be concluded that manuscripts submitted to journals by Turkish scholars need improvement in terms of accurate use of English. In addition, Turkish scholars were compared among themselves. The results showed that scholars in Communication committed the greatest number of errors while those in Economics made the fewest number of errors. The other disciplines were found to have comparable amounts of errors.

The investigation of Form errors shows that Turkish scholars made 0.96 form errors per 1000 tokens, with almost all of the errors in this category belonging to the spelling category. In addition, scholars in the Medicine discipline committed the greatest number of Form errors. The reason for this finding may be that they use more technical vocabulary than other disciplines. Examples of spelling errors are as follows:

- (1) "...can be defined as ... location based (location-based) ... application... ." (TR_)
- (2) "Age, Height (height) and Weight (weight) are the other details... ." (TR)
- (3) "... official narratives ... by Turkish serilas (serials)." (TR_COM)

The studies on second language learners' spelling errors have demonstrated that: i) there exists a wide variety of spelling error types ii) second language learners make more form errors than native speakers and iii) a great majority of form errors are caused by a lack of proficiency in the second language (Bestgen & Granger, 2011). Although there are a number of studies dealing with form errors in the literature (Al-Jarf, 2009; Botley & Dillah, 2007; Hovermale & Martin, 2008), a considerable amount of them did not report normalized numbers, making a possible comparison unreliable. To the best of my knowledge, only Bestgen and Granger (2011) reported normalized occurrences of errors in their study. They investigated spelling errors in 223 essays of high-intermediate or advanced EFL learners from 3 different L1 backgrounds and found that they made 10.7 errors per 1,000 words. In that sense, it can be argued that Turkish scholars made a fewer number of form errors.

In order to overcome form problems, spelling corrector tools such as Microsoft Word spell checker, Grammarly, or Ginger Software Spelling Check can be used since submitted articles are required to be error-free. However, scholars should be careful in using these tools as most of the tools are designed to address the needs of native speakers. However, the needs of native speakers and second language learners are different, and thus, the available tools for spelling correction may fail to deal with some of the errors international scholars make (Rimrott & Heift, 2005).

The second category was Grammar. The investigation of grammar errors demonstrated that Turkish scholars experienced significant problems in this category. They made 9.78 grammar errors per 1000 tokens. Grammar errors were found to be the most problematic category both in the main corpus and the subcorpora. Turkish scholars had difficulty, especially in articles (4.68 per 1000 tokens), verbs (3.04 per 1000 tokens), and nouns (1.14 per 1000 tokens), respectively. Such a finding is expected for a number of reasons.

First, it has been stated in the literature that mastering the English article system is among the most difficult problems experienced by a non-native speaker (Han, Chodorow & Leacock, 2006). It has also been stated that even highly proficient learners have considerable difficulty in acquiring the English article system (McEnery, Xiao, and Tono, 2006). Furthermore, it even becomes more problematic for learners whose mother tongue does not contain articles (Chaudron & Parker, 1990; Liu & Gleason, 2002). In this sense, contrary to many Indo-European languages, both definite and indefinite articles are absent from Turkish. Therefore, it is quite reasonable for Turkish scholars to have difficulty in the correct use of articles.

It was found that a great proportion of the article errors were missing errors. In other words, Turkish scholars did not use an article when they were supposed to use, as shown in (4) and (5).

(4) “(The) Outputs from the studies can be summarized as follows...” (TR_ENG)

(5) “(The) Lecturer expressed that the students...” (TR_EDU)

This finding is in line with other studies in the literature. For example, Kırkgöz (2010) investigated 120 essays written by beginner-level adult Turkish EFL learners and found that almost a quarter of the errors were article errors. Similarly, Barrett and Chen (2011) investigated the article errors of Taiwanese EFL learners and found that they experienced difficulties in using the English article system properly. Investigating 304 texts of Spanish learners of English, MacDonald (2016) found that the use of articles was among the most problematic language domains. Furthermore, Crompton (2011) concluded in his corpus-based study that advanced level Arabic learners of English did not fully master the use of articles though the Arabic language has an article system.

Second, verbs have specific functions in English academic writing. Academic writers have to use particular verb tenses or voices in different circumstances (Reid, 2000). For example, writers are required to use the present tense in the abstract or background information of the manuscripts and passive voice to present indirectness and objectivity in academic texts (Hinkel, 2013). In addition, guidelines on how to write a research paper recommend using past tense for describing the methods, presenting the results of the study, and for acknowledgement of sources, whereas present perfect tense is recommended for previous findings which have regarded as scholarly knowledge or general truth (Shaw, 2013). In this sense, verb errors are regarded as severe as they hinder communication in the manuscript (Reid, 2000). Therefore, using the appropriate verb tense or voice requires knowledge of research article rhetoric as well as mastery of English grammar.

As shown in (6), (7), and (8), Turkish scholars seem to suffer from a lack of rhetorical knowledge on verbs rather than grammatical knowledge. All of the samples below were excerpted from the methods section of the manuscripts. As is seen, Turkish scholars preferred to use the present tense while describing the method of their study where they were supposed to use the past tense. The errors below are clearly due not to lack of proficiency in English but to a lack of knowledge on research article rhetoric.

(6) “For this purpose, ... performance is (was) evaluated by ...” (TR_EAS)

(7) “Furthermore, males are (were) rougher and ...” (TR_EDU)

(8) “The participants consist (consisted) of 3 women and 6 men.” (TR_COM)

In this sense, academic writers should be provided with a detailed explanation about when and how specific tenses are used in academic writing. However, such an explanation is quite limited in L2 textbooks, which instead include a section on verb tenses focusing on description and exercises on tenses (Hinkel, 2013). Therefore, academic writing courses should be designed to teach the functions of verbs in academic writing.

Third, academic writing is characterized by nominalization, which is typically used while making generalizations, talking about actions, and describing processes (Hinkel, 2013). Therefore, academic writing requires a frequent use of nouns and noun phrases (Alderson, 2007). The investigation of noun errors shows that Turkish scholars had difficulty, especially in the number subcategory (0.98 per 1000 errors) as presented in (9) and (10).

(9) “In the present study, growth ... for females and males ... for all gender (genders).” (TR_EDU)

(10) “These subdivision (subdivisions) of ... more detailed analyses and ... results.” (TR_EAS)

This finding is in line with Kırkgöz’s study (2010) in which she reported that Turkish university students had difficulty in pluralization. In addition, Sajid (2016) found in his error analysis study that 7.85 % of the errors were related to pluralization. Similarly, Bond (2016) found that number errors ranked second in his error analysis study in which he investigated the errors Chinese students in the United Kingdom made in their writing. Furthermore, investigating a small corpus of Chinese EAP learners, Chuang and Nesi (2006) found that number errors accounted for 8.8% of all errors made in the corpus. The reason for this kind of error is related to the fact that the Turkish language does not require pluralization of nouns, whereas learners have to add an affix, in most case –s or –es, to the end of a noun to make that noun plural in English. Therefore, it is seen clearly in (9) and (10) that Turkish scholars applied the grammatical rules in their first languages to English while making nouns plural.

The third computer-aided error analysis category was Lexis. Lexical errors are regarded as significant in second language learning since they are among the most frequent errors in non-native English (Webber, 1993). In addition, it has been put forward that lexical structures are among the most problematic aspects in second language learning (Carrió-Pastor & Mestre-Mestre, 2014). The lexicon is even more important in academic writing as academic

achievement requires a rich vocabulary knowledge (Laufer & Nation, 1995; Leki & Carson, 1994). Lexical errors and errors in lexical items are also important in terms of text quality since using proper and correct lexical items is a sign of expertise a writer has in the scientific discipline (Kasperavičienė & Motiejūnienė, 2013).

It was found that Turkish scholars made 1.74 lexis errors per 1000 tokens. The comparison of the scholars revealed that scholars in the Communication field made the greatest number of lexis errors while those in Economics made the fewest number. It was found that the most frequent lexis errors belonged to lexical single category. In other words, the Turkish scholars had difficulty in using the proper words in their manuscripts. The errors in other lexical categories were relatively rare. (11), (12) and (13) are samples of erroneous use of single lexical items.

(11) "... is ensured with ... analytics (analysis)." (TR_ENG)

(12) "Participators (Participants) of this study consists ..." (TR_EDU)

(13) "... the analysis is made (repeated) by ..." (TR_EAS)

Especially (11) and (12) shows that Turkish scholars failed in using the correct lexical item, and instead, they used an item having a similar meaning to the one they should have used. This finding is consistent with other studies in the literature. For example, Carrió-Pastor and Mestre-Mestre (2014) reported similar findings and concluded that B2 level learners in their study had difficulty in finding the appropriate lexical item for the term even though they were proficient enough to use that word. In line with the results of the present study, Chuang and Nesi (2006) found that lexical errors accounted for 9.1% of the errors in their corpus and 63.1% of the lexical errors were due to misuse of a lexical item. Furthermore, MacDonald (2016) found that lexis was the second most problematic category for Spanish learners of English in her study. Contrary to the aforementioned studies, lexical errors were the fourth problematic category in the present study, which shows that Turkish scholars made relatively less lexical errors. Furthermore, both MacDonald (2016) and Mediero Durán (2013) revealed that learners with higher proficiency levels made less lexical errors than those with lower proficiency levels concluding that lexis errors tended to decrease in higher proficiency levels. Such a conclusion also explains the fewer number of lexical errors Turkish scholars made.

The other category was Register, which was found to be the third problematic category. Register, in its broad sense, refers to a variety related to a specific situation of use (Biber &

Conrad, 2009). In this sense, scholars are expected to adopt an academic register and use an academic language in their manuscripts. Therefore, the words, grammatical structures, and phrases they use need to fulfill the requirements of the academic language. In other words, scholars have to use more formal and precise words, grammatical structures, and phraseology in their research articles (Gilquin & Paquot, 2008).

Turkish scholars made 1.86 register errors per 1000 tokens a significant amount of which belonged to the vocabulary category. The results showed that the discipline that experienced the most difficulty in the Register category was Communication, while scholars in the Economics discipline made the fewest number of errors. (14), (15), and (16) show that Turkish scholars used informal words, structures, and phrases in their manuscripts.

(14) "..., we will examine the findings of the other researches (previous studies)" (TR_ENG)

(15) "Unlike other studies, ... rather than financial ratios to make it easy of (for the ease of) interpretation." (TR_EAS)

(16) "So (Hence) ... must be taken for ..." (TR_EDU)

This finding is in line with Hemchua and Schmitt's study (2006). They investigated 20 Thai EFL learners' lexical errors and found that the most frequent error type was *near synonyms*, which was defined as using informal and improper words. The present finding is also supported by Zughoul (1991), who stated that Arabic learners of English made the greatest number of errors in the assumed synonymy category. Furthermore, Severino and Prim (2015) reported that register errors made up 6% of word choice errors in their study. International scholars have to overcome the effects of their first languages and everyday conversational language in order to meet the requirements of the academic register (Yore & Treagust, 2006). Therefore, academic writing courses should focus on teaching academic register since as a result of globalization and the spread of technology, learners are exposed to different registers of English.

The fifth error category was Style. Turkish scholars made 0.45 errors per 1000 tokens in this category. Style errors were investigated using three subdomains: incomplete, unclear, and other. As shown in (17) and (18), the majority of the errors in this category were found to be in the unclear category. In addition, an example of the incomplete category is shown in (19).

(17) "This result indicates whether the firms are operating increasing or decreasing return to scale" (TR_EAS)

(18) “In the current study, 32 different ANN models which with four factors have 2 or 4 levels were trained and tested.” (TR_ENG)

(19) “The models, respectively:” (TR_COM)

Furthermore, investigation of the disciplines revealed that scholars in communication and education fields made the greatest number of style errors, respectively. Such a finding may be attributed to the fact that both communication and education are soft disciplines and thus rely more on qualitative studies. On the other hand, economics, engineering, and medicine are more procedural and deal extensively with numbers instead of language. Therefore, scholars in engineering and medicine made fewer style errors than those in the communication and education fields.

The findings of the present study are inconsistent with other studies in the literature. For example, Lopez (2009) investigated the Spanish part of the International Corpus of Learner English (ICLE) and found that almost 7% of the errors were Style errors. Similarly, Lee et al. (2016) analyzed TOEFL essays of learners with 5 different first languages and found an average number of 0.68 style errors. In this sense, it can be argued that Turkish scholars performed well in this category. As one of the participants in Huang’s study (2017) expressed, international scholars, even experienced professors, frequently make style errors and are asked to clarify the statements in their texts. However, the scholars seeking publication are expected to conform to the journals’ style norms, and therefore, style errors in the submitted manuscripts are not welcomed by journal editors and referees. In addition, style errors may be potentially regarded as an indicator of low-quality language in the text (Kasperavičienė & Motiejūnienė, 2013). Therefore, international scholars should be encouraged to avoid long unclear sentences causing ambiguity or misunderstanding.

Another category was the Word. It was found that Turkish scholars made 1.55-word errors per 1000 tokens. Excerpt (20) is an example of a missing error, and (21) and (22) show redundant errors.

(20) “ ... mine (is) still operational.” (TR_END)

(21) “... have some additions for this (0) Figure 1.” (TR_EAS)

(22) “... attitude level towards mathematics lesson (0), ...”. (TR_EDU)

The majority of the errors in this category were annotated in the Redundant subcategory (0.95 per 1000 tokens), whereas the errors in the Order and Missing subcategories were

comparable (0.30 and 0.29 errors per 1000 tokens, respectively). These findings are consistent with Lee et al.'s study (2016). Similar to the present study, the participants in their study made the greatest number of word errors in the redundant category. However, Lopez (2009) found that the participants in her study made the greatest number of word errors in the missing category.

Redundant errors are considered important in academic writing for a number of reasons. First, they are signs of low-quality writing and may confuse the readers. Second, they may potentially weaken the sentences as well as thoughts in the manuscripts (Xue & Hwa, 2014). Third, the flow of the writing may be interrupted by redundant words, paving the way for distraction and annoyance. Finally, they decrease the conciseness and precision of the texts (Dawson, 1992). The reason for redundant errors may be the rhetorical differences between Turkish and English. In Turkish academic convention, using numerous words and long sentences is common to define a term or concept (Güngör & Uysal, 2016). Similar to academic writing in Turkey, Turkish scholars may have thought that using numerous words would contribute to the formality and quality of their manuscripts. Therefore, Turkish scholars should be informed on the rhetorical requirements of writing a research article in English. The reason for order errors, on the other hand, maybe attributed to the fact that word order is different in Turkish and English. Therefore, Turkish scholars may have had difficulty in ordering words correctly. Finally, as shown in (21), a great majority of the missing errors were due to carelessness. Turkish scholars simply forgot using the required word.

Lexicogrammar was another computer-aided error analysis category. As a term used in Systemic Functional Linguistic (Halliday, 1961), lexicogrammar refers to a view in which lexicon and grammar are seen as two essentially associated parts of a particular body (Sinclair, 1991). Grammatical structures are potentially constrained by lexical items (Francis, 1993), and similarly, the lexicon is, most of the time, grammatical in that a grammatical consequence follows the use of a specific lexicon (Biber, Conrad & Reppen, 1998). In this context, it is argued that different words have distinctive associations, and each word has distinctive grammar, and thus typical contextual patterns are the only ways of acquiring the grammar (Aston, 2001). Furthermore, communicative assumptions of the discourse communities control lexicogrammatical choice, and thus each discourse community and each genre need to use appropriate lexicogrammatical units to communicate functionally (Park, 2010). In this sense, Flowerdew (2001) regards writers' insufficient

understanding of characteristic lexicogrammatical settings of words as one of the three areas of difficulty in academic writing. Therefore, scholars are expected to make correct and appropriate lexicogrammatical choices in their manuscript to communicate successfully with the rest of the discourse community. In this regard, even slight deviations from lexicogrammatical accuracy may lead to harm the expertise and comprehensibility of a manuscript and, in turn, may result in revision request or rejection (Ehara & Takahashi, 2007).

It was found that Turkish scholars made 1.33 lexicogrammatical errors per 1000 tokens. A detailed investigation of lexicogrammar errors revealed that Turkish scholars had difficulty in verb category as shown in (23) and (24) as well as using the correct dependent preposition as shown in (25) and (26).

(23) "... was applied for (applied to) examination of ..." (TR_EDU)

(24) "financial indicators are generally used to (used as) performance indicator." (TR_EAS)

(25) "Maybe the teacher of A is not aware (aware of) and ..." (TR_EDU)

(26) "...their model is superior on (superior to)" (TR_ENG)

These findings are supported by other studies in the literature. For example, Lee et al. (2016) found that verbs were the category where learners experienced the most difficulty in lexicogrammar. Similarly, Chan (2010) found that the misuse of prepositions was one of the difficulties Hong Kong Cantonese EFL learners experience. In addition, Chang and Nesi (2006) reported that preposition errors made up 8.1% of the errors in their study. More recently, McDowell and Liardét (2020) listed preposition errors as one of the areas international scholars experience difficulty. Based on these findings, it should be suggested that academic writing courses in Turkey should concentrate on teaching words with their dependent prepositions.

The final category was punctuation. Punctuation has, for a long time, been ignored by second language researchers due to the common belief that it is not important in language teaching or it can be learned without premeditation. (Hirvela, Nussbaum & Pierson, 2012; Waugh, 1998). Considering punctuation as boring (Gauthier, 1993), teachers tend to ignore teaching it, and learners do not pay attention to learn it due partly to modern social media interactions in which effective communication can be established and facilitated by using abbreviations and colloquialism without using proper punctuation (Johnson et al., 2017). However, as stressed by Alamin & Ahmed (2012), it is one of the most difficult language components in

second language teaching and learning. Punctuation even plays a more significant role in scientific writing in which the use of formal language is required, and incorrect use of punctuation marks may lead to a change in the meaning (McLaren, 2003).

Punctuation was found to be the second most difficult area for Turkish scholars, who made 2.41 punctuation errors per 1000 tokens in line with other studies in the literature. For example, Chuang and Nesi (2006) listed punctuation among the most problematic linguistic features in their error analysis study. Similarly, Buckingham and Aktuğ-Ekinci (2017) found that punctuation errors were the most frequent error type in their study. Furthermore, investigating writing samples of 60 university students from University College London, López and Manalastas (2017) reported that 50.3% of the total errors were punctuation errors. In addition, Hamed (2018) found that punctuation was among the most common errors detected in his participants' writing samples.

As shown in (27), (28), and (29) great majority of these errors were missing errors. In other words, Turkish scholars did not use punctuation marks where they were supposed to. In addition, they made a considerable amount of redundant errors. As shown in (30) and (31), they used a punctuation mark in a case where it was not needed. Furthermore, (32) and (33) shows confusion errors.

(27) "In this study (,) it is intended to ..." (TR_COM)

(28) "After getting these values (,) the scale efficiency ..." (TR_EAS)

(29) "Moreover (,) this technology initially appears" (TR_ENG)

(30) "... it can be said that, (0) teacher knowledge about ..." (TR_EDU)

(31) " Czech Republic, (0) in 2011." (TR_ENG)

(32) "...opportunities disappeared, (;) ... declined..." (TR_EAS)

(33) "In choosing ... in the study; (,) envying the journalists ... have been effective."
(TR_COM)

The reason for these findings may be attributed to cross-cultural differences in punctuation rules (Hirvela et al., 2012). In other words, second language writers may not differentiate the punctuation rules in their first languages from those in the target language and thus use the punctuation rules of their first languages in their writing. The excerpts above give support to this statement. As clearly seen, Turkish scholars applied the punctuation rules of Turkish in their manuscripts. Excerpt (31) is a typical example of such a situation. It shows that the

writer did not use a comma after *moreover* as punctuation mark is not used after conjunctions in Turkish. This conclusion is also supported by other studies in the literature. For example, both Elkılıç et al. (2009) and Kırkgöz (2010) attributed punctuation errors of their participants to the interference of Turkish, claiming that students used the punctuation rules of Turkish in their English production. In addition, Altunay (2009) acknowledged the role of first language transfer but also gave credit to the fact that the use of connectives in English is a complex issue. Another reason may be a lack of awareness with regard to punctuation since the proper use of punctuation rules is connected to awareness as well as knowledge (Planken et al., 2019). For example, participants in Benzer's (2010) study expressed that even though they had the knowledge of punctuation rules, they did not apply them. Finally, the reason may be attributed to the fact that punctuation errors are considered as treatable and being amenable to be corrected easily (Ferris, 2006). Therefore, this way of thinking may lead writers to underestimate the use of correct punctuation in their manuscripts.

In sum, the investigation of the errors made by Turkish scholars revealed that Turkish scholars had difficulty particularly in articles, verbs, punctuation, and using the correct lexical item. The findings also showed that Turkish scholars lacked adequate knowledge of the rhetoric of research articles in English, which paved the way for making errors. This conclusion is in line with other studies in the literature. For example, Çandarlı (2012) investigated the abstracts written by Turkish educational sciences researchers and found rhetorical differences. Similarly, Kafes (2018) identified conventions of the Turkish discourse community in research article abstracts written by Turkish scholars. In addition, Karsli et al. (2018) examined discussion sections of Turkish and American doctoral theses in the educational technology field and found that Turkish theses showed a number of divergences from their counterparts in the USA. More recently, Erdemir (2019) carried out a move analysis study to investigate the results section of Turkish research articles and concluded that Turkish scholars suffered from a lack of sufficient knowledge with regard to the rhetoric of research articles in English. On the basis of these findings, it is suggested in this study that EAP courses should focus on teaching rhetorical aspects of research articles as well as grammatical and linguistic features in order to overcome the difficulties Turkish scholars experience in writing for publication.

CHAPTER VI

CONCLUSION

This chapter concisely outlines the main contributions of the present study with regard to writing for publication literature. First, the conclusion will be drawn on the basis of the findings of the present study, and then the theoretical and pedagogical implications drawn from the present study will be mentioned. Further limitations recognized during the study will be considered. Finally, suggestions regarding possible topics for future studies will be touched upon.

6.1. Summary of the Study

The aim of the present study was to investigate the linguistics difficulties experienced by Turkish scholars in writing for publication by examining three linguistic features: lexical diversity, syntactic complexity, and grammatical errors. To do that, a corpus of 216 unedited research articles from 5 disciplines written by Turkish scholars was compiled. The selected disciplines were as follows: Communication, Economics, Education, Engineering, and Medicine. Overall, a corpus having 886,482 words was obtained as the MCWP. In order to make a reliable comparison, a reference corpus of 163 published research articles containing 885,791 words were collected as well. Lexical diversity was investigated using three different measures: type/token ratio, standardized type/token ratio, and moving average type/token ratio. In the analysis of syntactic complexity, Lu's (2010) L2 Syntactic Complexity Analyzer (L2SCA) was used. Each article in both the study and the reference corpus was investigated 14 measures in 5 different categories. The categories and measures used in the present study were as follows: Length of production unit: mean length of clause (MLC), mean length of sentence (MLS), mean length of T-unit (MLT); Amount of subordination: clauses per T-unit (C/T), complex T-units per T-unit (CT/T), dependent clauses per clause (DC/C) and dependent clauses per T-unit (DC/T); Amount of coordination: coordinate

phrases per clause (CP/C), coordinate phrases per T-unit (CP/T) and T-units per sentence (T/S); Degree of phrasal sophistication: complex nominals per clause (CN/C), complex nominals per T-unit (CN/T), and verb phrases per T-unit (VP/T); Overall sentence complexity: clauses per sentence (C/S). Finally, grammatical errors were investigated using computer-aided error analysis. Using Louvain Error Tagging Taxonomy, grammatical errors were analyzed in 8 broad categories which were broken down into 54 subdomains. The investigation was twofold. First, the MCWP and the reference corpus were treated as separate bodies and compared. Second, each subcorpus in the MCWP were compared with its counterpart in the reference corpus. In addition, Turkish subcorpora were compared among themselves in order to investigate whether any disciplinary variation existed.

The results of the type/token ratio analyses revealed that scholars in both the MCWP (TTR=0.06, STTR=0.68, MATTR=0.68) and the reference corpus (TTR=0.052, STTR=0.70, MATTR=0.70) used a wide variety of vocabulary in their manuscripts though the MCWP had slightly less lexical variation than the reference corpus. The comparison of Turkish subcorpora with each other demonstrated that scholars in TR_COM (TTR=0.17, STTR=0.70, MATTR=0.70) subcorpus had the highest TTR values while TR_EDU subcorpus showed the lowest values (TTR=0.08, STTR=0.67, MATTR=0.67).

The analysis of syntactic complexity showed that native scholars used significantly more syntactically complex structures in all of the measures except for two amount of coordination indices (CP/C and CP/T) in which Turkish scholars used comparable amounts of coordinate phrases. The investigation of the length of the production unit showed that native writers produced significantly longer sentences, clauses, and T-units than Turkish scholars. The amount of subordination results revealed that the native writers used more subordination in their manuscripts than Turkish scholars. The results regarding the amount of coordination category were somewhat different in that only one measure was found to be statistically significant. It was found that native writers used significantly more T-Units per sentence in their research articles. The degree of phrasal sophistication results indicated that native speakers used significantly more complex phrases than Turkish scholars in all of the indices. Finally, overall sentence complexity results revealed that the native scholars employed significantly more clauses per sentence. The results of the present study support the claim that the first language of the writers is a significant factor in syntactic complexity, and thus it should be taken into account in syntactic complexity research (Lu & Ai, 2015). In this regard, Turkish scholars were found to have remarkable similarities with upper intermediate

Japanese and Chinese learners in Ai and Lu's study (2013). The present study also revealed disciplinary variation among Turkish scholars. It was found that scholars in education and economics used significantly more syntactically complex features than those in engineering and medicine. In this regard, this study supports the claim that disciplinary variation exists in scholarly writing (Biber & Conrad, 2009; Bulte & Housen, 2012; Karakaya, 2017; Lu et al., 2020).

The computer-aided error analysis showed that the most problematic three categories were Grammar, Punctuation, and Register, respectively. In addition, it was found that Turkish scholars made 20.08 errors per 1000 tokens with an average of 110.15 errors per text. The examination of Form errors demonstrates that Turkish scholars committed 0.96 form errors per 1000 tokens. In the grammar category 9.78 errors per 1000 tokens were annotated. Grammar was found to be the most problematic category for Turkish scholars. They experienced difficulty, especially in articles (4.68 per 1000 tokens), verbs (3.04 per 1000 tokens), and nouns (1.14 per 1000 tokens), respectively. Turkish scholars were found to make 1.74 lexis errors per 1000 tokens, and the category they experienced the most difficulty was using the proper words in their manuscripts. The results of the register category showed that 1.86 register errors per 1000 tokens were made. In addition, Turkish scholars committed 0.45 errors per 1000 tokens in the Style category. Regarding the Word category, Turkish scholars were found to make 1.55-word errors per 1000 tokens. It was also revealed that Turkish scholars made 1.33 lexicogrammatical errors per 1000 tokens. A further examination of errors in this category revealed that Turkish scholars experienced difficulty particularly in verb and using the correct dependent preposition sub-categories. Finally, it was found that Turkish scholars made 2.41 punctuation errors per 1000 tokens in the Punctuation category. A detailed investigation of Punctuation errors showed that a significant amount of errors was missing errors, which indicated that Turkish scholars experienced difficulty in using punctuation marks where they were supposed to.

Overall, the examination of the errors committed by Turkish scholars demonstrated that Turkish scholars experienced difficulty, especially in articles, verbs, punctuation, and using the correct lexical item. The findings also indicated that Turkish scholars suffered from sufficient knowledge of the rhetorical aspect of writing a research article in English, which lead to more errors in their manuscripts.

6.2. Theoretical and Pedagogical Implications

The present study has a number of significant theoretical and pedagogical implications for researchers, scholars, instructors, as well as material developers and course designers. To begin with theoretical implications, it was one of the few attempts to investigate unpublished research articles produced by international scholars in the writing for the publication literature. Although there have been several empirical studies dealing with the problems of international scholars in the writing for the publication process, most of these studies examined published research articles, which are brokered by editors and reviewers (Lillis & Curry, 2010). However, in order to discover and understand the conventions and writing practices of international scholars, it is required to investigate unpublished articles that may best characterize the use of English by international scholars (Wu et al., 2020). Furthermore, this study presented one of the first error-tagged corpus in the writing for the publication literature. Tagging all the errors in a corpus is a laborious work (Andersen, 2011); and, for this reason, a great majority of the existing corpora do not have any error-related annotation (Lüdeling et al., 2005). In this sense, the error-tagged corpus presented in this study had a great potential to reveal the language-related difficulties in a more accurate way. In addition, the results of the present study may be used as need analysis. Besides, five different disciplines were included in the present study. The studies in writing for the publication literature have usually investigated one or two disciplines. However, since there is a disciplinary variation among disciplines (Biber & Conrad, 2014), investigating different as many disciplines as possible allowed us to identify the needs of scholars in different disciplines. The results of this study can be used to meet their need in a more accurate way since academic writing instruction become most successful in cases when it is designed to satisfy the needs of particular situations of learners (Sloane & Porter, 2010).

There are also some implications for the instructors. It was found in this study that the participants had insufficient knowledge regarding the conventions and rhetorical requirements of writing a research article in English, which paved the way for them to make linguistic errors in their manuscripts. Furthermore, it was identified that the participants used less syntactically complex sentences than native speakers of English except for two amount of coordination indices (CP/C and CP/T). Such findings may contribute to the issue of whether language teachers or subject teachers should deliver academic writing instruction (Spack, 1988). It has been, most of the time, revealed that subject teachers have insufficient knowledge of both the skills and aspiration to instruct literacy skills, and they are even not

willing to provide feedback on the manuscripts of learners (Hyland, 2013). Therefore, especially in terms of language-related problems, language teachers should be a part of academic writing instruction, at least in the context of Turkey. Furthermore, the present study revealed variation among Turkish scholars in terms of discipline as well. It was found that scholars in education and economics used significantly more syntactically complex features than those in engineering and medicine. In this sense, academic writing instructors should take disciplinary variation into consideration and tailor their teaching in the light of the needs of scholars (Hyland, 2016).

The present study also presents some implications for scholars. It was found in the present study that the scholars made a considerable number of spelling errors. It was stated in the literature that spell checker tools may play a role in correcting spelling errors (Bestgen & Granger, 2011). Therefore, scholars should use spell checkers to overcome their problems in spelling. One easily accessible tool is Microsoft Word. When the scholars activate the spell-checking function, most of their spelling errors can be easily identified by the tool.

Last but not least, the present study offers some implications for material developers and course designers. For example, with regard to the findings of the present study, it can be postulated that a great majority of the verb errors were a result of a lack of rhetorical knowledge. Therefore, academic writing courses should be designed to teach the functions of verbs in academic writing. Furthermore, material developers should consider providing a detailed explanation about when and how specific tenses are used in academic writing in L2 textbooks (Hinkel, 2013). In addition, it was found in the present study that one of the areas the participants had difficulty was punctuation. However, teaching punctuation is an underestimated part of English language teaching as it is often considered as insignificant or being learned spontaneously (Bakla, 2019). Nevertheless, international scholars seeking publication are required to employ punctuation rules in English in an accurate way. Therefore, both material developers and course designers should pay attention to include explicit teaching of punctuation in the curriculum and the textbooks. Also, the present study identified the areas international scholars experienced the most difficulty. Therefore, course designers and material developers should take the findings of the present study into consideration and deal with the problematic areas and language features in particular academic contexts (Hyland, 2016).

6.3. Suggestions for Future Research

The present study has made important implications to the writing for the publication literature by revealing some of the linguistic problems faced by international scholars. However, a number of questions still need to be answered, and prospective studies are recommended to investigate further statements (which can also be regarded as the limitations of the present study).

First, this study focused only on research articles. Nevertheless, as mentioned in the literature review section, writing for the publication covers other kinds of writing, such as books and conference proceedings. Therefore, future studies should consider examining other types of publications in order to explore the phenomena in detail.

Second, the present study used unedited research articles written by native speakers for comparison. However, international scholars also participate in the writing for the publication arena as professionals and compete with native speakers to publish in high-ranked international journals. For that reason, published research articles of international scholars should also be included in future studies.

Third, this study dealt with only three linguistic features: lexical diversity, syntactic complexity, and error analysis. Yet, there is a wide variety of features that can be considered linguistic problems. Therefore, future studies should concentrate on investigating other aspects of linguistic problems.

Forth, the aim of the computer-aided error analysis carried out in the present study was to provide a general picture of the problems of international scholars; and thus, the sources of errors were not investigated. Future studies may address the source of the errors in order to provide a deeper understanding of the errors.

Fifth, five disciplines were included in this study. It has been stated that disciplinary variation exists in writing for the publication literature (Biber & Conrad, 2014). Therefore, in order to reveal the nature of this variation, more disciplines should be compared and contrasted in future studies.

Last, although The Multidisciplinary Corpus of Writing for Publication is representative for English for Research Publication Purposes in that it contains 886,482 words in total, two subcorpora, namely communication and economics, are relatively small in size. Hence, future studies should consider building larger corpora for these disciplines in order to make the size more representative of the writing practices of the scholars in these disciplines.

In conclusion, few studies have attempted to investigate unedited research articles produced by international scholars. Thus, there is inadequate information about the nature of the writing for the publication process from the perspective of international scholars. Therefore, more studies, including larger and better-balanced data regarding international scholars, are needed in order to provide a more reliable and detailed analysis of the phenomena.





REFERENCES

- Adams, R., Alwi, N. A. N. M., & Newton, J. (2015). Task complexity effects on the complexity and accuracy of writing via text chat. *Journal of second language writing, 29*, 64-81.
- Adel, A. (2008). Involvement features in writing: Do time and interaction trump register awareness?. *Language and computers studies in practical linguistics, 66*, 35.
- Ädel, A. (2010). Using corpora to teach academic writing: Challenges for the direct approach. In M. C. Compoy-Cubillo, B. Belles-Fortunato & M. L. Gea-Valor (Eds.), *Corpus based approach to English language teaching* (pp. 39-54). London: Continuum International.
- Ahmad, U. K. (1997). Research article introductions in Malay: Rhetoric in an emerging research community . In A. Duszak (ed.), *Culture and Styles of Academic Discourse*. 273 – 301. Berlin : Mouton de Gruyter.
- Ai, H., & Lu, X. (2013). A corpus-based comparison of syntactic complexity in NNS and NS university students' writing. In A. Díaz-Negrillo, N. Ballier, & P. Thompson (Eds.), *Automatic treatment and analysis of learner corpus data* (pp. 249-264). Amsterdam/Philadelphia: John Benjamins.
- Aijmer, K. (2002). Modality in advanced Swedish learners' written interlanguage. In S. Granger, J. Hung & S. Petch-Tyson (Eds.) *Computer Learner Corpora, Second Language Acquisition and Foreign Language Teaching* (pp. 55-76). Amsterdam & Philadelphia: Benjamins.
- Aijmer, K. (2018). The Swedish modal auxiliary ska/skall seen through its English translations. *Bergen Language and Linguistics Studies, 9*(1), 139-154.

- Alamin, A., & Ahmed, S. (2012). Syntactical and punctuation errors: An analysis of technical writing of university students science college, Taif University, KSA. *English Language Teaching*, 5(5), 2-8.
- Alderson, J. C. (2007). Judging the frequency of English words. *Applied Linguistics*, 28(3), 383-409.
- Alfaifi, A. Y. G. (2015). *Building the Arabic learner corpus and a system for Arabic error annotation*. Doctoral Dissertation. University of Leeds.
- Al-Jarf, R. (2009). Enhancing freshman students' writing skills with a mind-mapping software. In *Conference proceedings of eLearning and Software for Education «(eLSE)»* (No. 01, pp. 375-382). "Carol I" National Defence University House.
- Allen, J. P., & Widdowson, H. G. (1974). Teaching the communicative use of English. *IRAL-International Review of Applied Linguistics in Language Teaching*, 12(1-4), 1-22.
- Altunay, D. (2009). *Use of connectives in written discourse: A study at an ELT department in Turkey*. Doctoral Dissertation, Anadolu University, Turkey.
- Ammon, U. (2000). Towards more fairness in international English: Linguistic rights of non-native speakers? In R. Philipson (Ed.), *Rights to language: Equity, power and education* (pp. 111-116). Mahwah, NJ: Lawrence Erlbaum.
- Ammon, U. (2012). Linguistic inequality and its effects on participation in scientific discourse and on global knowledge accumulation—with a closer look at the problems of the second-rank language communities. *Applied Linguistics Review*, 3(2), 333-355.
- Andersen, Ø. E. (2011). *Grammatical error prediction* (No. UCAM-CL-TR-794). University of Cambridge, Computer Laboratory. Retrieved from <https://www.cl.cam.ac.uk/techreports/UCAM-CL-TR-794.pdf>
- Anderson, J., & Poole, M. (2009). *Assignment and thesis writing*. Cape Town: Juta.
- Andrews, D. C., & Blicke, M. D. (1982). *Technical writing: Principles and forms*. Macmillan.
- Ansarifar, A., Shahriari, H., & Pishghadam, R. (2018). Phrasal complexity in academic writing: A comparison of abstracts written by graduate students and expert writers in applied linguistics. *Journal of English for Academic Purposes*, 31, 58-71.

- Anthony, L. (2014). *AntConc (Version 3.4. 3)[Computer software]*. Tokyo, Japan: Waseda University.
- Arik, B. T. & Arik, E. (2014). The role and status of English in Turkish higher education. *English Today*, 30(4), 5-10. doi:10.1017/S0266078414000339
- Armstrong, T. (2015). Peer feedback in disciplinary writing for publication in English: The case of 'rolli', a German-L1 novice scholar. *Journal of Academic Writing*, 5(1), 86-105.
- Aston, G. (2001). Learning with corpora: An overview. In G. Aston (Ed.), *Learning with corpora* (pp. 6–45). Houston, TX: Athelstan.
- Aston, G. (Ed.). (2001). *Learning with corpora*. Houston, TX: Athelstan.
- Atkinson, D. (1997). A critical approach to critical thinking in TESOL. *TESOL Quarterly*, 31, 71–94.
- Atkinson, D. (2003). Writing and culture in the post-process era. *Journal of Second Language Writing*, 12(1), 49–63.
- Atkinson, D. (2003). Writing and culture in the post-process era. *Journal of Second Language Writing*, 12(2003), 49–63.
- Aull, L. L., & Lancaster, Z. (2014). Linguistic markers of stance in early and advanced academic writing: A corpus-based comparison. *Written Communication*, 31(2), 151-183.
- Awelu, A. (2011). *Academic writing in English*. Stockholm, Sweden: Lund University.
- Aydinli, E., & Mathews, J. (2000). Are the core and periphery irreconcilable? The curious world of publishing in contemporary international relations. *International Studies Perspectives*, 1(3), 289-303.
- Aziz, S., Kashif, M., & Aijaz, M. (2016). English grammar problems seen in the original articles. *Journal of the College of Physicians and Surgeons Pakistan*, 26(8), 681-684.
- Badger, R., & White, G. (2000). A process genre approach to teaching writing. *ELT journal*, 54(2), 153-160.
- Bakla, A. (2019). A mixed-methods study of tailor-made animated cartoons in teaching punctuation in EFL writing. *ReCALL*, 31(1), 75-91.

- Baldwin, C., & Chandler, G. E. (2002). Improving faculty publication output: The role of a writing coach. *Journal of Professional Nursing, 18*(1), 8-15.
- Bardi, M. (2015). Learning the practice of scholarly publication in English—A Romanian perspective. *English for Specific Purposes, 37*, 98-111.
- Bardovi-Harlig, K. (1992). The relationship of form and meaning: A cross-sectional study of tense and aspect in the interlanguage of learners of English as a second language. *Applied Psycholinguistics, 13*(3), 253-278.
- Bardovi-Harlig, K. (1995). A narrative perspective on the development of the tense/aspect system in second language acquisition. *Studies in Second Language Acquisition, 17*(2), 263-91.
- Barrett, N. E., & Chen, L. M. (2011). English article errors in Taiwanese college students' EFL writing. In *International Journal of Computational Linguistics & Chinese Language Processing, 16*(3-4), September/December 2011.
- Basturkmen, H. (2009). Commenting on results in published research articles and masters dissertations in language teaching. *Journal of English for Academic Purposes, 8*(4), 241-251.
- Basturkmen, H. (2012). A genre-based investigation of discussion sections of research articles in dentistry and disciplinary variation. *Journal of English for Academic Purposes, 11*(2), 134-144.
- Bazerman, C. (1985). Physicists reading physics: Schema-laden purposes and purpose-laden schema. *Written communication, 2*(1), 3-23.
- Bazerman, C. (1988). *Shaping written knowledge: The genre and activity of the experimental article in science*. Madison, WI: University of Wisconsin.
- Bazerman, C. (1992). The generic performance of ownership: The patent claim and grant. *Paper presented in Re-thinking Genre Seminar, Carleton University, Ottawa*.
- Belcher, D. (1994). The apprenticeship approach to advanced academic literacy: Graduate students and their mentors. *English for specific purposes, 13*(1), 23-34.
- Belcher, D. D. (2007). Seeking acceptance in an English-only research world. *Journal of Second Language Writing, 16*(1), 1-22.

- Bennett, K. (2010). Academic discourse in Portugal: A whole different ballgame?. *Journal of English for Academic Purposes*, 9(1), 21-32.
- Bennett, K. (2011). *Academic writing in Portugal: I-discourses in conflict*. Imprensa da Universidade de Coimbra/Coimbra University.
- Benzer, A. (2010). Prospective teachers' proficiency in punctuation rules and opinions related to punctuation problems. *Procedia-Social and Behavioral Sciences*, 2(2), 1878-1883.
- Berkenkotter, C. , Huckin, T. N. , & Ackerman, J. (1991). Social contexts and socially constructed texts: The initiation of a graduate student into a writing research community. In C. Bazerman & J. Paradis (Eds.), *Textual dynamics of the professions: Historical and contemporary studies of writing in academic and other professional communities* (pp. 191-215). Madison: University of Wisconsin.
- Berkenkotter, C., & Huckin, T. N. (1995). *Genre knowledge in disciplinary communication: Cognition/culture/power*. Hillsdale, NJ: Lawrence Erlbaum.
- Bestgen, Y., & Granger, S. (2011). Categorising spelling errors to assess L2 writing. *International Journal of Continuing Engineering Education and Life Long Learning*, 21(2-3), 235-252.
- Bhatia, V. K. (1998). Generic conflicts in academic discourse. In I. Fortanet, S. Posteguillo, J.C. Palmer, & J. F. Coll (Eds.), *Genre studies in English for academic purposes* (pp. 15-28). Castellon, Spain: Universitat Jaume I.
- Bhatia, V. K. (2004). *Worlds of written discourse: A genre-based view*. London: Continuum.
- Bhatia, V. K. (2008). Genre analysis, ESP and professional practice. *English for Specific Purposes*, 27(1), 161-174.
- Biber, D. (1991). *Variation across speech and writing*. Cambridge: Cambridge University.
- Biber, D. (2006). *University language: A corpus-based study of spoken and written registers* (Vol. 23). Amsterdam: John Benjamins.
- Biber, D. E. (2012). Corpus-based and corpus-driven analyses of language variation and use. In *The Oxford handbook of linguistic analysis*. Oxford University.
- Biber, D., & Conrad, S. (2009). *Register, genre, and style*. Cambridge: Cambridge University.

- Biber, D., & Conrad, S. (2014). *Variation in English: Multi-dimensional studies*. New York: Routledge.
- Biber, D., & Gray, B. (2010). Challenging stereotypes about academic writing: Complexity, elaboration, explicitness. *Journal of English for Academic Purposes*, 9(1), 2-20.
- Biber, D., & Gray, B. (2013). Discourse characteristics of writing and speaking task types on the TOEFL IBT® test: A lexico-grammatical analysis. *ETS Research Report Series*, 2013(1), i-128.
- Biber, D., Conrad, S., & Leech, G. (2002). *Student grammar of spoken and written English*. London, England: Longman.
- Biber, D., Douglas, B., Conrad, S., & Reppen, R. (1998). *Corpus linguistics: Investigating language structure and use*. Cambridge: Cambridge University.
- Biber, D., Gray, B., & Poonpon, K. (2011). Should we use characteristics of conversation to measure grammatical complexity in L2 writing development? *Tesol Quarterly*, 45, 5-35.
- Biber, D., Gray, B., & Staples, S. (2014). Predicting patterns of grammatical complexity across language exam task types and proficiency levels. *Applied Linguistics*, 37(5), 639-668.
- Bocanegra-Valle, A. (2014). 'English is my default academic language': Voices from LSP scholars publishing in a multilingual journal. *Journal of English for Academic Purposes*, 13, 65-77.
- Boersma, P., & Weenink, D. (2014). *Praat: Doing phonetics by computer*. Retrieved from <http://www.fon.hum.uva.nl/praat>.
- Bond, R. (2016). *Main grammatical writing errors of Chinese undergraduate students in UK universities*. (Unpublished master's thesis). Leiden University, The Netherlands.
- Bondi, M. (2012). Voice in textbooks: Between exposition and argument. In K. Hyland & C. Sancho Guinda (Eds.), *Stance and voice in written academic genres* (pp. 101–117). London: Palgrave.
- Bondi, M., & Scott, M. (Eds.). (2010). *Keyness in texts* (Vol. 41). John Benjamins.

- Bondi, M., Silver, M. (2004). Textual voices: A cross-disciplinary study of attribution in academic discourse. In Anderson, L., Bamford, J. (Eds.), *Evaluation in spoken and written academic discourse* (pp. 121-141). Roma, Italy: Officina.
- Bordage, G. (2001). Reasons reviewers reject and accept manuscripts: the strengths and weaknesses in medical education reports. *Academic Medicine*, 76(9), 889-896.
- Borglin, G., & Fagerström, C. (2012). Nursing students' understanding of critical thinking and appraisal and academic writing: a descriptive, qualitative study. *Nurse Education in Practice*, 12(6), 356-360.
- Botley, S. P., & Dillah, D. (2007). Investigating spelling errors in a Malaysian learner corpus. *Malaysian Journal of ELT Research*, 3, 74–93.
- Bourdieu, P. (1986). The forms of capital. In Richardson, J. (Ed), *Handbook of Theory and Research for the Sociology of Education* (pp. 241–258). Westport, CT: Greenwood.
- Brazil, D. (1995). *A grammar of speech*. USA: Oxford University.
- Brezina, V. (2018). *Statistics in corpus linguistics: A practical guide*. Cambridge: Cambridge University.
- Brown, D. (2001). *Teaching by principles: An interactive approach to language pedagogy*. Addison Wesley Longman.
- Brown, D., & Yule, G. (1983). *Teaching the spoken language*. Cambridge, UK: Cambridge University.
- Bruce, I. (2009). Results sections in sociology and organic chemistry articles: A genre analysis. *English for Specific Purposes*, 28(2), 105-124.
- Bruffee, K. A. (1986). Social construction, language, and the authority of knowledge: A bibliographical essay. *College English*, 48(8), 773-790.
- Buckingham, L., & Aktuğ-Ekinci, D. (2017). Interpreting coded feedback on writing: Turkish EFL students' approaches to revision. *Journal of English for Academic Purposes*, 26, 1-16.
- Bui, G. (2018). A lexical approach to teaching formality in freshman L2 academic writing. In L. T. Wong & W. L. Wong (Eds.), *Teaching and learning English for academic purposes: Current research and practices* (pp.111-124). New York, NY: Nova Science.

- Bulté, B., & Housen, A. (2012). Defining and operationalising L2 complexity. In Housen, A., Kuiken, F. & Vedder, I. (Eds), *Dimensions of L2 performance and proficiency: Complexity, accuracy and fluency in SLA*, (pp.23-46). Amsterdam: John Benjamins.
- Bulté, B., & Housen, A. (2014). Conceptualizing and measuring short-term changes in L2 writing complexity. *Journal of second language writing*, 26, 42-65.
- Bulté, B., & Housen, A. (2018). Syntactic complexity in L2 writing: Individual pathways and emerging group trends. *International Journal of Applied Linguistics*, 28(1), 147-164.
- Bunton, D. (2014). Generic moves in Ph. D. thesis introductions. In J. Flowerdew (Ed), *Academic discourse* (pp. 67-85). Routledge.
- Burgess, S. (2002). Packed houses and intimate gatherings: Audience and rhetorical structure. In J. Flowerdew (Ed.), *Academic Discourse* (pp. 196 – 215). London: Pearson Education.
- Burgess, S., Gea-Valor, M. L., Moreno, A. I., & Rey-Rocha, J. (2014). Affordances and constraints on research publication: A comparative study of the language choices of Spanish historians and psychologists. *Journal of English for Academic Purposes*, 14, 72-83.
- Buttery, P., & Caines, A. (2012). Normalising frequency counts to account for ‘opportunity of use’ in learner corpora. In Y. Tono, Y. Kawaguchi, & M. Minegishi (Eds.), *Developmental and cross-linguistic perspectives in learner corpus research* (pp. 187–204). Amsterdam, the Netherlands: Benjamins.
- Callies, M. (2009). ‘What is even more alarming is...’ - A contrastive learner-corpus study of what-clefts in advanced German and Polish L2 writing. In M. Wysocka (Ed.) *On language structure, acquisition and teaching. Studies in honour of Janusz Arabski on the occasion of his 70th birthday* (pp. 283-292). Katowice: Wydawnictwo Uniwersytetu Slaskiego.
- Canagarajah, A. S. (2003). A somewhat legitimate and very peripheral participation. In C. P. Casanave & S.Vandrick (Eds.), *Writing for scholarly publication: Behind the scenes in language education* (pp. 197–210). Mahwah, NJ: Lawrence Erlbaum Associates.

- Canagarajah, S. A. (1996). "Nondiscursive" requirements in academic publishing, material resources of periphery scholars, and the politics of knowledge production. *Written Communication*, 13(4), 435–472.
- Cargill, M., & Burgess, S. (2008). Introduction to the special issue: English for research publication purposes. *Journal of English for Academic Purposes*, 2(7), 75-76.
- Cargill, O., Charvat, W., & Walsh, D. (1966). The publication of academic writing. *PMLA*, 81(4), 39–45. <https://doi.org/10.2307/1261170>
- Carneiro, I. A. (2016). *The information age*. Clube de Autores.
- Carrió-Pastor, M. L., & Mestre-Mestre, E. M. (2014). Lexical errors in second language scientific writing: Some conceptual implications. *International journal of English studies*, 14(1), 97-108.
- Carter, R., & Adolphs, S. (2008). Linking the verbal and visual: new directions for corpus linguistics. *Language and computers studies in practical linguistics*, 64, 275.
- Casanave, C. P. (1994). Language development in students' journals. *Journal of Second Language Writing*, 3(3), 179-201.
- Casanave, C. P. (1998). Transitions: The balancing act of bilingual academics. *Journal of Second Language Writing*, 7(2), 175-203.
- Celce-Murcia, M. (1991). Grammar pedagogy in second and foreign language teaching. *TESOL quarterly*, 25(3), 459-480.
- Chan, A. Y. (2010). Toward a taxonomy of written errors: Investigation into the written errors of Hong Kong Cantonese ESL learners. *TESOL Quarterly*, 44(2), 295-319.
- Chang, P., & Schleppegrell, M. (2011). Taking an effective authorial stance in academic writing: Making the linguistic resources explicit for L2 writers in the social sciences. *Journal of English for academic purposes*, 10(3), 140-151.
- Charles, M. (2003). 'This mystery...': a corpus-based study of the use of nouns to construct stance in theses from two contrasting disciplines. *Journal of English for Academic Purposes*, 2(4), 313-326.
- Chartrand, R. (2016). *Extraction and analysis of modal auxiliaries in consecutive clauses from a corpus*. Cambridge Scholars.

- Chaudron, C., & Parker, K. (1990). Discourse markedness and structural markedness: The acquisition of English noun phrases. *Studies in Second Language Acquisition*, 12(1), 43-64.
- Chen, B., Deng, D., Zhong, Z., & Zhang, C. (2020). Exploring linguistic characteristics of highly browsed and downloaded academic articles. *Scientometrics*, 122(3), 1769-1790.
- Chen, Y. H., & Baker, P. (2010). Lexical bundles in L1 and L2 academic writing. *Language learning & technology*, 14(2), 30-49.
- Chen, Y. H., & Baker, P. (2016). Investigating criterial discourse features across second language development: Lexical bundles in rated learner essays, CEFR B1, B2 and C1. *Applied Linguistics*, 37(6), 849-880.
- Cheng, W. (2011). *Exploring corpus linguistics: Language in action*. New York: Routledge.
- Cho, D. W. (2009). Science journal paper writing in an EFL context: The case of Korea. *English for Specific Purposes*, 28(4), 230-239.
- Cho, S. (2004). Challenges of entering discourse communities through publishing in English: Perspectives of nonnative-speaking doctoral students in the United States of America. *Journal of Language, Identity, and Education*, 3(1), 47-72.
- Cho, S. (2005). *International graduate students in us-based tesol discourse communities: finding and creating a space*. (Unpublished doctoral dissertation). University of Tennessee, Knoxville.
- Cho, S. (2009). Disciplinary enculturation experiences of five East Asian doctoral students in US-based second language studies programmes. *Asia Pacific Journal of Education*, 29(3), 295-310.
- Chuang, F. Y., & Nesi, H. (2006). An analysis of formal errors in a corpus of L2 English produced by Chinese students. *Corpora*, 1(2), 251-271.
- Clarivate Analytics. (2016) *Journal selection process*. (n.d.). Retrieved from <https://clarivate.com/essays/journal-selection-process/>
- Cohen, J. (1988). *Statistical power analyses for the behavioral sciences (2nd ed.)*. Hillsdale, NJ: Erlbaum.

- Collins, P. (2004). Reversed what-clefts in English. *Australian Review of Applied Linguistics*, 27(2), 63-74.
- Connor, U. (1996). *Contrastive rhetoric: Cross-cultural aspects of second-language writing*. New York: Cambridge University.
- Connor, U., & Upton, T. (2004). The genre of grant proposals: A corpus linguistic analysis. In U. Connor, & T. Upton (Eds.), *Discourse in the professions: Perspectives from corpus linguistics* (pp. 235-255). Amsterdam: John Benjamins.
- Cooper, T. C. (1976). Measuring written syntactic patterns of second language learners of German. *The Journal of Educational Research*, 69(5), 176-183.
- Cortes, V. (2013). The purpose of this study is to: Connecting lexical bundles and moves in research article introductions. *Journal of English for Academic Purposes*, 12(1), 33-43.
- Covington, M. A., & McFall, J. D. (2010). Cutting the Gordian knot: The moving-average type-token ratio (MATTR). *Journal of Quantitative Linguistics*, 17(2), 94-100.
- Crompton, P. (2011). Article errors in the English writing of advanced L1 Arabic learners: The role of transfer. *Asian EFL Journal*, 50(1), 4-35.
- Crossley, S. A., & McNamara, D. S. (2014). Does writing development equal writing quality? A computational investigation of syntactic complexity in L2 learners. *Journal of Second Language Writing*, 26, 66-79.
- Crossley, S. A., Allen, L. K., Kyle, K., & McNamara, D. S. (2014). Analyzing discourse processing using a simple natural language processing tool. *Discourse Processes*, 51(5-6), 511-534.
- Cumming, A. (1989). Writing expertise and second language proficiency. *Language Learning* 39(1), 81-141.
- Cumming, A., Kantor, R., Baba, K., Erdosy, U., Eouanzoui, K., & James, M. (2005). Differences in written discourse in independent and integrated prototype tasks for next generation TOEFL. *Assessing Writing*, 10(1), 5-43.
- Curry, M. J., & Lillis, T. (2004). Multilingual scholars and the imperative to publish in English: Negotiating interests, demands, and rewards. *TESOL quarterly*, 38(4), 663-688.

- Curry, M. J., & Lillis, T. M. (2010). Academic research networks: Accessing resources for English-medium publishing. *English for Specific Purposes*, 29(4), 281-295.
- Çandarlı, D. (2012). A cross-cultural investigation of English and Turkish research article abstracts in educational sciences. *Studies About Languages*, (20), 12-17.
- Dagneaux, E., Denness, S., & Granger, S. (1998). Computer-aided error analysis. *System*, 26(2), 163-174.
- Dagneaux, E., Denness, S., Granger, S., & Meunier, F. (1996). *Error tagging manual version 1.1. Louvain-la-neuve: Centre for English corpus linguistics*, Universite catholique de Louvain.
- Dahl, T. (2004a). Pragmatics of discourse. *Journal of Pragmatics*, 36, 1807-1825.
- Dahl, T. (2004b). Textual metadiscourse in research articles: a marker of national culture or of academic discipline?. *Journal of pragmatics*, 36(10), 1807-1825.
- Das, D., & Taboada, M. (2018). Signalling of coherence relations in discourse, beyond discourse markers. *Discourse Processes*, 55(8), 743-770.
- Dawson, J. H. (1992). Avoid redundancy in writing. *Weed Technology*, 6(3), 782-782.
- De Cock, S. (2002). Pragmatic prefabs in learners dictionaries. In *Proceedings of the Tenth EURALEX International Congress* (pp. 13-17), Copenhagen, Denmark, August.
- Demir, S. B. (2018). Predatory journals: Who publishes in them and why?. *Journal of Informetrics*, 12(4), 1296-1311.
- Demol, A., & Hadermann, P. (2008). An exploratory study of discourse organisation in French L1, Dutch L1, French L2 and Dutch L2 written narratives. In G. Gilquin, S. Papp, & Díez-Bedmar María (Eds.), *Linking up contrastive and learner corpus research* (pp. 255-282). Brill Rodopi.
- Dhia, A. (2006). *The information age and diplomacy: An emerging strategic vision in world affairs*. Florida: Dissertation.
- Díez-Bedmar M. B. & Papp, S. (2008). The use of the English article system by Chinese and Spanish learners. In G. Gilquin, S. Papp, & Díez-Bedmar María Belén (Eds.), *Linking up contrastive and learner corpus research* (pp. 147–175). Amsterdam & Atlanta: Rodopi.

- Ding, H. (2007). Genre analysis of personal statements: Analysis of moves in application essays to medical and dental schools. *English for Specific Purposes*, 26(3), 368-392.
- Dogancay-Aktuna, S. (1998). The spread of English in Turkey and its current sociolinguistic profile. *Journal of Multilingual and Multicultural Development*, 19(1), 24-39.
- Donato, R. (1994). Collective scaffolding in second language learning. In J. P. Lantolf, & G. Appel (Eds.), *Vygostkian approaches to second language research* (pp. 33-56). New Jersey: Ablex.
- Dos Santos, M. B. (1996). The textual organization of research paper abstracts in applied linguistics. *Text-Interdisciplinary Journal for the Study of Discourse*, 16(4), 481-500.
- Drott, M. C. (1995). Reexamining the role of conference papers in scholarly communication. *Journal of the American Society for Information Science*, 46(4), 299-305.
- Dudley-Evans, T. (2001). English for specific purposes. In R. Carter & D. Nunan (Eds.). *Teaching English to speakers of other languages* (131-136). Cambridge: Cambridge University.
- Dudley-Evans, T. (1994). Genre analysis: An approach to text analysis for ESP. In M. Coulthard (Ed.), *Advances in written text analysis* (pp. 219 – 228). London: Routledge.
- Dudley-Evans, T. (2002). The teaching of the academic essay: Is a genre approach possible? In Johns, A.M. (Ed.), *Genre in the classroom: Multiple perspectives* (pp. 225-235). London: Lawrence Erlbaum.
- Dudley-Evans, T., & St John, M. J. (1998). *Developments in ESP: A multi-disciplinary approach* (No. 428.007 D849). Cambridge University.
- Dueñas, P. M. (2007). ‘I/we focus on...’: A cross-cultural analysis of self-mentions in business management research articles. *Journal of English for Academic Purposes*, 6(2), 143-162.
- Durgun, Ş. (2010). *Modernleşme ve siyaset*. Ankara: A Kitap.
- Duszak, A., & Lewkowicz, J. (2008). Publishing academic texts in English: A Polish perspective. *Journal of English for Academic Purposes*, 7(2), 108-120.

- Ehara, S., & Takahashi, K. (2007). Reasons for rejection of manuscripts submitted to AJR by international authors. *American Journal of Roentgenology*, 188(2), W113-W116.
- El Malik, A. T., & Nesi, H. (2008). Publishing research in a second language: the case of Sudanese contributors to international medical journals. *Journal of English for Academic Purposes*, 7(2), 87-96.
- Elkilic, G., Han, T., & Aydin, S. (2009). Punctuation and capitalisation errors of Turkish EFL students in composition classes: An evidence of L1 interference. In *1st International Symposium on Sustainable Development, Sarajevo, Bosnia and Herzegovina* (pp. 270-284).
- Ellis, R. (1994). A theory of instructed second language acquisition. *Implicit and Explicit Learning of Languages*, 79-114.
- Ellis, R., & Yuan, F. (2004). The effects of planning on fluency, complexity, and accuracy in second language narrative writing. *Studies in Second Language Acquisition*, 26(1), 59-84.
- Englander, K. (2009). Transformation of the identities of nonnative English-speaking scientists as a consequence of the social construction of revision. *Journal of Language, Identity, and Education*, 8(1), 35-53.
- Englander, K., & Uzuner-Smith, S. (2013). The role of policy in constructing the peripheral scientist in the era of globalization. *Language Policy*, 12(3), 231-250.
- Erdemir, N. (2019). *Corpus-based rhetorical move analysis of the results sections in L1 English, L2 English and L1 Turkish research articles*. (Unpublished doctoral dissertation). Gazi University Graduate School of Educational Sciences, Ankara.
- Farwaneh, S., & Tamimi, M. (2012). *Arabic learners written corpus: A resource for research and learning*. Retrieved from the the University of Arizona, the Center for Educational Resources in Culture, Language and Literacy web site: <http://l2arabiccorpus.cercll.arizona.edu/?q=homepage>
- Ferguson, G., Pérez-Llantada, C., & Plo, R. (2011). English as an international language of scientific publication: A study of attitudes. *World Englishes*, 30(1), 41-59.
- Ferris, D. R. (2006). Does error feedback help student writers? New evidence on the short- and long-term effects of written error correction. *Feedback in Second Language Writing: Contexts and Issues*, 81-104.

- Ferris, D., & Hedgcock, J. S. (2005). Teacher response to student writing: Issues in oral and written feedback. In D. Ferris & J. S. Hedgcock (Eds.) *Teaching ESL composition: Purpose, process and practice* (pp. 184-222) Mahwah, NJ: Lawrence Erlbaum.
- Field, A. (2009). Non-parametric tests. *Discovering statistics using SPSS*, 2, 521-570.
- Fischer, R. A. (1984). Testing written communicative competence in French. *The Modern Language Journal*, 68(1), 13-20.
- Fisher R.A. (1992) Statistical methods for research workers. In Kotz S. & Johnson N.L. (Eds.), *Breakthroughs in Statistics. Springer series in statistics (Perspectives in Statistics)*. New York, NY: Springer. https://doi.org/10.1007/978-1-4612-4380-9_6
- Flowerdew, J. (1993). An educational, or process, approach to the teaching of professional genres. *ELT journal*, 47(4), 305-316.
- Flowerdew, J. (1999a). Problems in writing for scholarly publication in English: The case of Hong Kong. *Journal of Second Language Writing*, 8(3), 243-264.
- Flowerdew, J. (1999b). Writing for scholarly publication in English: The case of Hong Kong. *Journal of Second Language Writing*, 8(2), 123-145.
- Flowerdew, J. (2000). Discourse community, legitimate peripheral participation, and the nonnative-English-speaking scholar. *TESOL quarterly*, 34(1), 127-150.
- Flowerdew, J. (2001). Attitudes of journal editors to nonnative speaker contributions. *TESOL Quarterly*, 35(1), 121-150.
- Flowerdew, J. (2008). Scholarly writers who use English as an additional language: what can Goffman's "stigma" tell us?. *Journal of English for Academic Purposes*, 7(2), 77-86.
- Flowerdew, J. (2013). English for research publication purposes. In B. Paltridge & S. Starfield (Eds.), *The handbook of English for specific purposes* (pp. 301-321). Malden, MA: Wiley- Blackwell.
- Flowerdew, J. (2014). *Academic discourse*. New York: Routledge.
- Flowerdew, J. (2015). Some thoughts on English for research publication purposes (ERPP) and related issues. *Language Teaching*, 48(2), 250-262.
- Flowerdew, J., & Wang, S. H. (2015). Identity in academic discourse. *Annual Review of Applied Linguistics*, 35, 81-99.

- Foster, P., & Skehan, P. (1996). The influence of planning and task type on second language performance. *Studies in Second Language Acquisition*, 18(3), 299-323.
- Foster, P., & Tavakoli, P. (2009). Native speakers and task performance: Comparing effects on complexity, fluency, and lexical diversity. *Language Learning*, 59(4), 866-896.
- Francis, G. (1993). A corpus-driven approach to grammar: Principles, methods and examples. *Text and technology: In honour of John Sinclair*, 1, 137-156.
- Frantzen, D. (1995). The effects of grammar supplementation on written accuracy in an intermediate Spanish content course. *Modern Language Journal*, 79, 329–344.
- Freedman, A., & Medway, P. (Eds.). (1994). *Genre and the new rhetoric*. London: Taylor & Francis.
- Fukao, A., & Fujii, T. (2001). Investigating difficulties in the academic writing process: Interview as a research tool. *ICU Language Research Bulletin*, 16, 29-40.
- Fuoli, M. (2018). Building a trustworthy corporate identity: A corpus-based analysis of stance in annual and corporate social responsibility reports. *Applied Linguistics*, 39(6), 846-885.
- Gass, S. M. (2013). *Second language acquisition: An introductory course*. New York: Routledge.
- Gauthier, L. R. (1993). A strategy to increase punctuation awareness. *Journal of Reading*, 36(5), 401-402.
- Gee, J. P. (2014). *An introduction to discourse analysis: Theory and method*. New York: Routledge.
- George, D., & Mallery, P. (2012). *SPSS statistics 21: Step by step*. Boston: Allynand Bacon.
- Gilquin, G. & Granger, S. (2015). Learner language. In: Douglas Biber & Randi Reppen (Eds.), *The Cambridge handbook of English corpus linguistics* (pp. 418-435), Cambridge: Cambridge University.
- Gilquin, G. (2000). The integrated contrastive model: Spicing up your data. *Languages in Contrast*, 3(1), 95-123.
- Gilquin, G. (2008). Combining contrastive and interlanguage analysis to apprehend transfer: detection, explanation, evaluation. *Language and Computers Studies in Practical Linguistics*, 66, 3.

- Gilquin, G. (2012). Lexical infelicity in English causative constructions. Comparing native and learner collocations. In Leino, J. (Ed), *Analytical causatives. From 'give' and 'come' to 'let' and 'make'* (pp.41-63). München: Lincom Europa.
- Gilquin, G. (2015). Contrastive collocation analysis: Causative constructions in English and French. *Zeitschrift für Anglistik und Amerikanistik*, 63(3), 253-272.
- Gilquin, G. (2016). Input-dependent L2 acquisition: Causative constructions in English as a foreign and second language. In De Knop, S. & Gilquin, G. (Eds.), *Applied construction grammar* (pp. 115–148). Berlin, Boston: De Gruyter Mouton. doi: <https://doi.org/10.1515/9783110458268>
- Gilquin, G., & Paquot, M. (2008). Too chatty: Learner academic writing and register variation. *English Text Construction*, 1(1), 41-61.
- Gilquin, G., Granger, S., & Paquot, M. (2007). Learner corpora: The missing link in EAP pedagogy. *Journal of English for Academic Purposes*, 6(4), 319-335.
- Gilquin, G., Papp, S., & Dez-Bedmar, M. B. (Eds.). (2008). *Linking up contrastive and learner corpus research* (No. 66). Amsterdam: Rodopi.
- Glasman-Deal, H. (2010). *Science research writing for non-native speakers of English*. London: World Scientific.
- Gonzales, L. D., Martinez, E., & Ordu, C. (2014). Exploring faculty experiences in a striving university through the lens of academic capitalism. *Studies in Higher Education*, 39(7), 1097-1115.
- Gorrell, D. K. (1981). *Controlled composition for teaching basic writing to college freshmen: A comparison with grammar lessons*. Doctoral dissertation, Illinois: Illinois State University.
- Gosden, H. (1992). Research writing and NNSs: From the editors. *Journal of Second Language Writing*, 1(2), 123-139.
- Gosden, H. (1995). Success in research article writing and revision: A social-constructionist perspective. *English for Specific Purposes*, 14(1), 37-57.
- Gosden, H. (1996). Verbal reports of Japanese novices' research writing practices in English. *Journal of Second Language Writing*, 5(2), 109-128.

- Gosden, H. (2003). 'Why not give us the full story?': Functions of referees' comments in peer reviews of scientific research papers. *Journal of English for Academic Purposes*, 2(2), 87-101.
- Gould, S. J. (1995). Ladders and cones: Constraining evolution by canonical icons. In R.B. Silvers (Ed.), *Hidden histories of science* (pp. 37-67). New York: A New York Review.
- Grabe, W. (1988). English, information access, and technology transfer: A rationale for English as an international language. *World Englishes*, 7(1), 63-72.
- Grabowski, Ł. (2015). Keywords and lexical bundles within English pharmaceutical discourse: A corpus-driven description. *English for Specific Purposes*, 38, 23-33.
- Granger S. (2017b) Learner corpora in foreign language education. In Thorne. S. & May S. (Eds), *Language, education and technology. Encyclopedia of language and education (3rd ed.)*. New York: Springer.
- Granger, S. (1993). International corpus of learner English. In J. Aarts, P. de Haan & N. Oostdijk (Eds.), *English language corpora: Design, analysis and exploitation* (pp. 57-69). Amsterdam: Rodopi.
- Granger, S. (1996). "From CA to CIA and back: An integrated contrastive approach to computerized bilingual and learner corpora". In K. Aijmer, B. Altenberg & M. Johansson (Eds.), *Languages in contrast. Text-based cross-linguistic studies* (pp. 37-51). Lund: Lund University.
- Granger, S. (1997). On identifying the syntactic and discourse features of participle clauses in academic English: Native and non-native writers compared. In: Aarts, J., I. de Moënnink, E., & Wekker, H. (Eds.), *Studies in English Language Research and Teaching* (pp. 185–198). Amsterdam: Rodopi.
- Granger, S. (1999). Use of tenses by advanced EFL learners: Evidence from an error-tagged computer corpus. In H. Hasselgard & S. Oksefjell (Eds.), *Out of corpora. Studies in honour of Stig Johansson* (pp. 191-202). Amsterdam: Rodopi.
- Granger, S. (2002). A bird's eye view of learner corpus research. In S. Granger, J. Hung, & S. Petch-Tyson (Eds.), *Computer learner corpora, second language acquisition and foreign Language teaching* (pp. 3-33). Philadelphia, PA: John Benjamins.

- Granger, S. (2003). Error-tagged learner corpora and CALL: A promising synergy. *CALICO Journal*, 20, 465-480.
- Granger, S. (2004). Computer learner corpus research: Current status and future prospects. In *Applied Corpus Linguistics* (pp. 123-145). Amsterdam: Rodopi.
- Granger, S. (2007). The computer learner corpus: A versatile new source of data for SLA research. In: T. W. Krishnamurthy (Ed.), *Corpus linguistics: Critical concepts in linguistics* (pp. 166-182). London: Routledge.
- Granger, S. (2015). Contrastive interlanguage analysis: A reappraisal. *International Journal of Learner Corpus Research*, 1(1), 7-24.
- Granger, S. (2017). Academic phraseology: A key ingredient in successful L2 academic literacy. *Oslo Studies in Language*, 9(3), 9-27
- Granger, S., & Bestgen, Y. (2014). The use of collocations by intermediate vs. advanced non-native writers: A bigram-based study. *International Review of Applied Linguistics in Language Teaching*, 52(3), 229-252.
- Granger, S., & Leech, G. (Eds.). (2014). *Learner English on computer*. New York: Routledge.
- Granger, S., & Meunier, F. (Eds.). (2008). *Phraseology: An interdisciplinary perspective*. Amsterdam: John Benjamins.
- Granger, S., & Tyson, S. (1996). Connector usage in the English essay writing of native and non-native EFL speakers of English. *World Englishes*, 15(1), 17-27.
- Granger, S., Dagneaux, E., Meunier, F., & Paquot, M. (2009). *ICLE: International corpus of learner English*. Louvain-la-Neuve, Belgium: Presses Universitaires de Louvain.
- Gray, B. (2010). On the use of demonstrative pronouns and determiners as cohesive devices: A focus on sentence-initial this/these in academic prose. *Journal of English for Academic Purposes*, 9(3), 167-183.
- Gray, B., & Biber, D. (2015). Stance markers. In Aijmer, K., & Rühlemann, C. (Eds.). *Corpus pragmatics*. Cambridge: Cambridge University.
- Gray, B., & Cortes, V. (2011). Perception vs. evidence: An analysis of this and these in academic prose. *English for Specific Purposes*, 30(1), 31-43.

- Gries, S.T. (2013). Statistical tests for the analysis of learner corpus data. In Díaz-Negrillo, A., Ballier, N. & Thompson, P. (Eds.), *Automatic treatment and analysis of learner corpus data* (pp. 287-309). Amsterdam: John Benjamins.
- Grosz, P. G. (2016). Bridging uses of demonstrative pronouns in German. *Linguistics and Philosophy*, 41, 367-421.
- Güngör, F., & Uysal, H. H. (2016). A comparative analysis of lexical bundles used by native and non-native scholars. *English Language Teaching*, 9(6), 176-188.
- Hair, J. F., Ringle, C. M., & Sarstedt, M. (2013). Partial least squares structural equation modeling: Rigorous applications, better results and higher acceptance. *Long Range Planning*, 46(1-2), 1-12.
- Halliday, M. A. K. (1961). Categories of the theory of grammar. *Word*, 17(2), 241-292.
- Halliday, M. A. K. (1994). *An introduction to functional grammar* (2nd ed.). London: Edward Arnold.
- Halliday, M. A. K., & R. Hasan. (1976). *Cohesion in English*. London: Longman.
- Halliday, M.A.K. & J.R. Martin. (1993). *Writing science: Literacy and discourse power*. London: Falmer.
- Hamed, M. (2018). Common linguistic errors among non-English major Libyan students writing. *Arab World English Journal (AWEJ)*, 9(3), 219-232.
- Hamel, R. E. (2007). The dominance of English in the international scientific periodical literature and the future of language use in science. *Aila Review*, 20(1), 53-71.
- Hammarberg, B. (2010). *Introduction to the ASU Corpus : a longitudinal oral and written text corpus of adult learner Swedish with a corresponding part from native Swedes*. Retrieved from <http://urn.kb.se/resolve?urn=urn:nbn:se:su:diva-112130>
- Hammond, J. (1987). An overview of the genre-based approach to the teaching of writing in Australia. *Australian Review of Applied Linguistics*, 10(2), 163-181.
- Hamp-Lyons, L., & Henning, G. (1991). Communicative writing profiles: An investigation of the transferability of a multiple-trait scoring instrument across ESL writing assessment contexts. *Language Learning*, 41(3), 337-373.
- Han, N. R., Chodorow, M., & Leacock, C. (2006). Detecting errors in English article usage by non-native speakers. *Natural Language Engineering*, 12(2), 115.

- Harley, B. (1980). Interlanguage units and their relations. *Interlanguage Studies Bulletin*, 5, 3-30.
- Hartig, A. J., & Lu, X. (2014). Plain English and legal writing: Comparing expert and novice writers. *English for Specific Purposes*, 33, 87-96.
- Harwood, N. (2005). 'We do not seem to have a theory... The theory I present here attempts to fill this gap': Inclusive and exclusive pronouns in academic writing. *Applied Linguistics*, 26(3), 343-375.
- Hasselgård, H. (2009). Thematic choice and expressions of stance in English argumentative texts by Norwegian learners. *Corpora and Language Teaching*, 33, 121-139.
- Hasselgard, H., & Iohansson, S. (2011). Learner corpora and contrastive interlanguage analysis. In Meunier F., De Cock S., Gilquin G. & Paquot M. (Eds), *A taste for corpora. In honour of Sylviane Granger* (pp. 33-62). Amsterdam: Benjamins.
- Hazelkorn, E. (2011). *Rankings and the reshaping of higher education: The battle for world-class excellence*. Basingstoke: Palgrave Macmillan.
- Hedgcock, J., & Lefkowitz, N. (1992). Collaborative oral/aural revision in foreign language writing instruction. *Journal of Second Language Writing*, 1(3), 255-276.
- Hellwig, B. (2019). *ELAN - linguistic annotator*. Retrieved from <https://www.mpi.nl/corpus/html/elan/>
- Hellwig, B., & Van Uytvanck, D. (2005). *Eudico linguistic annotator (ELAN) version 2.4 manual*. Technical report, MaxPlanck Institute for Psycholinguistics, Nijmegen, The Netherlands.
- Henry, A., & Roseberry, R. L. (2007). Language errors in the genre-based writing of advanced academic ESL students. *RELC Journal*, 38(2), 171-198.
- Herbel-Eisenmann, B. A. (2007). From intended curriculum to written curriculum: Examining the "voice" of a mathematics textbook. *Journal for Research in Mathematics Education*, 38(4), 344-369.
- Heuboeck, A., Holmes, J., & Nesi, H. (2008). *The BAWE corpus manual*. Retrieved from <https://www.reading.ac.uk/internal/appling/bawe/BAWE.documentation.pdf>
- Hewings, M. (2006). English language standards in academic articles: Attitudes of peer reviewers. *Revista Canaria de Estudios Ingleses*, 53, 47-62.

- Hill, S. S., Soppelsa, B. F., & West, G. K. (1982). Teaching ESL students to read and write experimental-research papers. *TESOL quarterly*, 16(3), 333-347.
- Hilton, A., & Armstrong, R. A. (2006). Statnote 6: post-hoc ANOVA tests. *Microbiologist*, 2006, 34-36.
- Hinkel, E. (2002). *Second language writers' text*. Mahwah, NJ: Lawrence Erlbaum.
- Hinkel, E. (2011). What research on second language writing tells us and what it doesn't. In E. Hinkel (Ed.), *Handbook of research in second language teaching and learning volume 2* (pp 523-538). New York: Routledge.
- Hinkel, E. (2013). Research findings on teaching grammar for academic writing. *English Teaching*, 68(4), 3-21.
- Hirvela, A., Nussbaum, A., & Pierson, H. (2012). ESL students' attitudes toward punctuation. *System*, 40(1), 11-23.
- Ho, M. C. (2017). Navigating scholarly writing and international publishing: Individual agency of Taiwanese EAL doctoral students. *Journal of English for Academic Purposes*, 27, 1-13.
- Hollander, M., & Wolfe, D. A. (1999). *Nonparametric statistical methods*. New York: John Wiley & Sons.
- Homburg, T. J. (1984). Holistic evaluation of ESL compositions: Can it be validated objectively?. *TESOL Quarterly*, 18(1), 87-107.
- Hood, S. (2010). *Appraising research: Evaluation in academic writing*. London: Palgrave Macmillan.
- Horowitz, D. M. (1986). What professors actually require: Academic tasks for the ESL classroom. *TESOL Quarterly*, 20(3), 445-462.
- Housen, A. (2002). A corpus-based study of the L2-acquisition of the English verb system. In: Granger, S., Hung, J. & Petch-Tyson, S. (Eds) *Computer learner corpora, second language acquisition and foreign language teaching* (pp. 77-116). Amsterdam/Philadelphia, PA: John Benjamins.
- Hovermale, D. J., & Martin, S. (2008). Developing an annotation scheme for ELL spelling errors. In *Proceedings of MCLC-5 (Midwest Computational Linguistics Colloquium)*. East Lansing, MI.

- Hu, G., & Cao, F. (2011). Hedging and boosting in abstracts of applied linguistics articles: A comparative study of English-and Chinese-medium journals. *Journal of Pragmatics*, 43(11), 2795-2809.
- Huang, J. C. (2010). Publishing and learning writing for publication in English: Perspectives of NNES PhD students in science. *Journal of English for Academic Purposes*, 9(1), 33-44.
- Huang, J. C. (2017). What do subject experts teach about writing research articles? An exploratory study. *Journal of English for Academic Purposes*, 25, 18-29.
- Hundt, M., Schneider, G., & Seoane, E. (2016). The use of the be-passive in academic Englishes: Local versus global usage in an international language. *Corpora*, 11(1), 29-61.
- Hunt, K. W. (1970). Recent measures in syntactic development. In M. Lester (ed.) *Reading in applied transformational grammar* (pp. 179-192). Nueva York: Holt, Rinehart and Wiston.
- Hunt, K.W. (1965). *Grammatical structures written at three grade levels*. Research Report 3. Urbana (Ill): NCTE.
- Hutchinson, T., & Waters, A. (1987). *English for specific purposes*. Cambridge: Cambridge University.
- Hyland, K. & Tse, P. (2005). Evaluative that constructions: Signaling stance in research abstracts. *Functions of Language*, 12, 39 – 64.
- Hyland, K. & Tse, P. (2007). Is there an “academic vocabulary”? *TESOL Quarterly*, 41(2), 235–254.
- Hyland, K. (1998). Persuasion and context: The pragmatics of academic metadiscourse. *Journal of Pragmatics*, 30(4), 437 – 455.
- Hyland, K. (2000). *Disciplinary discourses: Social interactions in academic writing*. Essex: Pearson Education.
- Hyland, K. (2001). Humble servants of the discipline? Self-mention in research articles. *English for Specific Purposes*, 20(3), 207-226.
- Hyland, K. (2002a). 6. Genre: Language, context, and literacy. *Annual Review of Applied Linguistics*, 22, 113-135. doi:10.1017/S0267190502000065

- Hyland, K. (2002b). Authority and invisibility: Authorial identity in academic writing. *Journal of Pragmatics*, 34(8), 1091-1112.
- Hyland, K. (2003). Genre-based pedagogies: A social response to process. *Journal of Second Language Writing*, 12(1), 17-29.
- Hyland, K. (2004a). *Disciplinary discourses: social interactions in academic writing*. Ann Arbor, MI: University of Michigan.
- Hyland, K. (2004b). Graduates' gratitude: The generic structure of dissertation acknowledgements. *English for Specific Purposes*, 23(3), 303-324.
- Hyland, K. (2006). Representing readers in writing: Student and expert practices. *Linguistics and Education*, 16, 363-377.
- Hyland, K. (2007). Applying a gloss: Exemplifying and reformulating in academic discourse. *Applied Linguistics*, 28, 266-285.
- Hyland, K. (2007a). Genre pedagogy: Language, literacy and L2 writing instruction. *Journal of Second Language Writing*, 16(3), 148-164.
- Hyland, K. (2007b). English for professional academic purposes: Writing for scholarly publication. In D. Belcher (Ed.), *English for specific purposes in theory and practice*. (pp. 17-38). Ann Arbor, MI: University of Michigan.
- Hyland, K. (2008a). Academic clusters: Text patterning in published and postgraduate writing. *International Journal of Applied Linguistics*, 18(1), 41-62.
- Hyland, K. (2008b). Genre and academic writing in the disciplines. *Language Teaching*, 41(4), 543-562.
- Hyland, K. (2014). Activity and evaluation: Reporting practices in academic writing. In J. Flowerdew (Ed), *Academic discourse* (pp. 125-140). London: Longman.
- Hyland, K. (2015a). Corpora and written academic English. In D. Biber & R. Reppen (Eds.), *The Cambridge handbook of English corpus linguistics* (pp. 292-308). Cambridge: Cambridge University. doi:10.1017/CBO9781139764377.017
- Hyland, K. (2015b). *Teaching and researching writing*. New York: Routledge.
- Hyland, K. (2016). Academic publishing and the myth of linguistic injustice. *Journal of Second Language Writing*, 31, 58-69.

- Hyland, K. L. (2013). Writing in the university: Education, knowledge and reputation. *Language Teaching*, 46(1), 53-70.
- Hyland, K., & Diani, G. (Eds.). (2009). *Academic evaluation: Review genres in university settings*. London: Palgrave Macmillan.
- Hyland, K., & Guinda, C. S. (Eds.). (2012). *Stance and voice in written academic genres*. Houndmills, UK: Palgrave Macmillan.
- Hyland, K., & Hamp-Lyons, L. (2002). EAP: Issues and directions. *Journal of English for Academic Purposes*, 1(1), 1-12.
- Hyland, K., & Jiang, F. K. (2016). "We must conclude that...": A diachronic study of academic engagement. *Journal of English for Academic Purposes*, 24, 29-42.
- Hyland, K., & Salager-Meyer, F. (2008). Scientific writing. *Annual Review of Information Science and Technology*, 42(1), 297.
- Hyland, K., & Tse, P. (2004). Metadiscourse in academic writing: A reappraisal. *Applied Linguistics*, 25(2), 156-177.
- Hyland, K., Huat, C. M., & Handford, M. (Eds.). (2012). *Corpus applications in applied linguistics*. London: Continuum.
- Hyon, S. (1996). Genre in three traditions: Implications for ESL. *TESOL Quarterly*, 30(4), 693-722.
- Irvin, L. L. (2010). What is "academic" writing?. *Writing spaces: Readings on Writing*, 1, 3-17.
- Ishikawa, S. (1995). Objective measurement of low-proficiency EFL narrative writing. *Journal of Second Language Writing*, 4(1), 51-69.
- Ivanič, R., & Camps, D. (2001). I am how I sound: Voice as self-representation in L2 writing. *Journal of Second Language Writing*, 10(1-2), 3-33.
- Izumi, E., Saiga, T., Supnithi, T., Uchimoto, K., & Isahara, H. (2003). The development of the spoken corpus of Japanese learner English and the applications in collaboration with NLP techniques. In *Proc. of the Corpus Linguistics 2003 Conference* (pp. 359-366).

- Izumi, E., Uchimoto, K., & Isahara, H. (2004). The NICT JLE Corpus: Exploiting the language learners' speech database for research and education. *International Journal of the Computer, the Internet and Management*, 12(2), 119-125.
- Jalongo, M. R. (2002). *Writing for publication: A practical guide for educators*. Norwood, MA: Christopher-Gordon.
- James C. (1998). *Errors in language learning and use. Exploring error analysis*. London & New York: Longman.
- Jaroongkhongdach, W., Todd, R. W., Keyuravong, S., & Hall, D. (2012). Differences in quality between Thai and international research articles in ELT. *Journal of English for Academic Purposes*, 11(3), 194-209.
- Jiang, F. K. (2015). Nominal stance construction in L1 and L2 students' writing. *Journal of English for Academic Purposes*, 20, 90-102.
- Jiang, F. K., & Ma, X. (2018). 'As we can see': Reader engagement in PhD candidature confirmation reports. *Journal of English for Academic Purposes*, 35, 1-15.
- Jiang, F., & Wang, F. (2018). 'This is because...': Authorial practice of (un) attending this in academic prose across disciplines. *Australian Journal of Linguistics*, 38(2), 162-182.
- Jiang, X., Borg, E., & Borg, M. (2017). Challenges and coping strategies for international publication: Perceptions of young scholars in China. *Studies in Higher Education*, 42(3), 428-444.
- Johns, A. M. (2013). The history of English for specific purposes research. In B. Paltridge & S. Starfield (Eds.), *The handbook of English for specific purposes*, (pp. 5-30). Sussex: Wiley-Blackwell.
- Johnson, A. C., Wilson, J., & Roscoe, R. D. (2017). College student perceptions of writing errors, text quality, and author characteristics. *Assessing Writing*, 34, 72-87.
- Jones, R. (2015). *Academic writing plagiarism across Europe and beyond*. Conference Proceedings. Brno: Mendelu.
- Jurafsky, D. & Martin, J.H. (2009). *Speech and language processing*. 2nd edition. New Jersey: Prentice Hall.

- Kachru, B. (1986). *The alchemy of English: The spread, functions and models of non-native Englishes*. Oxford: Pergamon.
- Kachru, B. B. (1985). Standards, codification and sociolinguistic realism: The English language in the outer circle. In Q. Randolph & H. G. Widdowson, (Eds.), *English in the World: Teaching and learning the language and literatures* (pp. 11–30). Cambridge: Cambridge University.
- Kafes, H. (2018). A genre analysis of English and Turkish research article introductions. *Novitas-ROYAL (Research on Youth and Language)*, 12(1), 66-79.
- Kameen, P. (1979). Syntactic skill and ESL writing quality. *On TESOL*, 79, 343-364.
- Kan, M. O. (2017). Türkiye’de akademik yazma alanında yapılan tezler. *İnsan ve Toplum Bilimleri Araştırmaları Dergisi*, 6(2), 1037-1048.
- Kang, O., & Wang, L. (2014). Impact of different task types on candidates’ speaking performances and interactive features that distinguish between CEFR levels. *Research Notes*, 57, 40-49.
- Kanoksilapathan, B. (2007). Rhetorical moves in biochemistry research articles. In D. Biber, U. Connor, & T. A. Upton (Eds.). *Discourse on the move: Using corpus analysis to describe discourse structure. Studies in corpus linguistics* (pp. 73–119). Amsterdam: John Benjamins.
- Kaplan, R. B. (1966). Cultural thought patterns in inter-cultural education. *Language Learning*, 16(1-2), 1-20.
- Kaplan, R. B., & Baldauf Jr, R. B. (2005). Editing contributed scholarly articles from a language management perspective. *Journal of Second Language Writing*, 14(1), 47-62.
- Karakaya, K. (2017). *A corpus-based and systemic functional analysis of syntactic complexity and nominal modification in academic writing*. Doctoral Dissertation. Iowa State University.
- Karpov, Y. V., & Haywood, H. C. (1998). Two ways to elaborate Vygotsky's concept of mediation. *American Psychologist*, 53(1), 27.

- Karsli, M. B., Karabey, S., Cagiltay, N. E., & Goktas, Y. (2018). Comparison of the discussion sections of PhD dissertations in educational technology: The case of Turkey and the USA. *Scientometrics*, *117*(3), 1381-1403.
- Kasperavičienė, R., & Motiejūnienė, J. (2013). On language editing of research articles translated from Lithuanian to English. *Kalbu Studijos*, (22), 78-85.
- Kawase, T. (2015). Metadiscourse in the introductions of PhD theses and research articles. *Journal of English for Academic Purposes*, *20*, 114-124.
- Kennedy, G. (2014). *An introduction to corpus linguistics*. New York: Routledge.
- Kırkgöz, Y. (2005). English language teaching in Turkey: Challenges for the 21st Century. In G. Braine (Ed), *Teaching English to the world: History, curriculum and practice* (pp. 159-175). London: Lawrence Erlbaum Associates.
- Kırkgöz, Y. (2010). An analysis of written errors of Turkish adult learners of English. *Procedia-Social and Behavioral Sciences*, *2*(2), 4352-4358.
- Klein, D., & Manning, C. D. (2003). Fast exact inference with a factored model for natural language parsing. In *Advances in Neural Information Processing Systems 15 (NIPS 2002)* (pp. 3-10), Cambridge, MA: MIT.
- Knoch, U., Rouhshad, A., & Storch, N. (2014). Does the writing of undergraduate ESL students develop after one year of study in an English-medium university?. *Assessing Writing*, *21*, 1-17.
- Korthagen, F. A. (2010). Situated learning theory and the pedagogy of teacher education: Towards an integrative view of teacher behavior and teacher learning. *Teaching and Teacher Education*, *26*(1), 98-106.
- Koyalán, A., & Mumford, S. (2011). Changes to English as an Additional Language writers' research articles: From spoken to written register. *English for Specific Purposes*, *30*(2), 113-123.
- Kroll, B. (1990). What does time buy? ESL student performance on home versus class compositions. In B. Kroll (Ed.), *Second language writing: Research insights for the classroom* (pp. 140-154). Cambridge: Cambridge University.

- Kuiken, F., & Vedder, I. (2012). Syntactic complexity, lexical variation and accuracy as a function of task complexity and proficiency level in L2 writing and speaking. In A. Housen, F. Kuiken, & I. Vedder (Eds.), *Dimensions of L2 performance and proficiency: Complexity, accuracy and fluency in SLA* (pp. 143– 169). Philadelphia/Amsterdam: John Benjamins.
- Kukla, A. (2000). *Social constructivism and the philosophy of science*. London: Routledge.
- Kübler, S., & Zinsmeister, H. (2015). *Corpus linguistics and linguistically annotated corpora*. New York: Bloomsbury.
- Kyle, K., & Crossley, S. (2017). Assessing syntactic sophistication in L2 writing: A usage-based approach. *Language Testing*, 34(4), 513-535.
- Lancaster, Z. (2016). Using corpus results to guide the discourse-based interview: A study of one student's awareness of stance in academic writing in philosophy. *Journal of Writing Research*, 8(1), 119-148.
- Larsen-Freeman, D. (1978). An ESL index of development. *TESOL Quarterly*, 12(4), 439-448.
- Larsen-Freeman, D., & Long, M. H. (2014). *An introduction to second language acquisition research*. New York: Routledge.
- Laufer, B. (1991). The development of L2 lexis in the expression of the advanced learner. *The Modern Language Journal*, 75(4), 440-448.
- Laufer, B., & Nation, P. (1995). Vocabulary size and use: Lexical richness in L2 written production. *Applied linguistics*, 16(3), 307-322.
- Lave, J., & Wenger, E. (1991). *Situated learning: Legitimate peripheral participation*. Cambridge: Cambridge University.
- Lee, D. Y. W. 2010, 'What corpora are available?', in M. McCarthy & A. O'Keeffe (eds), *The Routledge handbook of corpus linguistics* (pp. 107-121), Abingdon: Routledge.
- Lee, D. Y., & Chen, S. X. (2009). Making a bigger deal of the smaller words: Function words and other key items in research writing by Chinese learners. *Journal of Second Language Writing*, 18(4), 281-296.

- Lee, Y. W., Chodorow, M., & Gentile, C. (2016). Investigating patterns of writing errors for different L1 groups through error-coded ESL learners' essays. *Foreign Languages Education*, 23(1), 169-190.
- Leech, G. (1992). Corpora and theories of linguistic performance. In J. Svartvik (Ed.), *Directions in corpus linguistics* (pp. 105–122). Berlin: Mouton de Gruyter.
- Leech, G., Garside, R., & Bryant, M. (1994). The large-scale grammatical tagging of text: Experience with the British National Corpus. In Oostdijk, N. & De Haan, P. (Eds.), *Corpus-based research into language* (pp. 47-63). Amsterdam: Rodopi.
- Leibowitz, B. (2000). The importance of writing and teaching writing in the academy. In B. Leibowitz & Y. Mohamed (Eds.), *Routes to writing in Southern Africa* (pp. 15-41). Cape Town: Silk Road International.
- Leki, I. (1991). The preferences of ESL students for error correction in college-level writing classes. *Foreign Language Annals*, 24(3), 203-218.
- Leki, I., & Carson, J. G. (1994). Students' perceptions of EAP writing instruction and writing needs across the disciplines. *TESOL Quarterly*, 28(1), 81-101.
- Lenko-Szymanska, A. (2003). "Lexical problem areas in the advanced learner corpus of written data". In B. Lewandowska-Tomaszczyk (Ed.), *PALC' 2001. Practical applications in language corpora* (pp. 505-520). Frankfurt am Main: Peter Lang.
- Leńko-Szymańska, A. (2008). Non-native or non-expert? The use of connectors in native and foreign language learners' texts. *Acquisition et Interaction en Langue Etrangère*, 27, 91-108.
- Lent, R. C. (2015). *This is disciplinary literacy: Reading, writing, thinking, and doing... content area by content area*. California: Corwin.
- Leontiev, A. N. (1978). *Activity, consciousness, and personality*. Englewood Cliffs: Prentice-Hall.
- Levine, D. (2014). *Even you can learn statistics and analytics: An easy to understand guide to statistics and analytics 3rd edition*. New York: Pearson.
- Levy, R. & Andrew, G. (2006). Tregex and Tsurgeon: Tools for querying and manipulating tree data structures. In *Proceedings of the Fifth International Conference on Language Resources and Evaluation*. Genoa, Italy: ELRA, 2231–2234.

- Lewin, B. A. & Fine, J. (1996). The writing of research texts: Genre analysis and its applications. In G. Rijlaarsdam, H. van den Bergh & M. Couzijn (eds.), *Theories, models and methodology in writing research* (pp. 423 – 444). Amsterdam: Amsterdam University.
- Li, Y. (2002). Writing for international publication: The perception of Chinese doctoral researchers. *Asian Journal of English Language Teaching*, 12, 179-193.
- Li, Y. (2005). Multidimensional enculturation: The case of an EFL Chinese doctoral student. *Journal of Asian Pacific Communication*, 15(1), 153-170.
- Li, Y. (2006). Negotiating knowledge contribution to multiple discourse communities: A doctoral student of computer science writing for publication. *Journal of Second Language Writing*, 15(3), 159-178.
- Li, Y. (2007). Apprentice scholarly writing in a community of practice: An intraview of an NNES graduate student writing a research article. *TESOL Quarterly*, 41(1), 55-79.
- Li, Y., & Flowerdew, J. (2007). Shaping Chinese novice scientists' manuscripts for publication. *Journal of Second Language Writing*, 28(3), 100-117.
- Li, Y., & Flowerdew, J. (2020). Teaching English for research publication purposes (ERRP): A review of language teachers' pedagogical initiatives. *English for Specific Purposes*, 59, 29-41.
- Lillis, T. M., & Curry, M. J. (2010). *Academic writing in global context*. London: Routledge.
- Lillis, T., & Curry, M. J. (2006). Professional academic writing by multilingual scholars: Interactions with literacy brokers in the production of English-medium texts. *Written Communication*, 23(1), 3-35.
- Liu, D., & Gleason, J. L. (2002). Acquisition of the article the by nonnative speakers of English: An analysis of four nongeneric uses. *Studies in Second Language Acquisition*, 24(1) 1-26.
- Liu, J. (2004). Co-constructing academic discourse from the periphery: Chinese applied linguists' centripetal participation in scholarly publication. *Asian Journal of English Language Teaching*, 14, 1-22.
- Liu, L. (2018). *Examining syntactic complexity in EFL academic writing*. Doctoral dissertation, The Hong Kong Polytechnic University, Hong Kong.

- Loi, C. K. (2010). Research article introductions in Chinese and English: A comparative genre-based study. *Journal of English for Academic Purposes*, 9(4), 267–279.
- Lombard, M., Snyder-Duch, J., & Bracken, C. C. (2004). A call for standardization in content analysis reliability. *Human Communication Research*, 30(3), 434.
- López, C. L., & Manalastas, G. (2017). Errors in L1 and L2 university students' writing in English: Grammar, spelling and punctuation. *RAEL: Revista Electrónica de Lingüística Aplicada*, 16(1), 118-134.
- López, W. C. (2009). *Error analysis in a learner corpus: what are the learners' strategies*. Retrieved from <https://www.um.es/lacell/aelinco/contenido/pdf/45.pdf>
- Lu, X. (2010). Automatic analysis of syntactic complexity in second language writing. *International Journal of Corpus Linguistics*, 15(4), 474-496.
- Lu, X. (2011). A corpus-based evaluation of syntactic complexity measures as indices of college-level ESL writers' language development. *TESOL Quarterly*, 45(1), 36-62.
- Lu, X. (2012). The relationship of lexical richness to the quality of ESL learners' oral narratives. *The Modern Language Journal*, 96(2), 190-208.
- Lu, X., & Ai, H. (2015). Syntactic complexity in college-level English writing: Differences among writers with diverse L1 backgrounds. *Journal of Second Language Writing*, 29, 16-27.
- Lu, X., Casal, J. E., & Liu, Y. (2020). The rhetorical functions of syntactically complex sentences in social science research article introductions. *Journal of English for Academic Purposes*, 44, 100832.
- Lüdeling, A., Walter, M., Kroymann, E. and Adolphs, P. (2005, July). *Multi-level annotation error annotation in learner corpora*. Paper presented at Proceedings of Corpus Linguistics 2005, Birmingham.
- Mac Donald, P. (2016). We all make mistakes!. Analysing an error-coded corpus of Spanish university students written English. *Complutense Journal of English Studies*, 24, 103-129.
- MacDonald, S. (1994). *Professional academic writing in the humanities and social sciences*. Illinois: Southern Illinois University.

- Mackey, A., & Gass, S. M. (2015). *Second language research: Methodology and design*. New York: Routledge.
- Mackey, A., & Gass, S. M. (Eds.). (2011). *Research methods in second language acquisition: A practical guide* (Vol. 7). New York: John Wiley & Sons.
- Maleki, H. (2008). *Textbook (Design and writing)*. Tehran: Allameh Tabatabai University.
- Malhotra, A., & Krishna, S. (2018). Release velocities and bowler performance in cricket. *Journal of Applied Statistics*, 45(9), 1616-1627.
- Man, J. P., Weinkauff, J. G., Tsang, M., & Sin, J. H. D. D. (2004). Why do some countries publish more than others? An international comparison of research funding, English proficiency and publication output in highly ranked general medical journals. *European Journal of Epidemiology*, 19(8), 811-817.
- Mancilla, R. L., Polat, N., & Akcay, A. O. (2017). An investigation of native and nonnative English speakers' levels of written syntactic complexity in asynchronous online discussions. *Applied Linguistics*, 38(1), 112-134.
- Mann, H. B., & Whitney, D. R. (1947). On a test of whether one of two random variables is stochastically larger than the other. *The Annals of Mathematical Statistics*, 18(1), 50-60.
- Mansourizadeh, K., & Ahmad, U. K. (2011). Citation practices among non-native expert and novice scientific writers. *Journal of English for Academic Purposes*, 10(3), 152-161.
- Marcus, M., Santorini, B., & Marcinkiewicz, M. A. (1993). Building a large annotated corpus of English: The Penn Treebank. *Computational Linguistics*, 19(2), 313-330.
- Mark, K. L. (1998). The significance of learner corpus data in relation to the problems of language teaching. *Bulletin of General Education*, 312, 77-90.
- Martin, J. R. 1992. *English for text: system and structure*. Amsterdam: John Benjamins.
- Martín, P., & Pérez, I. K. L. (2014). Convincing peers of the value of one's research: A genre analysis of rhetorical promotion in academic texts. *English for Specific Purposes*, 34, 1-13.
- Martínez, I. A. (2005). Native and non-native writers' use of first person pronouns in the different sections of biology research articles in English. *Journal of Second Language Writing*, 14(3), 174-190.

- Maschler, Y. & Schiffrin, D. (2015). Discourse markers: Language, meaning, and context. In D. T., Heidi, E. Hamilton & D. Schiffrin (Eds.), *The handbook of discourse analysis* (pp. 189 - 221). Chichester, UK: John Wiley & Sons.
- Mason, O., & Uzar, R. (2000). *NLP meets TEFL: Tracing the zero article*. Paper presented in PALC'99: Practical Applications in Language Corpora International Conference, the University of Łódź.
- Matsuda, P. & Jeffery, J. (2012). Voice in student essays. In K. Hyland & C. S. Guinda (Eds.), *Stance and voice in written academic genres*, (pp. 151–165). London: Palgrave.
- Matsuda, P. K. (1997). Contrastive rhetoric in context: A dynamic model of L2 writing. *Journal of Second Language Writing*, 6(1), 45-60.
- Matsuda, P. K. (2003). Second language writing in the twentieth century: A situated historical perspective. *Exploring the Dynamics of Second Language Writing*, 1, 15-34.
- Matsuda, P. K. (2015). Identity in written discourse. *Annual Review of Applied Linguistics*, 35, 140-159.
- Matthiessen, C. M., & Halliday, M. A. (2009). *Systemic functional grammar: A first step into the theory*. Beijing: Higher Education.
- McCarthy, M., & Carter, R. (1995). Spoken grammar: What is it and how can we teach it?. *ELT Journal*, 49(3), 207-218.
- McDowell, L. (2020). Error analysis: A methodological exploration and application. In P. Clements, A. Krause, & R. Gentry (Eds.), *Teacher efficacy, learner agency*. Tokio: JALT. <https://doi.org/10.37546/JALTPCP2019-53>
- McDowell, L., & Liardét, C. (2020). Towards specialized language support: An elaborated framework for Error Analysis. *English for Specific Purposes*, 57, 16-28.
- McEnery, T., & Hardie, A. (2011). *Corpus linguistics: Method, theory and practice*. Cambridge: Cambridge University.
- McEnery, T., & Wilson, A. (1997). Teaching and language corpora (TALC). *ReCALL*, 9(1), 5-14.

- McEnery, T., & Wilson, A. (2003). Corpus linguistics. In R. Mitkov (Ed.), *The Oxford handbook of computational linguistics* (pp. 448-463). Oxford: Oxford University.
- McEnery, T., Xiao, R., & Tono, Y. (2006). *Corpus-based language studies: An advanced resource book*. New York: Taylor & Francis.
- McIntosh, K., Connor, U., & Gokpinar-Shelton, E. (2017). What intercultural rhetoric can bring to EAP/ESP writing studies in an English as a lingua franca world. *Journal of English for Academic Purposes*, 29, 12-20.
- Mclaren, S. (2003). *Easy writer students guide to writing essays and report*. New Delhi: Viva.
- Mediero Durá n, Esther (2013). *Lexical errors and language acquisition by university level EFL students*. Paper presented at 10th TESOL/Applied Linguistics Graduate Students Conference, 16th February, 2013, Greenville, NC, USA.
- Melander, B., Swales, J. M., and Fredrickson, K. M. (1997) Journal abstracts from three academic fields in the United States and Sweden: National or disciplinary proclivities? In A. Duszak, (Ed.), *Culture and Styles of Academic Discourse* (pp. 251 – 272). Berlin: Mouton de Gruyter.
- Mestre, E. M. (2011). *CEFR & Error analysis in second language teaching at university level*. Saarbrücken: Lambert Academic.
- Meunier, F., & Gouverneur, C. (2009). New types of corpora for new educational challenges. In Aijmer K. (Ed.), *Corpora and language teaching* (pp. 179-201). Amsterdam: Benjamins.
- Min, H. T. (2014). Participating in international academic publishing: A Taiwan perspective. *TESOL Quarterly*, 48(1), 188-200.
- Mišak, A., Marušić, M., & Marušić, A. (2005). Manuscript editing as a way of teaching academic writing: Experience from a small scientific journal. *Journal of Second Language Writing*, 14(2), 122-131.
- Mollin, S. (2007). The Hansard hazard: Gauging the accuracy of British parliamentary transcripts. *Corpora*, 2(2), 187-210.
- Monroe, J. H. (1975). Measuring and enhancing syntactic fluency in French. *The French Review*, 48(6), 1023-1031.

- Moreno, A. I., & Swales, J. M. (2018). Strengthening move analysis methodology towards bridging the function-form gap. *English for Specific Purposes*, 50, 40-63.
- Moreno, A., & Suarez, L. (2008). A study of critical attitude across English and Spanish academic book reviews. *Journal of English for Academic Purposes*, 7(1), 15–26.
- Morton, J., Storch, N., & Thompson, C. (2015). What our students tell us: Perceptions of three multilingual students on their academic writing in first year. *Journal of Second Language Writing*, 30, 1-13.
- Mukherjee, J. (2004). The state of the art in corpus linguistics: Three book-length perspectives. *English Language & Linguistics*, 8(1), 103-119.
- Mukherjee, J. (2005). The native speaker is alive and kicking: Linguistic and language-pedagogical perspectives. *Anglistik*, 16(2), 7-23.
- Mukherjee, J., & Rohrbach, J.-M. (2006). Rethinking applied corpus linguistics from a language-pedagogical perspective: New departures in learner corpus research. In B. Kettemann & G. Marko (Eds.), *Planning, gluing and painting corpora: Inside the applied corpus linguist's workshop* (pp. 205-232). Frankfurt: Peter Lang.
- Mungra, P., & Webber, P. (2010). Peer review process in medical research publications: Language and content comments. *English for Specific Purposes*, 29(1), 43-53.
- Muresan, L. M., & Pérez-Llantada, C. (2014). English for research publication and dissemination in bi-/multiliterate environments: The case of Romanian academics. *Journal of English for Academic Purposes*, 13, 53-64.
- Müller, S. (2005). *Discourse markers in native and non-native English discourse*. Amsterdam: John Benjamins.
- Myers, G. (1998). Displaying opinions: Topics and disagreement in focus groups. *Language in Society*, 27(1), 85-111.
- Nakov, P. (2013). On the interpretation of noun compounds: Syntax, semantics, and entailment. *Natural Language Engineering*, 19(3), 291-330.
- Nation, I. S. (2013). *Learning vocabulary in another language Google eBook*. Cambridge: Cambridge University.
- Neary-Sundquist, C. A. (2017). Syntactic complexity at multiple proficiency levels of L2 German speech. *International Journal of Applied Linguistics*, 27(1), 242-262.

- Neff, J., Ballesteros, F., Dafouz, E., Martínez, F., Rica, J. P., Díez, M., & Prieto, R. (2007). A contrastive functional analysis of errors in Spanish EFL university writers' argumentative texts: Corpus-based study. In E. Fitzpatrick (Ed.) *Corpus Linguistics Beyond the Word* (pp. 203-225). Rodopi: Brill.
- Nesi, H., & Gardner, S. (2012). *Genres across the disciplines: Student writing in higher education*. Cambridge: Cambridge University.
- Nesselhauf, N. (2004a). How learner corpus analysis can contribute to language teaching: A study of support verb constructions. *Corpora and Language Learners*, 17, 109-124.
- Nesselhauf, N. (2004b). Learner corpora and their potential for language teaching. *How to Use Corpora in Language Teaching*, 12, 125-156.
- Nesselhauf, N. (2005). *Collocations in a learner corpus (Vol. 14)*. Amsterdam: John Benjamins.
- Nga, N. T. H. (2009). Academic English at tertiary level: What, why and how. *VNU Journal of Foreign Studies*, 25(2), 112-117.
- Nicholls, D. (2003). The Cambridge learner corpus: Error coding and analysis for lexicography and ELT. In: Archer D., Rayon P., Wilson A. & McEnery T. (eds.) *Proceedings of the Corpus Linguistics 2003 conference UCREL* (pp. 572-581). Lancaster University.
- Norrby, C., & Håkansson, G. (2007). The interaction of complexity and grammatical processability: The case of Swedish as a foreign language. *International Review of Applied Linguistics in Language Teaching*, 45(1), 45-68.
- Norris, J. M., & Ortega, L. (2009). Towards an organic approach to investigating CAF in instructed SLA: The case of complexity. *Applied linguistics*, 30(4), 555-578.
- Nugues, P. M. (2006). *An introduction to language processing with perl and prolog*. Berlin: Springer.
- Nwogu, K. N. (1989). *Discourse variation in medical texts: Schema, theme and cohesion in professional and journalistic accounts*. Doctoral Dissertation, University of Aston, Birmingham.
- Nystrand, M. (1986). *The structure of written communication: Studies in reciprocity between writers and readers*. Rodopi: Brill.

- Nystrand, M. (2006). The social and historical context for writing research. In C. A. MacArthur, S. Graham, & J. Fitzgerald (Eds.), *Handbook of writing research* (pp. 11–27). New York: Guilford.
- Nystrand, M. (Ed.). (1982). *What writers know: The language, process, and structure of written discourse*. New York: Academic.
- O'Donnell, M. B., & Römer, U. (2009b). *Michigan corpus of upper-level student papers*. Retrieved from: <http://micusp.elicorpora.info/>
- O'Donnell, M. B., & Römer, U. (2009a). *From student hard drive to web corpus: The design, compilation, annotation and online distribution of the MICUSP corpus*. A poster presented in *ICAME 30*, United Kingdom.
- O'Donnell, M. (2008). The UAM CorpusTool: Software for corpus annotation and exploration". In B. Callejas, M. Carmen et al. (Eds), *Applied Linguistics Now: Understanding Language and Mind* (pp. 1433-1447) Almería : Universidad de Almería.
- O'Keeffe, A., & McCarthy, M. (2010). Historical perspective: what are corpora and how have they evolved? In M. McCarthy & A. O'Keeffe (Eds.), *The Routledge Handbook of Corpus Linguistics* (pp. 31-41). New York: Routledge.
- Ortega, L (2012) Interlanguage complexity: A construct in search of theoretical renewal. In B. Kortmann & B. Szmrecsanyi (Eds.), *Linguistic complexity: Second language acquisition, indigenization, contact* (pp. 127-155). Berlin / Boston, MA: Mouton de Gruyter.
- Ortega, L. (2003). Syntactic complexity measures and their relationship to L2 proficiency: A research synthesis of college-level L2 writing. *Applied linguistics*, 24(4), 492-518.
- Oshima, A. & Hogue, A. (2006). *Writing academic English*. New York: Pearson Education.
- Ozturk, I. (2007). The textual organisation of research article introductions in applied linguistics: Variability within a single discipline. *English for Specific Purposes*, 26(1), 25-38.
- Öztuna, D., Elhan, A. H., & Tüccar, E. (2006). Investigation of four different normality tests in terms of type 1 error rate and power under different distributions. *Turkish Journal of Medical Sciences*, 36(3), 171-176.

- Pallotti, G. (2009). CAF: Defining, refining and differentiating constructs. *Applied Linguistics*, 30(4), 590-601.
- Paltridge, B. (1994). Genre analysis and the identification of textual boundaries. *Applied Linguistics*, 15(3), 288 – 299.
- Paltridge, B. (2001). *Genre and the language learning classroom*. Michigan: University of Michigan.
- Paltridge, B. (2012). *Discourse analysis: An introduction*. London: Bloomsbury.
- Paltridge, B. (2014). Genre and second-language academic writing. *Language Teaching*, 47(3), 303-318.
- Pan, F., Reppen, R., & Biber, D. (2016). Comparing patterns of L1 versus L2 English academic professionals: Lexical bundles in telecommunications research journals. *Journal of English for Academic Purposes*, 21, 60-71.
- Paquot, M. (2017). L1 frequency in foreign language acquisition: Recurrent word combinations in French and Spanish EFL learner writing. *Second Language Research*, 33(1), 13-32.
- Paquot, M. (2018). Phraseological competence: A useful toolbox to delimitate CEFR levels in higher education? Insights from a study of EFL learners' use of statistical collocations. *Language Assessment Quarterly*, 15(1), 29-43.
- Park, K. (2010). *Enhancing lexicogrammatical performance through corpus-based mediation in L2 academic writing instruction*. Unpublished doctoral dissertation, The Pennsylvania State University, USA.
- Parrish, S. M. (1962). Problems in the making of computer concordances. *Studies in Bibliography*, 15, 1-14.
- Peacock, M. (2002). Communicative moves in the discussion section of research articles. *System*, 30(4), 479-497.
- Pecorari, D. (2002). *Original reproductions: An investigation of the source use of postgraduate second-language writers*. Doctoral dissertation, University of Birmingham, Birmingham.
- Péry-Woodley, M. P. (1990). Contrasting discourses: Contrastive analysis and a discourse approach to writing. *Language Teaching*, 23(3), 143-151.

- Petch-Tyson, S. (2000). Demonstrative expressions in argumentative discourse. A computerbased comparison of non-native and native English. In S. Botley & A. M. McEnery (Eds.) *Corpus-based and computational approaches to discourse anaphora* (pp. 43-64). Amsterdam: John Benjamins.
- Phillipson, R. (1992). *Linguistic imperialism*. Oxford: Oxford University.
- Planken, B. C., van Meurs, W. F. J., & Maria, K. (2019). Do errors matter? The effects of actual and perceived L2 English errors in writing on native and non-native English speakers' evaluations of the text, the writer and the persuasiveness of the text. *International Journal of English Language Teaching*, 6(1), 1-13.
- Polat, N., Mahalingappa, L., & Mancilla, R. L. (2019). Longitudinal growth trajectories of written syntactic complexity: The case of Turkish learners in an intensive English program. *Applied Linguistics*, 41(5), 988-711.
- Polio, C. G. (1997). Measures of linguistic accuracy in second language writing research. *Language learning*, 47(1), 101-143.
- Rasier, L., & Hiligsmann, P. (2007). Prosodic transfer from L1 to L2. Theoretical and methodological issues. *Nouveaux Cahiers de Linguistique Française*, 28(2007), 41-66.
- Reid, J. M. (2000). *The process of composition*. New York: Longman.
- Richards, B. (1987). Type/token ratios: What do they really tell us?. *Journal of Child Language*, 14(2), 201-209.
- Richards, J. C. (1972). Social factors, interlanguage, and language learning. *Language Learning*, 22(2), 159-188.
- Richards, J. C., & Renandya, W. A. (Eds.). (2002). *Methodology in language teaching: An anthology of current practice*. Cambridge: Cambridge University.
- Rimrott, A., & Heift, T. (2005). Language learners and generic spell checkers in CALL. *CALICO journal*, 23(1) 17-48.
- Ringbom, H. (1987). *The role of the first language in foreign language learning*. Philadelphia: Multilingual Matters.
- Robb, T., Ross, S., & Shortreed, I. (1986). Salience of feedback on error and its effect on EFL writing quality. *TESOL Quarterly*, 20(1), 83-96.

- Römer, U. (2011). Corpus research applications in second language teaching. *Annual Review of Applied Linguistics*, 31, 205-225.
- Rundell, M., & Stock, P. (1992). The corpus revolution. *English Today*, 8(2), 9-14.
- Safari, I. (2018). A corpus-based contrastive study of code glosses used in English academic articles written by authors of politics and applied linguistics. *International Journal of Linguistics*, 10(2), 39-47.
- Sajid, M. (2016). Diction and expression in error analysis can enhance academic writing of L2 university students. *Advances in Language and Literary Studies*, 7(3), 71-79.
- Salager-Meyer, F. (1992). A text-type and move analysis study of verb tense and modality distribution in medical English abstracts. *English for Specific Purposes*, 11(2), 93-113.
- Salager-Meyer, F. (2008). Scientific publishing in developing countries: Challenges for the future. *Journal of English for Academic Purposes*, 7(2), 121-132.
- Samraj, B. (2002). Introductions in research articles: Variations across disciplines. *English for Specific Purposes*, 21(1), 1-17.
- Samraj, B. (2005). An exploration of a genre set: Research article abstracts and introductions in two disciplines. *English for Specific Purposes*, 24(2), 141–156.
- Schachter, J. (1983). A new account of language transfer. In S. M. Gass & L. Selinker (Eds.) *Language transfer in language learning. Issues in second language research* (pp. 98-111). Rowley: Newbury House.
- Schembri, A. C., Johnston, T., & Goswell, D. (2006). *NAME dropping: Location variation in Australian sign language*. Washington: Gallaudet University.
- Schmitt, N., & Hemchua, S. (2006). An analysis of lexical errors in the English composition of Thai learners: 2698. *Prospect: An Australian Journal of TESOL*, 21(3), 3-25.
- SCImago. (2020). *SJR — SCImago journal & country rank*. Retrieved from <http://www.scimagojr.com/countryrank.php>
- Scott, M. (2004). *WordSmith tools*. Oxford: Oxford University.
- Scott, M. (2016). *WordSmith tools version 7*. Stroud: Lexical Analysis Software.
- Selinker, L. (1989). CA/EA/IL: The earliest experimental record. *IRAL: International Review of Applied Linguistics in Language Teaching*, 27(4), 267.

- Selvi, A. F. (2011). World Englishes in the Turkish sociolinguistic context. *World Englishes*, 30(2), 182-199.
- Seoane, E., & Hundt, M. (2018). Voice alternation and authorial presence: Variation across disciplinary areas in academic English. *Journal of English Linguistics*, 46(1), 3-22.
- Severino, C., & Prim, S. N. (2015). Word choice errors in Chinese students' English writing and how online writing center tutors respond to them. *The Writing Center Journal*, 34(2), 115-143.
- Shaw, P. (1991). Science research students' composing processes. *English for Specific Purposes*, 10(3), 189-206
- Shaw, P. M. (2013). Grammar in academic writing. In C. Chapelle (Ed), *The encyclopedia of applied linguistics*. Oxford: John Wiley & Sons.
- Sheldon, E. (2009). From one I to another: Discursive construction of self-representation in English and Castilian Spanish research articles. *English for Specific Purposes*, 28(4), 251-265.
- Shin, Y. K., Cortes, V., & Yoo, I. W. (2018). Using lexical bundles as a tool to analyze definite article use in L2 academic writing: An exploratory study. *Journal of Second Language Writing*, 39, 29-41.
- Silva, T. (1990). Second language composition instruction: Developments, issues, and directions in ESL. In B. Kroll (Ed.), *Second language writing: Research insights for the classroom* (pp. 11-23). New York: Cambridge University.
- Silva, T. (1993). Toward an understanding of the distinct nature of L2 writing: The ESL research and its implications. *TESOL quarterly*, 27(4), 657-677.
- Silva, T., & Leki, I. (2004). Family matters: The influence of applied linguistics and composition studies on second language writing studies—Past, present, and future. *The Modern Language Journal*, 88(1), 1-13.
- Sinclair, J. (1991). *Corpus, concordance, collocation*. Oxford: Oxford University.
- Sinclair, J. (2005). Corpus and text-basic principles. In M. Wynne (Ed.), *Developing linguistic corpora: A guide to good practice* (pp. 1-16). Oxford: Oxbow.
- Sinclair, J. M. (1996). The empty lexicon. *International Journal of Corpus Linguistics*, 1(1), 99-119.

- Sloan, D., & Porter, E. (2010). Changing international student and business staff perceptions of in-session EAP: Using the CEM model. *Journal of English for Academic Purposes*, 9(3), 198-210.
- Smoak, R. (2003). What is English for specific purposes. *English Teaching Forum*, 41(2), 22-27.
- Song, R., & Wang, H. (2019). An investigation into the syntactic complexity of Chinese doctoral dissertation abstracts. *Journal of PLA University of Foreign Languages*, 42, 84-91.
- Sotillo, S. M. (2000). Discourse functions and syntactic complexity in synchronous and asynchronous communication. *Language Learning & Technology*, 4(1), 77-110.
- Spack, R. (1988). Initiating ESL students into the academic discourse community: How far should we go?. *TESOL Quarterly*, 22(1), 29-51.
- St. John, M. J. (1987). Writing processes of Spanish scientists publishing in English. *English for Specific Purposes*, 6(2), 113-120.
- Stanton-Salazar, R. D. (1997). A social capital framework for understanding the socialization of racial minority children and youths. *Harvard Educational Review*, 67(1), 1-40.
- Staples, S., & Reppen, R. (2016). Understanding first-year L2 writing: A lexico-grammatical analysis across L1s, genres, and language ratings. *Journal of Second Language Writing*, 32, 17-35.
- Steele, V. (1992). *Product and process writing: A Comparison*. Rowley: Newbury House.
- Stemler, S. E. (2004). A comparison of consensus, consistency, and measurement approaches to estimating interrater reliability. *Practical Assessment, Research, and Evaluation*, 9(1), 4.
- Stockwell, G., & Harrington, M. (2003). The incidental development of L2 proficiency in NS-NNS email interactions. *CALICO Journal*, 20(2), 337-359.
- Stolk, P., Egberts, A. C. G., & Leufkens, H. G. M. (2002). Fate of abstracts presented at five international conferences on pharmacoepidemiology (ICPE): 1995–1999. *Pharmacoepidemiology and Drug Safety*, 11(2), 105-111.

- Sun, C., & Feng, G. (2009). Process approach to teaching writing applied in different teaching models. *English Language Teaching*, 2(1), 150-155.
- Swales, J. (1987). Utilizing the literatures in teaching the research paper. *TESOL quarterly*, 21(1), 41-68.
- Swales, J. (1990). *Genre analysis: English in academic and research settings*. Cambridge: Cambridge University.
- Swales, J. 1981. *Aspects of article introductions*. Aston ESP Research Report No. 1. Language Studies Unit. University of Aston in Birmingham.
- Swales, J. M. (1997). English as Tyrannosaurus rex. *World Englishes*, 16(3), 373-382.
- Swales, J. M. (2004). *Research genres: Explorations and applications*. Cambridge: Cambridge University.
- Swales, J. M., & Feak, C. B. (2011). *Navigating academia: Writing supporting genres*. Michigan: University of Michigan.
- Taguchi, N., Crawford, W., & Wetzel, D. Z. (2013). What linguistic features are indicative of writing quality? A case of argumentative essays in a college composition program. *TESOL Quarterly*, 47(2), 420-430.
- Tardy, C. M. (2005). "It's like a story": Rhetorical knowledge development in advanced academic literacy. *Journal of English for Academic Purposes*, 4(4), 325-338.
- Tardy, C. M. (2011). Genre analysis. In K. Hyland & B. Paltridge (Eds.), *Bloomsbury companion to discourse analysis* (pp. 54-68). London: Bloomsbury.
- Tardy, C. M., & Jwa, S. (2016). Composition studies and EAP. In K. Hyland & P. Shaw (Eds.), *The Routledge handbook of English for academic purposes* (pp. 80-92). New York: Routledge.
- Tarone, E., Downing, B., Cohen, A., Gillette, S., Murie, R., & Dailey, B. (1993). The writing of Southeast Asian-American students in secondary school and university. *Journal of Second Language Writing*, 2(2), 149-172.
- Teodorescu, A. (2010). Teaching English for specific purposes. *Petroleum-Gas University Of Ploiesti Bulletin, Philology Series*, 62(2).
- Thaiss, C., & Zawacki, T. M. (2006). *Engaged writers and dynamic disciplines: research on the academic writing life*. Portsmouth, NH: Boynton/Cook-Heinemann.

- Thoday, E. (2008). Issues in building learner corpora: An investigation into the acquisition of German passive constructions. *Newcastle Working Papers in Linguistics*, 14, 145-155.
- Thomas, D. (2005). *Type-token ratios in one teacher's classroom talk: An investigation of lexical complexity*. United Kingdom: University of Birmingham.
- Thompson, D. K. (1993). Arguing for experimental “facts” in science: A study of research article results sections in biochemistry. *Written communication*, 10(1), 106-128.
- Thompson, G. (2001). Interaction in academic writing: Learning to argue with the reader. *Applied linguistics*, 22(1), 58-78.
- Thonney, T. (2011). Teaching the conventions of academic discourse. *Teaching English in the Two Year College*, 38(4), 347.
- Tognini-Bonelli, E. (2001). *Corpus linguistics at work* (Vol. 6). Amsterdam: John Benjamins.
- Tono, Y. (2000). A corpus-based analysis of interlanguage development: Analysing POS tag sequences of EFL learner corpora. *Practical Applications in Language Corpora*, 123-132.
- Tse, P., & Hyland, K. (2008). ‘Robot kung fu’: Gender and professional identity in biology and philosophy reviews. *Journal of Pragmatics*, 40(7), 1232-1248.
- Tübitak Araştırma ve Yayın Etiği Kurulu Yönetmeliği. (2010). Retrieved from http://www.tubitak.gov.tr/tubitak_content_files/mevzuat/yonetmelik/YONETMELIK_III_9.pdf
- Tweedie, F. J., & Baayen, R. H. (1998). How variable may a constant be? Measures of lexical richness in perspective. *Computers and the Humanities*, 32(5), 323-352.
- Unicode Consortium. (2015). The unicode® standard version 8.0--core specification. *The Unicode Consortium*, Mountain View, CA.
- University Ranking by Academic Performance. (2019). *World Ranking*. Retrieved from https://www.urapcenter.org/Rankings/2019-2020/World_Ranking_2019-2020
- Uysal, H. H. (2008). Tracing the culture behind writing: Rhetorical patterns and bidirectional transfer in L1 and L2 essays of Turkish writers in relation to educational context. *Journal of Second Language Writing*, 17(3), 183-207.

- Uysal, H. H. (2012a). Argumentation across L1 and L2 writing: Exploring cultural influences and transfer issues. *Vigo International Journal of Applied Linguistics*, 9, 133-159.
- Uysal, H. H. (2012b). The critical role of journal selection in scholarly publishing: A search for journal options in language-related research areas and disciplines. *Journal of Language and Linguistic Studies*, 8(1), 50-95.
- Uysal, H. H. (2014). Turkish academic culture in transition: Centre-based state policies and semiperipheral practices of research, publishing and promotion. In K. Bennett (Ed.), *The semiperiphery of academic writing: Discourses, communities and practices* (pp. 165-188). London: Palgrave Macmillan UK.
- Uzuner, S. (2008). Multilingual scholars' participation in core/global academic communities: A literature review. *Journal of English for Academic Purposes*, 7(4), 250-263.
- van Dijk, T. A. (1994). Academic nationalism. *Discourse & Society*, 5(3), 275–276.
- Van Zundert, M., Sluijsmans, D., & Van Merriënboer, J. (2010). Effective peer assessment processes: Research findings and future directions. *Learning and Instruction*, 20(4), 270-279.
- Vanderbauwhede, G. (2012). The integrated contrastive model evaluated: the French and Dutch demonstrative determiner in L1 and L2. *International Journal of Applied Linguistics*, 22(3), 392-413.
- Vann, R. J. (1979). Oral and written syntactic relationships in second language learning. In C. Yorio, K. Perkins, & J. Schachter. (Eds.), *On TESOL '79: The learner in focus*. (pp. 322-329). Washington, DC: TESOL.
- Villamil, O. S., & de Guerrero, M. C. (2006). Sociocultural theory: A framework for understanding the social-cognitive dimensions of peer feedback. In K. Hyland & F. Copland (Eds.), *Feedback in second language writing: Contexts and issues* (pp. 23-41). Cambridge: Cambridge University.
- Vo, S. (2019). Use of lexical features in non-native academic writing. *Journal of Second Language Writing*, 44, 1-12.
- Vyatkina, N. (2013). Specific syntactic complexity: Developmental profiling of individuals based on an annotated learner corpus. *The Modern Language Journal*, 97(1), 11-30.

- Vygotsky, L. S. (1978). *Mind in society: The development of higher psychological processes*. Massachusetts: Harvard University.
- Walková, M. (2019). A three-dimensional model of personal self-mention in research papers. *English for Specific Purposes*, 53, 60-73.
- Wang, Y., Hu, R., & Liu, M. (2017). The geotemporal demographics of academic journals from 1950 to 2013 according to Ulrich's database. *Journal of Informetrics*, 11(3), 655-671.
- Waugh, D. (1998). Practical approaches to teaching punctuation in the primary school. *Reading*, 32(2), 14-17.
- Way, D. P., Joiner, E. G., & Seaman, M. A. (2000). Writing in the secondary foreign language classroom: The effects of prompts and tasks on novice learners of French. *The Modern Language Journal*, 84(2), 171-184.
- Webber, P. (1993). Writing medical articles: a discussion of common errors made by L2 authors and some particular features of discourse. *UNESCO-ALSED LSP Newsletter*, 15(2), 38-49.
- Weller, M. (2002). Delivering learning on the Net: The why, what & how of online education. *The Internet and Higher Education*, 6(3), 293–296.
- Wenger, E. (1998). Communities of practice: Learning as a social system. *Systems Thinker*, 9(5), 2-3.
- Whitaker, A. (2009). *Academic writing guide: A step-by-step-guide to writing academic papers*. City University of Seattle. Retrieved from <http://www.vsm.sk/Curriculum/academicsupport/academicwritingguide.pdf>
- Wigglesworth, G., & Storch, N. (2012). Feedback and writing development through collaboration: A socio-cultural approach. In R. Manchon (Ed.), *L2 writing development: Multiple perspectives* (pp. 69–101). New York: De Gruyter Mouton.
- Wolfe-Quintero, K., Inagaki, S., & Kim, H. Y. (1998). *Second language development in writing: Measures of fluency, accuracy, & complexity* (No. 17). Honolulu: University of Hawaii.
- Wood, D., Bruner, J. S., & Ross, G. (1976). The role of tutoring in problem solving. *Journal of Child Psychology and Psychiatry*, 17(2), 89-100.

- Wu, X., Mauranen, A., & Lei, L. (2020). Syntactic complexity in English as a lingua franca academic writing. *Journal of English for Academic Purposes*, 43, 100798.
- Wulff, S., Römer, U., & Swales, J. (2012). Attended/unattended this in academic student writing: Quantitative and qualitative perspectives. *Corpus Linguistics and Linguistic Theory*, 8(1), 129-157.
- Wynne, M. (2005). Stylistics: Corpus approaches. In K. Brown (Ed.), *Encyclopedia of Language & Linguistics (Second Edition)* (pp. 223-226). New York: Elsevier.
- Xiaoxiao, L., & Yan, L. (2010). A case study of dynamic assessment in EFL process writing. *Chinese Journal of Applied Linguistics*, 33(1), 24-40.
- Xue, H., & Hwa, R. (2014). Improved correction detection in revised ESL sentences. *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*, 2, 599-604.
- Yağız, O. (2009). *The academic writing of Turkish graduate students insocial sciences: Approaches, processes, needs and challenges*. Doctoral Dissertation. Atatürk Üniversitesi, Sosyal Bilimler Enstitüsü, Erzurum.
- Yang, M. N. (2015). A nursing academic word list. *English for Specific Purposes*, 37, 27-38.
- Yang, R. & Allison, D. (2003). Research articles in applied linguistics: Moving from results to conclusions. *English for Specific Purposes*, 22(4), 365 – 85.
- Yang, W., Lu, X., & Weigle, S. C. (2015). Different topics, different discourse: Relationships among writing topic, measures of syntactic complexity, and judgments of writing quality. *Journal of Second Language Writing*, 28, 53-67.
- Yang, Y. (2013). Exploring linguistic and cultural variations in the use of hedges in English and Chinese scientific discourse. *Journal of Pragmatics*, 50(1), 23-36.
- Yin, S., Gao, Y., & Lu, X. (2021). Syntactic complexity of research article part-genres: Differences between emerging and expert international publication writers. *System*, 102427.
- Yoon, H. (2008). More than a linguistic reference: The influence of corpus technology on L2 academic writing. *Language Learning & Technology*, (12)2, 31–48. <http://dx.doi.org/10125/44142>

- Yoon, H. J., & Polio, C. (2017). The linguistic development of students of English as a second language in two written genres. *TESOL Quarterly*, 51(2), 275-301.
- Yore, L. D., & Treagust, D. F. (2006). Current realities and future possibilities: Language and science literacy—empowering research and informing instruction. *International Journal of Science Education*, 28(2-3), 291-314.
- Yükseköğretim Kurulu. (2017a). *Akademik teşvik ödeneği yönetmeliği yayımlandı*. Retrieved from http://www.yok.gov.tr/documents/10279/30318223/Akademik_Tesvik_Odenegi_Yonetmeliği_31_12_2016_Resmi_Gazete_Yayin_Tarihi.pdf
- Yükseköğretim Kurulu. (2017b). *Araştırma üniversitesi olarak belirlenen üniversiteler cumhurbaşkanı Erdoğan tarafından açıklandı*. Retrieved from <http://www.yok.gov.tr/web/guest/arastirma-universiteleri-aciklandi>
- Yükseköğretim Kurulu. (2020). *Yükseköğretim bilgi yönetim sistemi*. Retrieved from <https://istatistik.yok.gov.tr/>
- Zall, P.M. (1963). *Elements of technical report writing*. New York: Harper and Row.
- Zamel, V. (1983). The composing processes of advanced ESL students: Six case studies. *TESOL quarterly*, 17(2), 165-188.
- Zhang, S. (1987). Cognitive complexity and written production in English as a second language. *Language Learning*, 37(4), 469-481.
- Zimmerman, R. (2000). L2 writing: Subprocesses, a model of formulating and empirical findings. *Learning and Instruction* 10(1), 73-99.
- Zughoul, M. R. (1991). Lexical choice: Towards writing problematic word lists. *International Review of Applied Linguistics*, 29(1), 45-60.



APPENDICES



Appendix 1. List of Communication Research Articles in the Reference Corpus

Discipline	Year	Journal	Country	Authors	Title
REF_COM_RA_001	2017	Journal of New Zealand Public Relations Research	New Zealand	Susan Fountaine	What's not to Like?: A Qualitative Study of Young Women Politicians' Self-Framing on Twitter
REF_COM_RA_002	2016	Public Relations Review	New Zealand	Chris Galloway	Media jihad: What PR can learn in Islamic State's public relations masterclass
REF_COM_RA_003	2015	Journalism Studies	UK	Daniel Jackson & Kevin Moloney	INSIDE CHURNALISM & PR, journalism and power relationships in flux
REF_COM_RA_004	2016	Public Relations Review	USA	Cayce Myers	Reconsidering early U.S. public relations institutions: An analysis of publicity and information bureaux 1891–1918
REF_COM_RA_005	2015	Journalism Practice	UK	Matthew Ricketson	When Slow News Is Good News Book-length journalism's role in extending and enlarging daily news

Appendix 2. List of Economics and Administrative Sciences Research Articles in the Reference Corpus

Discipline	Year	Journal	Country	Authors	Title
REF_EAS_RA_001	2016	The Econometrics Journal	USA	Christopher P. Adams	Finite mixture models with one exclusion restriction
REF_EAS_RA_002	2016	Econometrica	USA	Timothy B. Armstrong	Large Market Asymptotics for Differentiated Product Demand Estimators With Economic Models Of Supply
REF_EAS_RA_003	2016	The Journal of Finance	USA	Patrick Bolton	Debt and Money: Financial Constraints and Sovereign Finance
REF_EAS_RA_004	2015	The Journal of Finance	USA	Philip Bond & Itay Goldstein	Government Intervention and Information Aggregation by Prices
REF_EAS_RA_005	2016	Journal of Financial Economics	USA	Jeffrey R. Brown, Anne M. Farrel & Scott J. Weisbenner	Decision-making approaches and the propensity to default: Evidence and implications
REF_EAS_RA_006	2017	Econometrica	USA	Timothy M. Christensen	Nonparametric Stochastic Discount Factor Decomposition
REF_EAS_RA_007	2016	Journal of Financial Economics	USA	Michael Faulkender & Jason M. Smith	Taxes and leverage at multinational corporations
REF_EAS_RA_008	2016	British Journal of Management	UK	Bill Gerrard & Andy Lockett	Team-specific Human Capital and Performance
REF_EAS_RA_009	2017	Econometrica	USA	Bryan S. Graham	An Econometric Model of Network Formation with Degree Heterogeneity
REF_EAS_RA_010	2016	British Journal of Management	UK & Australia	Stewart Johnstone & Adrian Wilkinson	The Potential of Labour-Management Partnership: A Longitudinal Case Analysis
REF_EAS_RA_011	2016	Journal of Financial Economics	USA	Benjamin J. Keys, Devin G. Pope & Jaren C. Pope	Failure to refinance
REF_EAS_RA_012	2016	Applied Economics	USA	Karen Maguire	What's powering wind? The effect of the U.S. state renewable energy policies on wind capacity (1994–2012)

Appendix 3. List of Education Research Articles in the Reference Corpus

Discipline	Year	Journal	Country	Authors	Title
REF_EDU_RA_001	2016	Journal for Research in Mathematics Education	USA	Katherine E. Lewis & Marie B. Fisher	Taking Stock of 40 Years of Research on Mathematical Learning Disability: Methodological Issues and Future Directions
REF_EDU_RA_002	2016	Educational Sciences	USA	Mary Antony Bair	Professionalism: a comparative case study of teachers, nurses, and social workers
REF_EDU_RA_003	2016	American Educational Research Journal	USA	Christine Baron	Using Embedded Visual Coding to Support Contextualization of Historical Texts
REF_EDU_RA_004	2017	Studies in Science Education	USA	Bronwyn Bevan	The promise and the promises of Making in science education
REF_EDU_RA_005	2015	British Journal of Educational Technology	Australia	Matt Bower	Deriving a typology of Web 2.0 learning technologies
REF_EDU_RA_006	2015	European Journal of Education	USA	Lori Beslow	The Pedagogy and Pleasures of Teaching a 21st-Century Skill
REF_EDU_RA_007	2015	Studies in Science Education	UK	Richard Brock	Intuition and insight: two concepts that illuminate the tacit in science education
REF_EDU_RA_008	2016	Studies in Science Education	UK	Richard Brock & Keith S. Taber	The application of the microgenetic method to studies of learning in science education: characteristics of published studies, methodological issues and recommendations for future research
REF_EDU_RA_009	2015	Journal of Research in Science Teaching	USA	Stephen R. Burgin & Troy D. Sadler	Learning Nature of Science Concepts Through a Research Apprenticeship Program: A Comparative Study of Three Approaches
REF_EDU_RA_010	2016	The Curriculum Journal	UK	James Edward Carrol	Exploring historical 'frameworks' as a curriculum goal: a case study examining students' notions of historical significance when using millennia-wide time scales
REF_EDU_RA_011	2016	Theory & Research in Social Education	USA	Christopher H. Clark	Examining the Relationship Between Civic Education and Partisan Alignment in Young Voters

REF_EDU_RA_012	2015	British Journal of Educational Technology	Australia	Phillip Dawson		Five ways to hack and cheat with bring-your-own-device electronic examinations
REF_EDU_RA_013	2017	Educational Psychology An International Journal of Experimental Educational Psychology	USA	Daniel Dinsmore	L.	Examining the ontological and epistemic assumptions of research on metacognition, self-regulation and self-regulated learning
REF_EDU_RA_014	2017	Gifted Child Quarterly	USA	Matthew Edinger	J.	Online Teacher Professional Development for Gifted Education: Examining the Impact of a New Pedagogical Model
REF_EDU_RA_015	2015	Cambridge Journal of Education	UK	Howard Gibson & Jennifer England		The inclusion of pseudowords within the year one phonics 'Screening Check' in English primary schools
REF_EDU_RA_016	2016	Cambridge Journal of Education	USA	Frank Gresham	M.	Social skills assessment and intervention for children and youth
REF_EDU_RA_017	2017	Educational Review	UK	Aneta Hayes		Why international students have been "TEF-ed out"?
REF_EDU_RA_018	2017	Journal of Research in Science Teaching	USA	Benjamin Herman	C.	Students' Environmental NOS Views, Compassion, Intent, and Action: Impact of Place-Based Socioscientific Issues Instruction
REF_EDU_RA_019	2016	American Educational Research Journal	USA	Nicholas Hillman	W.	Geography of College Opportunity: The Case of Education Deserts
REF_EDU_RA_020	2015	British Journal of Educational Technology	New Zealand	Lucie Lindsay		Transformation of teacher practice using mobile technology with one-to-one classes: M-learning pedagogical approaches
REF_EDU_RA_021	2016	Theory & Research in Social Education	USA	Scott Metzger & Richard Paxton	Alan & J.	Gaming History: A Framework for What Video Games Teach About the Past
REF_EDU_RA_022	2016	Educational Review	Scotland	Rosie Mulholland, Andy McKinlay & John Sproule		Teachers in need of space: the content and changing context of work
REF_EDU_RA_023	2016	Computers & Education	Australia	Michelle Neumann	M.	Young children's use of touch screen tablets for writing and reading at home: Relationships with emergent literacy

REF_EDU_RA_024	2016	Educational Review	New Zealand	Karen Nicholas & Jo Fletcher	What is happening in the use of ICT mathematics to support young adolescent learners? A New Zealand experience
REF_EDU_RA_025	2016	Journal of Research in Science Teaching	USA	Jonathan Osborne	F. The Development and Validation of a Learning Progression for Argumentation in Science
REF_EDU_RA_026	2015	Educational Psychology An International Journal of Experimental Educational Psychology	USA	Jennifer Lee Petersen & Janet Shibley Hyde	Trajectories of self-perceived math ability, utility value and interest across middle school as predictors of high school math performance
REF_EDU_RA_027	2017	Physical Education and Sport Pedagogy	USA	Sharon R. Phillips, Kevin Mercier & Sarah Doolittle	Experiences of teacher evaluation systems on high school physical education programs
REF_EDU_RA_028	2017	British Educational Research Journal	New Zealand	Elizabeth Rata	Knowledge and teaching
REF_EDU_RA_029	2015	Harvard Educational Review	USA	Jessica Wolpaw Reyes	Lead Policy and Academic Performance: Insights from Massachusetts
REF_EDU_RA_030	2016	Journal of Research in Science Teaching	USA	William L. Romine, Troy D. Sadler, and Andrew T. Kinslow	Assessment of Scientific Literacy: Development and Validation of the Quantitative Assessment of Socio-Scientific Reasoning (QuASSR)
REF_EDU_RA_031	2015	Educational Psychology An International Journal of Experimental Educational Psychology	USA	Chad A. Rose, Cynthia G. Simpson & Stephanie K. Ellis	The relationship between school belonging, sibling aggression and bullying involvement: implications for students with and without disabilities
REF_EDU_RA_032	2017	Educational Research	Canada	Rachel Ryerson	Creating possibilities: studying the student experience
REF_EDU_RA_033	2016	James International Journal of Educational Technology in Higher Education	Australia	Rosalind James	Tertiary student attitudes to invigilated, online summative examinations

Appendix 4. List of Engineering Research Articles in the Reference Corpus

Discipline	Year	Journal	Country	Authors	Title
REF_ENG_RA_001	2017	Applied Energy	Ireland	R. O'Shea, D.M. Wall & J.D. Murphy	An energy and greenhouse gas comparison of centralised biogas production with road haulage of pig slurry, and decentralised biogas production with biogas transportation in a low-pressure pipe network
REF_ENG_RA_002	2016	International Journal of Plasticity	USA	Mattia Bacca & Robert M. McMeeking	Latent heat saturation in microstructural evolution by severe plastic deformation
REF_ENG_RA_003	2016	IEEE Transactions On Pattern Analysis And Machine Intelligence	USA	Jonathan T. Barron & Jitendra Malik	Intrinsic Scene Properties from a Single RGB-D Image
REF_ENG_RA_004	2017	Materials Characterization	USA	A.N. Black, S.M. Copley, J.A. Todd	A new method for isolating plasma interactions from those of the laser beam during plasma nitriding
REF_ENG_RA_005	2016	Materials Characterization	USA	Todd Book, Michael D. Sangid	Strain Localization in Ti-6Al-4V Widmanstätten Microstructures Produced by Additive Manufacturing
REF_ENG_RA_006	2017	Canadian Journal of Civil Engineering	Canada	David N. Bristow and Michele Bristow	Retrofitting for resiliency and sustainability of households
REF_ENG_RA_007	2017	Materials Characterization	USA	Jessica L. Buckner, Stephen W. Stafford & Darren M. Cone	Microstructural characterization of Ti-6Al-4V X-links from the Space Shuttle Columbia
REF_ENG_RA_008	2017	Materials Science & Engineering A	USA	D.W. Brown	Deformation Behavior of Additively Manufactures GP1 Stainless Steel
REF_ENG_RA_009	2015	Nature Biotechnology	USA	Brian Cleary	Detection of low-abundance bacterial strains in metagenomic datasets by eigengenome partitioning
REF_ENG_RA_010	2016	Journal of Electronic Materials	USA	MICHAEL K. CONNORS, JAMAL E. MILLSAPP, & GEORGE W. TURNER	Dielectric Coating Thermal Stabilization During GaAs-Based Laser Fabrication for Improved Device Yield
REF_ENG_RA_011	2017	Nature Energy	UK	Declan Conway, Carole Dalin, Willem A. Landman & Timothy J. Osborn	Hydropower plans in eastern and southern Africa increase risk of concurrent climate-related electricity supply disruption
REF_ENG_RA_012	2016	International Journal of Refractory Metals and Hard Materials	USA	Jeffrey J. Swab & Jared C. Wright	Application of ASTM C1421 to WC-Co fracture toughness measurement
REF_ENG_RA_013	2015	IEEE Transactions on Fuzzy Systems	USA	Scott Dick	On Pythagorean and Complex Fuzzy Set Operations

REF_ENG_RA_014	2017	Composites: Part A	Ireland	M. Flanagan		Permeability of carbon fibre PEEK composites for cryogenic storage tanks of future space launchers
REF_ENG_RA_015	2016	Information Fusion	USA	Jemin George		A finite point process approach to multi-target localization using transient measurements
REF_ENG_RA_016	2017	Canadian Metallurgical Quarterly	Australia	David Grimsey	E.	Key factors affecting nickel recovery during the segregation of laterite ores
REF_ENG_RA_017	2015	Journal of Manufacturing Science and Engineering	USA	Steven Hansen	R.	Impact Welding of Aluminum Alloys 6061 and 5052 by Vaporizing Foil Actuators: Heat-Affected Zone Size and Peel Strength
REF_ENG_RA_018	2015	Computers and Chemical Engineering	USA	Robert Herring III, Mario R. Eden	H.	Evolutionary algorithm for de novo molecular design with multi-dimensional constraints
REF_ENG_RA_019	2017	Journal of Materials Research and Technology	USA	Lawrence Murr	E.	3D metal droplet printing development and advanced materials additive manufacturing
REF_ENG_RA_020	2016	Additive Manufacturing	USA	N.E. Hodge		Experimental Comparison of Residual Stresses for a Thermomechanical Model for the Simulation of Selective Laser Melting
REF_ENG_RA_021	2016	IEEE Transactions on Pattern Analysis and Machine Intelligence	UK	David Hofmeyr	P.	Clustering by Minimum Cut Hyperplanes
REF_ENG_RA_022	2017	Materials Science & Engineering A	USA	T.R. Jacobs, D.K. Matlock & K.O. Findley		Fractographic analysis of anisotropic deformation behavior after tensile testing of pipeline steels at elevated temperatures
REF_ENG_RA_023	2015	Journal of Manufacturing Science and Engineering	USA	Sagil James & Murali Sundaram	M.	Modeling of Material Removal Rate in Vibration Assisted Nano Impact-Machining by Loose Abrasives
REF_ENG_RA_024	2017	Materials Science & Engineering A	USA	Jayme S. Keist & Todd A. Palmer	A.	Development of strength-hardness relationships in additively manufactured titanium alloys
REF_ENG_RA_025	2015	Journal of Manufacturing Science and Engineering	USA	James Magarian, Robert White & Douglas Matson	N. D. M.	Real-Time Acoustic and Pressure Characterization of Two-Phase Flow for Quality Control of Expanded Polystyrene Injection Molding Processes
REF_ENG_RA_026	2017	Chemical Engineering Journal	USA	Katherine Manz, Kimberly Carter	E. E.	Investigating the effects of heat activated persulfate on the degradation of furfural, a component of hydraulic fracturing fluid chemical additives

REF_ENG_RA_027	2017	IEEE Transactions On Components, Packaging And Manufacturing Technology	USA	Scott McCann	Experimental and Theoretical Assessment of Thin Glass Substrate for Low Warpage
REF_ENG_RA_028	2017	Nature Energy	USA	Dev Millstein	The climate and air-quality benefits of wind and solar power in the United States
REF_ENG_RA_029	2017	International Journal of Machine Tools and Manufacture	UK	Jonathon Mitchell-Smith	Energy distribution modulation by mechanical design for electrochemical jet processing techniques
REF_ENG_RA_030	2017	Materials Science & Engineering A	USA	Philip J. Noell	Growth of Preexisting Abnormal Grains in Molybdenum under Static and Dynamic Conditions
REF_ENG_RA_031	2016	Information Fusion	Australia	Andrew W. Palmer	Applying Gaussian distributed constraints to Gaussian distributed variables
REF_ENG_RA_032	2016	Energy & Environmental Science	UK	Giles Richardson	Can slow-moving ions explain hysteresis in the current-voltage curves of perovskite solar cells?†
REF_ENG_RA_033	2016	Additive Manufacturing	USA	J. Robbins	An efficient and scalable approach for generating topologically optimized cellular structures for additive manufacturing
REF_ENG_RA_034	2017	Materials Science & Engineering A	USA	Andrew L. Robertson	Microscale fracture mechanisms of a Cr ₃ C ₂ -NiCr HVOF coating
REF_ENG_RA_035	2017	Composites: Part A	USA	Cecily A. Ryan	Methodology to assess end-of-life anaerobic biodegradation kinetics and methane production potential for composite materials
REF_ENG_RA_036	2016	International Journal of Metalcasting	Canada	A. M. Samuel and F. H. Samuel	New Method Of Eutectic Silicon Modification In Cast Al-Si Alloys
REF_ENG_RA_037	2015	Computers and Chemical Engineering	UK	James S. Shepherd	Mathematical modelling of the pre-oxidation of auranium carbide fuel pellet
REF_ENG_RA_038	2017	Computers and Electrical Engineering	UK	Michael Short	Eligible earliest deadline first: Server-based scheduling for master-slave industrial wireless networks
REF_ENG_RA_039	2016	Materials Science and Engineering R	Australia	Glenn C. Sneddon	Transmission Kikuchi diffraction in a scanning electron microscope: A review
REF_ENG_RA_040	2015	International Journal of Data Mining, Modelling and Management	USA	Michelle M. Thompson	Public participation GIS and neighbourhood recovery: using community mapping for economic development
REF_ENG_RA_041	2015	International Society for Microbial Ecology	USA	Sarah J. Spencer	Massively parallel sequencing of single cells by epicPCR links functional genes with phylogenetic markers

REF_ENG_RA_042	2017	Materials Science & Engineering A	UK	M.J. Starink	Dislocation versus grain boundary strengthening in SPD processed metals: non-causal relation between grain size and strength of deformed polycrystals
REF_ENG_RA_043	2017	Nature Energy	USA	Leah C. Stokes	Renewable energy policy design and framing influence public support in the United States
REF_ENG_RA_044	2017	Nature Energy	USA	Benjamin Strom	Intracycle angular velocity control of cross-flow turbines
REF_ENG_RA_045	2017	Materials Science & Engineering A	USA	J. Telesman	Relationship between unusual high-temperature fatigue crack growth threshold behavior in superalloys and sudden failure mode transitions
REF_ENG_RA_046	2017	Journal of Materials in Civil Engineering	USA	Jason Weiss	Influence of Mechanically Induced Cracking on Chloride Ingress in Concrete
REF_ENG_RA_047	2017	Canadian Journal of Civil Engineering	Canada	Steven Wood	Regulations governing the operation of longer combination vehicles in Canada
REF_ENG_RA_048	2016	Energy & Environmental Science	Canada	Andrew G. Wright	Hexamethyl-p-terphenyl poly(benzimidazolium): a universal hydroxide-conducting polymer for energy conversion devices†
REF_ENG_RA_049	2015	IEEE Transactions on Fuzzy Systems	USA	Ronald R. Yager	Firing Fuzzy Rules with Measure Type Inputs
REF_ENG_RA_050	2015	Information Fusion	USA	Ronald R. Yager	Modeling Multi-Criteria Objective Functions Using Fuzzy Measures
REF_ENG_RA_051	2017	Applied Energy	UK	A.M. Zealand	Effect of feeding frequency and organic loading rate on biomethane production in the anaerobic digestion of rice straw
REF_ENG_RA_052	2016	Advances in Structural Engineering	Australia	David J Heath	Shaking table study of a brick veneer house subjected to blast vibrations
REF_ENG_RA_053	2017	Advances in Structural Engineering	Australia	Christoph Adam	Quick assessment of high-speed railway bridges based on a non-dimensional parameter representation
REF_ENG_RA_054	2016	Building Environment and	Ireland	Oliver Kinnane	Evaluation of passive ventilation provision in domestic housing retrofit
REF_ENG_RA_055	2016	Building Environment and	USA	Daniel Baseley	Hyperspectral analysis for standoff detection of dimethyl methylphosphonate on building materials
REF_ENG_RA_056	2016	Building Services Engineering Research & Technology	UK	David Johnston	Quantifying the aggregate thermal performance of UK holiday homes

REF_ENG_RA_057	2017	Building Services Engineering Research & Technology	UK	James Parker	Accounting for refrigeration heat exchange in energy performance simulations of large food retail buildings
REF_ENG_RA_058	2017	Building Services Engineering Research & Technology	UK	Anne Stafford	An exploration of load-shifting potential in real in-situ heat-pump/gas-boiler hybrid systems
REF_ENG_RA_059	2017	Coastal Engineering	UK	A. S.J. Foster	An experimentally validated approach for evaluating tsunami inundation forces on rectangular buildings
REF_ENG_RA_060	2017	Coastal Engineering	UK	Mark A. Davidson	Annual prediction of shoreline erosion and subsequent recovery
REF_ENG_RA_061	2016	Concurrent Engineering: Research and Applications	USA	James Mathieson	A protocol for modeling and tracking engineering design process through structural complexity metrics applied against communication networks
REF_ENG_RA_062	2017	Engineering, Construction and Architectural Management	New Zealand	Kathryn Davies	Making friends with Frankenstein: hybrid practice in BIM
REF_ENG_RA_063	2017	Engineering, Construction and Architectural Management	Australia	George Denny-Smith	Integrating Indigenous enterprises into the Australian construction industry
REF_ENG_RA_064	2017	Engineering, Construction and Architectural Management	UK	Barry James Gledson	The adoption of 4D BIM in the UK construction industry: An Innovation Diffusion approach
REF_ENG_RA_065	2016	Hydrometallurgy	Australia	Peter Smith	Reactions of lime under high temperature Bayer digestion conditions
REF_ENG_RA_066	2016	Hydrometallurgy	Australia	G. Riley	Effect of titanium species on the precipitation of boehmite under digestion conditions
REF_ENG_RA_067	2015	Nature Nanotechnology	USA	B. H. Davison	Opto-nanomechanical spectroscopic material characterization
REF_ENG_RA_068	2016	Research in Engineering Design	USA	Christine A. Toh & Scarlett R. Miller	Choosing creativity: the role of individual risk and ambiguity aversion on creative concept selection in engineering design
REF_ENG_RA_069	2016	Rock Mechanics and Rock Engineering	Canada	Laura Brown	Identification of Stress Change Within a Rock Mass Through Apparent Stress of Local Seismic Events
REF_ENG_RA_070	2017	Big Data & Society	USA	Steve Mann	Big Data is a big lie without little data: Humanistic intelligence as a human right
REF_ENG_RA_071	2016	Big Data Analytics	UK	Kevin Swingler	Structure discovery in mixed order hyper networks

Appendix 5. List of Medicine Research Articles in the Reference Corpus

Discipline	Year	Journal	Country	Authors	Title
REF_MED_RA_001	2016	Physiotherapy	UK	Nicola Hancock	J. Towards Upright Pedalling to drive recovery in people who cannot walk in the first weeks after stroke: movement patterns and measurement
REF_MED_RA_002	2017	Clinical Neurophysiology	UK	Melanie Fleming	K. The effect of transcranial direct current stimulation on motor sequence learning and upper limb function after stroke
REF_MED_RA_003	2016	Journal of Physiotherapy	Australia	Lauren Miller	No difference between two types of exercise after proximal phalangeal fracture fixation: a randomised trial
REF_MED_RA_004	2016	Journal of Physiotherapy	Australia	Anne Harrison	L. Exercise improves glycaemic control in women diagnosed with gestational diabetes mellitus: a systematic review
REF_MED_RA_005	2016	Journal of Physiotherapy	Australia	Sonia Bird	Primary contact physiotherapy services reduce waiting and treatment times for patients presenting with musculoskeletal conditions in Australian emergency departments: an observational study
REF_MED_RA_006	2016	The Journal of Emergency Medicine	USA	Kelsey Miller	A. Factors Associated with Misplaced Endotracheal Tubes During Intubation In Pediatric Patients
REF_MED_RA_007	2016	American Journal of Physiology	USA	Lillianne Wright	H. A class of their own: exploring the nondeacetylase roles of class IIa HDACs in cardiovascular disease
REF_MED_RA_008	2016	American Journal of Physiology	USA	Daniel Machin	R. Exercise-induced brachial artery blood flow and vascular function is impaired in systemic sclerosis
REF_MED_RA_009	2015	The Annual Review of Physiology	USA	David Mount	B. The Molecular Physiology of Uric Acid Homeostasis
REF_MED_RA_010	2015	Oral and Maxillofacial Surgery Clinics of North America	USA	George Deeb	R. Soft Tissue Grafting Around Teeth and Implants

REF_MED_RA_011	2017	European Journal of Internal Medicine	Ireland	John McCabe	J.	Deprivation status and the hospital costs of an emergency medical admission
REF_MED_RA_012	2015	Experimental Physiology	UK	Susan Wray		Insights from physiology into myometrial function and dysfunction
REF_MED_RA_013	2015	Experimental Physiology	UK	Matthew Frise	C.	The pulmonary vasculature – lessons from Tibetans and from rare diseases of oxygen sensing
REF_MED_RA_014	2017	British Journal of Oral and Maxillofacial Surgery	UK	R.J. McGalliard		Ophthalmic outcomes of fractured zygomas
REF_MED_RA_015	2016	European Journal of Applied Physiology	UK	David Giles		Validity of the Polar V800 heart rate monitor to measure RR intervals at rest
REF_MED_RA_016	2016	European Journal of Clinical Microbiology & Infectious Diseases	UK	S. Goldenberg	D.	The impact of the introduction of fidaxomicin on the management of Clostridium difficile infection in seven NHS secondary care hospitals in England: a series of local service evaluations
REF_MED_RA_017	2016	Journal of Enzyme Inhibition and Medicinal Chemistry	USA	Neil McIntyre	R.	Inactivation of peptidylglycine a-hydroxylating monooxygenase by cinnamic acid analogs
REF_MED_RA_018	2017	Journal of Clinical Oncology	USA	Leora Horn		Nivolumab Versus Docetaxel in Previously Treated Patients With Advanced Non–Small-Cell Lung Cancer: Two-Year Outcomes From Two Randomized, Open-Label, Phase III Trials (CheckMate 017 and CheckMate 057)
REF_MED_RA_019	2017	Journal of Clinical Oncology	USA	Katherine Thornton		Modernizing Clinical Trial Eligibility: Recommendations of the American Society of Clinical Oncology–Friends of Cancer Research Minimum Age Working Group
REF_MED_RA_020	2016	The Journal of Experimental Medicine	USA	Andrew Getahun		Continuous inhibitory signaling by both SHP-1 and SHIP-1 pathways is required to maintain unresponsiveness of anergic B cells

REF_MED_RA_021	2016	The Journal of Experimental Medicine	USA	Shaun Jackson	W.	B cell IFN- γ receptor signaling promotes autoimmune germinal centers via cell-intrinsic induction of BCL-6
REF_MED_RA_022	2016	The Journal of General Physiology	USA	Franklin Mullins	M.	Orai1 pore residues control CRAC channel inactivation independently of calmodulin
REF_MED_RA_023	2016	The Journal of General Physiology	USA	Richard Lewis	S.	The inactivation domain of STIM1 is functionally coupled with the Orai1 pore to enable Ca ²⁺ -dependent inactivation
REF_MED_RA_024	2017	The Journal of Physiology	USA	David F. Wilson		Oxidative phosphorylation: regulation and role in cellular and tissue metabolism
REF_MED_RA_025	2017	The Journal of Emergency Medicine	USA	Mark Mason	D.	Anticoagulated Trauma Patients: A Level I Trauma Center's Response To A Growing Geriatric Population
REF_MED_RA_026	2015	The new England journal of medicine	USA	John C. Byrd		Acalabrutinib (ACP-196) in Relapsed Chronic Lymphocytic Leukemia
REF_MED_RA_027	2016	The new England journal of medicine	USA	Andrew Roberts	W.	Targeting BCL2 with Venetoclax in Relapsed Chronic Lymphocytic Leukemia
REF_MED_RA_028	2015	Comprehensive Physiology	USA	Stephen Cannon	C.	Channelopathies of skeletal muscle excitability
REF_MED_RA_029	2016	Obstetrics and Gynecology Clinics of North America	USA	Elizabeth Bonney	A.	Immune regulation in pregnancy: a matter of perspective?
REF_MED_RA_030	2017	Physiotherapy	USA	Richard Bohannon	W.	Functional Reach of Older Adults: Normative Reference Values Based on New and Published Data
REF_MED_RA_031	2017	The Journal of Emergency Medicine	USA	Sean Murphy	M.	A Randomized Controlled Trial of a Citywide Emergency Department Care Coordination Program to Reduce Prescription Opioid Related Visits: An Economic Evaluation
REF_MED_RA_032	2017	Clinical Neurophysiology	USA	Jeremy Harper		Testing the effects of adolescent alcohol use on adult conflict-related theta dynamics

REF_MED_RA_033	2017	Nature Medicine	USA	E. Kaitlynn Allen	SNP-mediated disruption of CTCF binding at the IFITM3 promoter is associated with severe influenza risk in humans
REF_MED_RA_034	2017	Journal of Medicinal Chemistry	USA	Andrew J. Robles	Structure–Activity Relationships of New Natural Product-Based Diaryloxazoles with Selective Activity against Androgen Receptor-Positive Breast Cancer Cells
REF_MED_RA_035	2016	Journal of Orthopaedic Trauma	Canada	Joseph Westgeest	Factors Associated with Development of Nonunion or Delayed Healing After an Open Long Bone Fracture: A Prospective Cohort Study of 736 Subjects
REF_MED_RA_036	2015	Applied Nursing Research	USA	Nicole Adams	A review of Yellow Dirt: A Poisoned Land and the Betrayal of the Navajos
REF_MED_RA_037	2017	The Lancet Child & Adolescent Health	UK	Clare Murray	Diagnosis of asthma in symptomatic children based on measures of lung function: an analysis of data from a population-based birth cohort study
REF_MED_RA_038	2015	BMC Ophthalmology	Australia	Christopher J. Layton	Diabetic levels of glucose increase cellular reducing equivalents but reduce survival in three models of 661W photoreceptor-like cell injury
REF_MED_RA_039	2015	British Dental Journal	UK	F. J. T. Burke	A practice-based assessment of patients' knowledge of dental materials
REF_MED_RA_040	2016	American Journal of Clinical Dermatology	USA	Linda Stein Gold	Moderate and Severe Inflammatory Acne Vulgaris Effectively Treated with Single-Agent Therapy by a New Fixed-Dose Combination Adapalene 0.3 %/Benzoyl Peroxide 2.5 % Gel: A Randomized, Double-Blind, Parallel-Group, Controlled Study

REF_MED_RA_041	2017	Physiotherapy	UK	Kay Stevenson	Development and delivery of a physiotherapist-led exercise intervention in a randomised controlled trial for subacromial impingement syndrome (the SUPPORT trial)
REF_MED_RA_042	2016	European Journal of Applied Physiology	UK	Conor Taylor	W. Exercise duration-matched interval and continuous sprint cycling induce similar increases in AMPK phosphorylation, PGC-1 α and VEGF mRNA expression in trained individuals
REF_MED_RA_043	2017	European Journal of Clinical Microbiology & Infectious Diseases	UK	C. E. Berry	Is there a relationship between the presence of the binary toxin genes in Clostridium difficile strains and the severity of C. difficile infection (CDI)?
REF_MED_RA_044	2016	JAMA Psychiatry	UK	Josephine Mollon	Psychotic Experiences and Neuropsychological Functioning in a Population-based Sample
REF_MED_RA_045	2016	JAMA Psychiatry	New Zealand	Kate M. Scott	Association of Mental Disorders With Subsequent Chronic Physical Conditions World Mental Health Surveys From 17 Countries
REF_MED_RA_046	2016	Journal of Orthopaedics, Trauma and Rehabilitation	UK	Carter Thomas & Jefferies James G.	Bilateral Simultaneous Neck of Femur Fractures Arising from a Simple Mechanical Fall: A Case Report for Guidance on Safe Surgical Management
REF_MED_RA_047	2017	BJS Open	New Zealand	L. Clarke	Impact of restenting for recurrent colonic obstruction due to tumour ingrowth
REF_MED_RA_048	2016	Clinical Chemistry	USA	G. Terrance Walker	Analytical Performance of Multiplexed Screening Test for 10 Antibiotic Resistance Genes from Perianal Swab Samples

REF_MED_RA_049	2015	Allergy, Asthma & Clinical Immunology	Canada	Anne K. Ellis	Clinical validation of controlled grass pollen challenge in the Environmental Exposure Unit (EEU)
REF_MED_RA_050	2016	Allergy, Asthma & Clinical Immunology	USA	Bryan L. Love	Antibiotic prescription and food allergy in young children
REF_MED_RA_051	2015	Applied Nursing Research	Australia	Joanne Siffleet	Delivering best care and maintaining emotional wellbeing in the intensive care unit: the perspective of experienced nurses
REF_MED_RA_052	2015	Applied Nursing Research	USA	Carolyn Miller Reilly	Single subject design: Use of time series analyses in a small cohort to understand adherence with a prescribed fluid restriction

Appendix 6. Skewness Kurtosis Values of Amount of Coordination for Subcorpora

Measure	Corpus	Discipline	Skewness		Kurtosis	
			Statistic	Std. Error	Statistic	Std. Error
Coordinate Phrases per Clause (CP/C)	MCWP	TR_COM	1.332	.794	2.280	1.587
		TR_EAS	-.185	.597	-1.016	1.154
		TR_EDU	.406	.340	.124	.668
		TR_ENG	.646	.257	-.335	.508
		TR_MED	.640	.325	.451	.639
	Reference Corpus	REF_COM	-.067	.913	1.794	2.000
		REF_EAS	.259	.637	.146	1.232
		REF_EDU	.726	.409	-.036	.798
		REF_ENG	.858	.285	.589	.563
		REF_MED	.263	.365	-.195	.717
Coordinate Phrases per T-unit (CP/T)	MCWP	TR_COM	.513	.794	-1.126	1.587
		TR_EAS	-.208	.597	-1.322	1.154
		TR_EDU	.427	.340	.126	.688
		TR_ENG	.718	.257	.023	.508
		TR_MED	.633	.325	.212	.639
	Reference Corpus	REF_COM	-1.765	.913	3.395	2.000
		REF_EAS	.331	.637	-.330	1.232
		REF_EDU	.704	.409	-.070	.798
		REF_ENG	.544	.285	.000	.563
		REF_MED	.578	.365	-.404	.717
T-units per Sentence (T/S)	MCWP	TR_COM	.852	.794	1.304	1.587
		TR_EAS	.649	.597	-.597	1.154
		TR_EDU	.217	.340	-.751	.668
		TR_ENG	.391	.257	-.054	.508
		TR_MED	-.008	.325	.077	.639
	Reference Corpus	REF_COM	1.404	.913	1.265	2.000
		REF_EAS	.366	.637	-1.383	1.232
		REF_EDU	.588	.409	-.185	.798
		REF_ENG	.622	.285	.865	.563
		REF_MED	-.324	.365	.544	.717

Appendix 7. Skewness Kurtosis Values of Amount of Subordination for the Subcorpora

Measure	Corpus	Discipline	Skewness		Kurtosis	
			Statistic	Std. Error	Statistic	Std. Error
Clauses per T-unit (C/T)	MCWP	TR_COM	-.570	.794	-1.901	1.587
		TR_EAS	.991	.597	.105	1.154
		TR_EDU	.132	.340	-.427	.668
		TR_ENG	.529	.257	.097	.508
		TR_MED	.896	.325	.160	.639
	Reference Corpus	REF_COM	.602	.913	-2.312	2.000
		REF_EAS	-.032	.637	-.856	1.232
		REF_EDU	.528	.409	.070	.798
		REF_ENG	.532	.285	.833	.563
		REF_MED	.390	.365	.372	.717
Complex T-units per T-unit (CT/T)	MCWP	TR_COM	-.155	.794	-1.980	1.587
		TR_EAS	-.273	.597	-.434	1.154
		TR_EDU	-.185	.340	-.423	.668
		TR_ENG	.283	.257	.197	.508
		TR_MED	.955	.325	.810	.639
	Reference Corpus	REF_COM	.824	.913	-.638	2.000
		REF_EAS	-.147	.637	-1.010	1.232
		REF_EDU	.818	.409	.287	.798
		REF_ENG	-.181	.285	.523	.563
		REF_MED	.314	.365	-.089	.717
Dependent Clauses per Clause (DC/C)	MCWP	TR_COM	-.576	.794	-1.277	1.587
		TR_EAS	.099	.597	-.858	1.154
		TR_EDU	-.424	.40	.241	.668
		TR_ENG	.358	.257	.098	.508
		TR_MED	.624	.325	-.248	.639
	Reference Corpus	REF_COM	.614	.913	-2.809	2.000
		REF_EAS	.026	.637	-1.625	1.232
		REF_EDU	.207	.409	-.558	.798
		REF_ENG	-.547	.285	.686	.563
		REF_MED	.062	.365	-.440	.717
Dependent Clauses per T-unit (DC/T)	MCWP	TR_COM	-.775	.794	-1.302	1.587
		TR_EAS	.435	.597	-.791	1.154
		TR_EDU	.076	.340	-.531	.668
		TR_ENG	.943	.257	.997	.508
		TR_MED	.799	.325	-.354	.639
	Reference Corpus	REF_COM	.706	.913	-2.692	2.000
		REF_EAS	.072	.637	-1.628	1.232
		REF_EDU	-.100	.409	-.855	.798
		REF_ENG	.350	.285	.653	.563
		REF_MED	.688	.365	.615	.717

Appendix 8. Skewness Kurtosis Values of Degree of Phrasal Sophistication for the Subcorpora

Measure	Corpus	Discipline	Skewness		Kurtosis	
			Statistic	Std. Error	Statistic	Std. Error
Complex Nominals per Clause (CN/C)	MCWP	TR_COM	-.271	.794	.626	1.587
		TR_EAS	.945	.597	.825	1.154
		TR_EDU	.374	.340	-.328	.668
		TR_ENG	-.136	.257	-.430	.508
		TR_MED	.154	.325	.026	.639
	Reference Corpus	REF_COM	-1.189	.913	.101	2.000
		REF_EAS	-.177	.637	-.459	1.232
		REF_EDU	.212	.409	.179	.798
		REF_ENG	-.303	.285	.103	.563
		REF_MED	.471	.365	-.148	.717
Complex Nominals per T- unit (CN/T)	MCWP	TR_COM	.839	.794	.088	1.587
		TR_EAS	.485	.597	-1.188	1.154
		TR_EDU	.435	.340	.052	.668
		TR_ENG	-.171	.257	-.552	.508
		TR_MED	.330	.325	-.491	.639
	Reference Corpus	REF_COM	-.486	.913	-1.867	2.000
		REF_EAS	-.163	.637	.078	1.232
		REF_EDU	.650	.409	.145	.798
		REF_ENG	.116	.285	-.366	.563
		REF_MED	-.194	.65	-.388	.717
Verb Phrases per T-unit (VP/T)	MCWP	TR_COM	-.346	.794	1.161	1.587
		TR_EAS	.161	.597	-.934	1.154
		TR_EDU	-.478	.340	.153	.668
		TR_ENG	.720	.257	.098	.508
		TR_MED	.286	.325	.690	.639
	Reference Corpus	REF_COM	-.412	.913	.245	2.000
		REF_EAS	.248	.637	-.228	1.232
		REF_EDU	.308	.409	-.361	.798
		REF_ENG	.538	.285	.607	.563
		REF_MED	.389	.365	-.526	.717

Appendix 9. Skewness Kurtosis Values of Length of Production Unit

Measure	Corpus	Discipline	Skewness		Kurtosis	
			Statistic	Std. Error	Statistic	Std. Error
Mean Length of Sentence (MLS)	MCWP	TR_COM	.675	.794	-.838	1.587
		TR_EAS	.281	.794	-.838	1.587
		TR_EDU	.281	.597	-1.295	1.154
		TR_ENG	.629	.257	-.120	.508
		TR_MED	.389	.325	.027	.639
	Reference Corpus	REF_COM	-1.619	.913	2.834	2.000
		REF_EAS	-.516	.637	-.659	1.232
		REF_EDU	.917	.409	-.218	.798
		REF_ENG	.511	.285	-.028	.563
		REF_MED	.282	.365	-.787	.717
Mean Length of T-Unit (MLT)	MCWP	TR_COM	.232	.794	-.873	1.587
		TR_EAS	.818	.597	-1.014	1.154
		TR_EDU	.219	.340	.471	.668
		TR_ENG	.756	.257	.618	.508
		TR_MED	.331	.325	.509	.639
	Reference Corpus	REF_COM	-.513	.913	-3.139	2.000
		REF_EAS	-.531	.637	-1.161	1.232
		REF_EDU	.799	.409	-.172	.798
		REF_ENG	.521	.285	.094	.563
		REF_MED	-.047	.365	-.956	.717
Mean Length of Clause (MLC)	MCWP	TR_COM	2.351	.794	5.747	1.587
		TR_EAS	.056	.597	-1.024	1.154
		TR_EDU	.431	.340	.825	.668
		TR_ENG	.737	.257	.090	.508
		TR_MED	.089	.325	.483	.639
	Reference Corpus	REF_COM	.464	.913	-2.470	2.00
		REF_EAS	.193	.637	-.452	1.232
		REF_EDU	.916	.409	-.118	.798
		REF_ENG	.150	.285	-.218	.563
		REF_MED	.069	.365	-.450	.717

Appendix 10. Skewness Kurtosis Values of Overall Sentence Complexity for the Subcorpora

Measure	Corpus	Discipline	Skewness		Kurtosis	
			Statistic	Std. Error	Statistic	Std. Error
Clauses per Sentence (C/S)	MCWP	TR_COM	.291	.794	-2.383	1.587
		TR_EAS	.875	.597	-.617	1.54
		TR_EDU	.148	.340	-.359	.668
		TR_ENG	.610	.257	.184	.508
		TR_MED	.851	.325	1.741	.639
		REF_COM	.742	.913	-.800	2.000
	Reference Corpus	REF_EAS	-.136	.637	-1.484	1.232
		REF_EDU	.403	.409	-.907	.798
		REF_ENG	.644	.285	.324	.563
		REF_MED	.348	.365	-.500	.717

Appendix 11. Skewness Kurtosis Values of Amount of Coordination

Measure	Corpus	Skewness		Kurtosis	
		Statistic	Std. Error	Statistic	Std. Error
Coordinate Phrases per Clause (CP/C)	MCWP	.561	.167	-.133	.333
	Reference Corpus	.641	.190	.018	.378
Coordinate Phrases per T-unit (CP/T)	MCWP	.640	.167	.114	.333
	Reference Corpus	.614	.190	.185	.378
T-units per Sentence (T/S)	MCWP	.145	.167	-.113	.333
	Reference Corpus	.350	.190	.144	.378



Appendix 12. Skewness Kurtosis Values of Amount of Subordination

Measure	Corpus	Skewness		Kurtosis	
		Statistic	Std. Error	Statistic	Std. Error
Clauses per T-unit (C/T)	MCWP	.859	.167	.655	.333
	Reference Corpus	.406	.190	-.013	.378
Complex T-units per T-unit (CT/T)	MCWP	.402	.167	-.118	.333
	Reference Corpus	.036	.190	.056	.378
Dependent Clauses per Clause (DC/C)	MCWP	.428	.167	.226	.333
	Reference Corpus	-.069	.190	-.012	.378
Dependent Clauses per T-unit (DC/T)	MCWP	.824	.167	.355	.333
	Reference Corpus	.192	.190	-.525	.378



Appendix 13. Skewness Kurtosis Values of Degree of Phrasal Sophistication

Measure	Corpus	Skewness		Kurtosis	
		Statistic	Std. Error	Statistic	Std. Error
Complex Nominals per Clause (CN/C)	MCWP	.007	.167	.043	.333
	Reference Corpus	.563	.190	.991	.378
Complex Nominals per T-unit (CN/T)	MCWP	.531	.167	.540	.333
	Reference Corpus	.412	.190	.424	.378
Verb Phrases per T- unit (VP/T)	MCWP	.274	.167	-.243	.333
	Reference Corpus	.068	.190	-.596	.378



Appendix 14. Skewness Kurtosis Values of Length of Production Unit

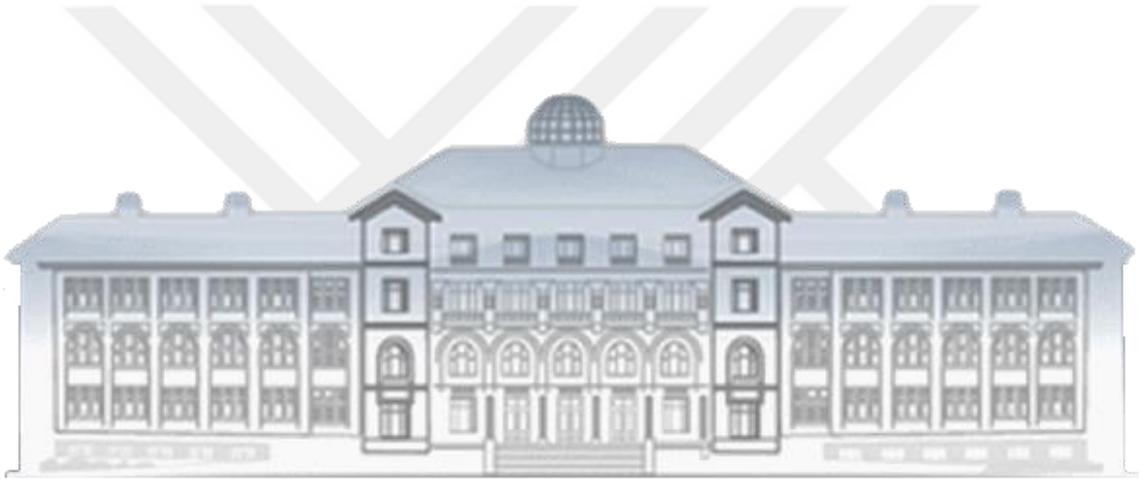
Measure	Corpus	Skewness		Kurtosis	
		Statistic	Std. Error	Statistic	Std. Error
Mean Length of Sentence (MLS)	MCWP	.423	.167	-.278	.333
	Reference Corpus	.677	.190	.104	.378
Mean Length of T-Unit (MLT)	MCWP	.406	.167	.024	.333
	Reference Corpus	.596	.190	-.096	.378
Mean Length of Clause (MLC)	MCWP	.576	.167	.244	.333
	Reference Corpus	.702	.190	.192	.378



Appendix 15. Skewness Kurtosis Values of Overall Sentence Complexity

Measure	Corpus	Skewness		Kurtosis	
		Statistic	Std. Error	Statistic	Std. Error
Clauses per Sentence (C/S)	MCWP	.918	.167	.888	.333
	Reference Corpus	.639	.190	.121	.378





GAZİLİ OLMAK AYRICALIKTIR..