





A virtual conference, under the aegis of the Learner Corpus Association.

Hosted by the Chair of English and Digital Linguistics, Chemnitz University of Technology

22-23-24 October 2025

BOOK OF ABSTRACTS

Editors: Cansu Akan, Sepideh Javdani Esfahani, Sasha Coelho, Mina Raeisi Nafchi, Mehrdad Vasheghani Farahani, Christina Sanchez-Stockhammer

Version: 01 ©LCRGradConf25 Organizing Committee, Chemnitz, Germany, 2025

Conference website: https://lcrgrad2025.tu-chemnitz.de

Table of Contents

ARNER CORPUS ASSOCIATION	5
RGrad 2025 CONFERENCE	6
RGrad 2025 CONFERENCE ORGANIZERS	7
IENTIFIC COMMITTEE	9
YNOTE PRESENTATIONS	. 10
The EMI Corpus of student academic writing: Addressing challenges in corpus design and data collection in a large-scale corpus development project Dana Gablasova (Lancaster University)	. 10
Patterns in corpora and in society Michaela Mahlberg (FAU Friedrich-Alexander-Universität Erlangen- Nürnberg)	. 11
Purposes and patterns in academic texts Randi Reppen (Northern Arizona University (NAU))	. 12
DRKSHOPS	. 13
Publishing during, as part of and from the doctorate Achilleas Kostoulas & Richard Fay (University of Thessaly & The University of Manchester)	. 13
Getting started with Open Science Elen Le Foll (University of Cologne)	. 14
Networking for researchers: A strategic approach John Kluempers	. 15
The lucky mindset - How to attract serendipitous opportunities in your career Anne Schreiter (GSO'Guidance, Skills & Opportunities for Researchers)	
ESENTATIONS	. 17
Combining corpus analysis and machine learning to predict the CEFR level of Estonian learner tex Kais Allkivi (Tallinn University School of Digital Technologies)	
Learner factors in Data-Driven Learning: A meta-analysis of population and treatment variables Reem Alojaimi (Lancaster University, King Saud University)	. 18
Uncovering the patterns beneath: Detecting suspicious essays in learner corpora by the analysis of keystrokes Ahood Al Sawar & Nicolas Ballier (Université Paris Cité)	
Longitudinal development of syntactic complexity in Chinese ESL MA students Liwen Bing (Universion of Birmingham)	-
Capturing underlying patterns in the acquisition of attribution strategies in L2 German Aylin Braunewell (Justus Liebig University Giessen)	. 22
Pauses in L2 writing: A study of keystroke logging data from the Process Corpus of English in Education Amanda Chng Yi, Anna-Katharina Scholz & Christina Sanchez-Stockhammer (Chemnitz	
University of Technology)	. 24

The impact of age on the acquisition of scalar implicatures, tense and aspect in L2 Spanish spoken production by native English students Nebojsa Damnjanovic (University of Belgrade)
Modal verb-argument construction development in Asian L2 English learners: A cross-register, corpus-based analysis Dilay Candan & Daniel Dixon (Georgia State University)
Analyzing collocational errors among Russian learners of Persian: Evidence from Learner Corpus Research Mehrdad Vasheghani Farahani& Zahra Haghighi Naseri (Leipzig University& Iranian Ministry of Education)
IceLC: A new corpus and new methods for analyzing linguistic features in Icelandic L2 writing Isidora Glišić (University of Iceland)
Turkish Sign Language Learner Corpus: TİD ∟2 -Corpus Yasemin Güçlütürk (Ankara University)
Al-based writing evaluation system for L1-specific graph description task Darya Kharlamova (National Research University Higher School of Economics)
Profiling L2 learner proficiency for AI-supported writing feedback: A corpus-based study using ICNALE Hyunhwa Kim (Georgia State University)
Toward dual-level modeling of conjunction: Developing a position-aware framework for L1/L2 writing Jose Lema-Alarcon (Universidad de las Fuerzas Armadas ESPE, University of Exeter)40
Development of noun phrase complexity in L2 spoken and written Chinese Yilei Li (University of Arizona)
Study of the various types of tasks in the compilation of the longitudinal spoken multi-L1 (Slavic) learner corpus of L2 Italian: pros and cons Kristýna Lorenzová (Masaryk University)
Applying learner corpus tools in the classroom: A data-driven approach to collocation development in EFL writing Bahareh Malmir (Boğaziçi University)44
Behind the scenes of linguistic accuracy: Investigating the impact of automated writing evaluation on L2 academic theses Katharina Maschke (Chemnitz University of Technology)
Investigating second language pragmatic competence across proficiency level and first language, using a multimodal corpus Gerard O'Hanlon (Mary Immaculate College)47
Do subpart frequency and phrase frequency affect articulatory durations of multi-word expressions produced by L2 English learners? Yi Qi & Anna Siyanova-Chanturia (Victoria University of Wellington) 49
Introducing an index for analysis of grammatical diversity in writing Christian Holmberg Sjöling & Taehyeong Kim (Luleå University of Technology & Northern Arizona University)50
The influence of grammatical complexity features on grade in Swedish high-stakes EFL exams Christian Holmberg Sjöling (Luleå University of Technology)
Communicative strategies in oral and written Romanian: A case study Isabella Şinca (University of Bucharest)
Analysis of metaphor production in Finnish L1 learner English Renata Turunen (University of Inland Norway)54
Validating automated measures of phraseological competence in L2 speaking: Evidence from human

P	OSTER PRESENTATIONS57
	LC-meta: A core metadata schema for L2 data Jennifer-Carmen Frey ¹ , Alexander König ² , Hubert Naets ³ , Egon W. Stemle ¹ and Magali Paquot ³ (¹ Institute for Applied Linguistics, Eurac Research, ² CLARIN ERIC, ³ UCLouvain)
	Separating agreement and assignment of gender marking in L2 Spanish writing Gabriela Sanchez (University of Texas at Arlington)
	Complexity development of Intermediate German learners of English: A longitudinal corpus analysis Philine Metzger (Philipps University Marburg)
	Comparative analysis of the modal verb 'could' among Indonesian and Japanese EFL learners Nida Ghaziyah (Kanazawa University)

LEARNER CORPUS ASSOCIATION

The Learner Corpus Association is an international association which aims to promote the field of learner corpus research and provide an interdisciplinary forum for all the researchers and professionals who are actively involved in the field or simply want to know more about it.

LCA supports the compilation of learner corpora (i.e. electronic collections of written and/or spoken language produced by foreign/second language learners) in a wide range of languages and the design of innovative methods and tools to analyze them. It seeks to link up learner corpus research to second language acquisition theory, first language acquisition theory and linguistic theory in general and to promote applications in fields including foreign language teaching, language testing and natural language processing (automated scoring, spell- and grammar-checking, L1 identification).

The founding members of the Association are Gaëtanelle Gilquin, Sylviane Granger, Fanny Meunier and Magali Paquot, all based at the Centre for English Corpus Linguistics, Université catholique de Louvain (Belgium).

LCA was officially launched at the Learner Corpus Research Conference (Bergen, 27-29 September 2013).

LCA Board (September 2025)

President: Sylvie De Cock, UCLouvain

Vice-Presidents: Bill Crawford, Northern Arizona University

Agnieszka Leńko-Szymańska, University of Warsaw

General Secretary: Tove Larsson, Northern Arizona University

Treasurer: Monika Kavalir, University of Ljubljana

Conference Officer: Agnieszka Leńko-Szymańska, University of Warsaw

Webmaster: Elen Le Foll, University of Cologne

Website: https://www.learnercorpusassociation.org/

LCRGrad 2025 CONFERENCE

The Pattern Beneath

Revealing the structures that shape and define learner language

The central theme of LCRGrad25 conference is "The Pattern Beneath". This theme celebrates the unique role of learner corpus research in uncovering the underlying structures and patterns of learner language through LCR. It emphasizes the field's potential to provide insights into second language acquisition, linguistic development, and the intricacies of language use in educational contexts. The theme also highlights how learner corpora serve as powerful tools for identifying trends, testing hypotheses, and advancing our understanding of learner-specific challenges and strategies.

In line with the conference theme, we welcome contributions addressing the following areas:

- 1. **Advancing Methodologies in Learner Corpus Research**: Cutting-edge methods in corpus design, annotation, and analysis.
- 2. **Technological Innovations in Learner Corpus Development**: Utilizing digital tools, apps, and virtual environments for corpus creation and use.
- 3. **Al and Data Science in Learner Corpus Analysis**: Harnessing AI, machine learning, NLP, and big data for learner corpus insights.
- 4. **Multimodal and Multilingual Learner Corpora**: Integrating diverse languages and modalities (e.g., text, audio, video) into corpus research.
- 5. Interdisciplinary Approaches to Learner Corpus Research: Contributions spanning linguistics, psychology, education, computational sciences, and more.
- 6. **Application of Learner Corpus Insights in Education**: Practical outcomes for teaching, curriculum design, materials development, and assessment.
- 7. Sociocultural, Cross-Cultural and Contextual Perspectives on Learner Data: Investigations into learner-specific, sociocultural, and contextual influences.
- 8. Corpus-Based Insights into Language Acquisition and Development: Advancing understanding of first and second language acquisition and linguistic development.
- 9. Ethical and Inclusive Practices in Learner Corpus Research: Addressing inclusivity, accessibility, and ethical considerations in corpus studies.
- 10. Future Directions in Learner Corpus Research: Examining emerging trends, technologies, and challenges in the field.

LCRGrad 2025 CONFERENCE ORGANIZERS

Conference Chair: Cansu Akan (English and Digital Linguistics at TU Chemnitz)

Conference Committee:

Christina Sanchez-Stockhammer (Chair of English and Digital Linguistics at TU Chemnitz)

Sasha Coelho (English and Digital Linguistics at TU Chemnitz)

Sepideh Javdani Esfahani (English and Digital Linguistics at TU Chemnitz)

Mina Raeisi Nafchi (TESOL at TU Chemnitz)

Mehrdad Vasheghani Farahani (English and Digital Linguistics at TU Chemnitz)

Host Department: English and Digital Linguistics at Chemnitz University of Technology)

English and Digital Linguistics at Chemnitz University of Technology comprises all aspects related to the English language, including its structure, use, acquisition, learning, processing, and investigation. We aim to cover the full range of English linguistics in our teaching, with a particular focus on digital research methods that allow for empirical investigation, such as corpus linguistics.

Corpus-based studies form a central part of our work, including projects that integrate learner corpora and digital methods to support English learning and teaching from an early age. These interdisciplinary projects employ principles of second language acquisition, corpus linguistics, teacher training, and multimodality in digital tools. Our work with corpora includes *TransGrimm*, a parallel corpus of Grimms' fairy tales and their English translations, and we are also involved in *PROCEED* (Process Corpus of English in Education), an international project in collaboration with the University of Louvain, which focuses on L2 writing processes.

In addition to corpus-based research, current projects include eye-tracking and skin-conductance studies on linguistic and stylistic processing, as well as the development of a linguistic naming system for archaeological human remains in collaboration with archaeologists. We are also actively engaged in creating digital tools, such as *Bridge of Knowledge VR*, a virtual-reality learning app developed with the Leibniz Supercomputing

Centre, and LeDiT (Learning with Digital Testimonies), which builds on work at LMU Munich to make Holocaust survivor testimonies accessible worldwide through subtitling and interactive dialogue. On a broader scale, we contribute to the international *LITHME* (Language in the Human-Machine Era) network, investigating how emerging technologies, including human-robot interaction, may reshape language use and learning.

At the Chair of English and Digital Linguistics, we place great importance on communicating our research to the wider public. To this end, we use press releases, video tutorials on YouTube, and a podcast. In cooperation with the University of Zaragoza, we have also developed best practices for *LingComm*—science communication in linguistics. The podcast *Linguistics Behind the Scenes* aims to present linguistic research to a general audience in a fun and engaging way.

Our goal is to enable students to reflect on language from a variety of perspectives. To this end, we provide intensive tutoring and a broad, solid education that prepares students for work in multiple domains. We also encourage and support international exchanges, such as a semester abroad at one of our many partner universities, and we welcome colleagues from around the world.

Website: https://www.tu-chemnitz.de/phil/english/sections/edling/

SCIENTIFIC COMMITTEE

Acknowledgment and thanks are given to the Scientific and Evaluation Committee

Akira Murakami University of Birmingham Agnieszka Leńko-Szymańska University of Warsaw

Aivars Glaznieks Eurac Research

Anke Lüdeling

Andrea Abel

Cansu Akan

Christina Sanchez-Stockhammer

Humboldt University of Berlin

Free University of Bolzano/Bozen

Chemnitz University of Technology

Erik Castello University of Padua

Gaëtanelle Gilquin Catholic University of Louvain

Hilde Hasselgård University of Oslo Jennifer Carmen Frey Eurac Research

Jennifer Thewissen

Katherine Ackerley

Luciana Forti

María Belén Díez-Bedmar

Pascual Pérez-Paredes

University of Antwerp
University of Padua
University of Perugia
Universidad de Jaén
Universidad de Murcia

Stefanie Wulff Bremen University

Susan Nacey Inland Norway University of Applied Sciences

Sylviane Granger Catholic University of Louvain
Sylvie De Cock Catholic University of Louvain
Tove Larsson Northern Arizona University

KEYNOTE PRESENTATIONS

The EMI Corpus of student academic writing: Addressing challenges in corpus design and data collection in a large-scale corpus development project Dana Gablasova (Lancaster University)

English-medium instruction (EMI) is a major pedagogical trend, reflecting and shaping the status of English as a global language. In EMI contexts, the ability to use academic English – the language through which academic subjects are taught and assessed – is crucial to academic success. However, many recent studies of EMI in higher education (e.g. Galloway et al., 2020; Macaro et al., 2018) report that students experience language-related difficulties affecting their academic performance. These studies rely largely on self-reported data from questionnaires, surveys and interviews, with surprisingly little evidence available about actual language use in EMI contexts. To address this, we have developed the EMI Corpus – a corpus representing student writing across multiple EMI contexts, with data collected at eight universities in China, Italy, Thailand, Austria and the UK (Gablasova et al., 2024). Currently, the corpus contains over 4.5 million words from over 2,000 student texts representing three disciplinary areas: i) Social Sciences and Humanities (e.g. History, Education), ii) Science and Technology (e.g. Computer Science, Engineering), and iii) Business and Management.

The talk first offers the description of the EMI Corpus and its future uses in English for Academic Purposes (EAP) research and practice. We will then focus on the process and challenges in the design and development of a large-scale corpus that involves multi-site data creation, and we will discuss the strategies adopted to address these challenges in our project.

Short Bio:

Dana Gablasova is Senior Lecturer in the Department of Linguistics and English Language, Lancaster University. Her research interests are in the development of L2 corpora and application of corpus methods to language learning, teaching and assessment, with particular attention to formulaic language and pragmatics. She is the director of the Corpus for Schools project.

Patterns in corpora and in society *Michaela Mahlberg* (FAU Friedrich-Alexander-Universität Erlangen-Nürnberg)

The identification and description of patterns is fundamental to corpus linguistic research. The ubiquity of language patterns is further highlighted by the AI boom that has put language into the spotlight like never before. From a methodological point of view, concordance software has always played a central role in the study of patterns. From a theoretical point of view, lexico-grammatical patterns that are identifiable on the basis of concordance data need to be interpreted in their social contexts, including the contexts of language learners. In this talk, I will look at both these perspectives. For the practical analysis of patterns, I will illustrate how the new FlexiConc (Dykes et al. 2025) functionality in the web application CLiC (Mahlberg et al. 2020) can support the reading of concordance lines. To look at the bigger picture of patterns in society, I will discuss case studies of patterns in fictional and non-fictional texts. To bring the arguments together, I will propose some initial steps towards a big-picture corpus linguistic approach that can capture the patterns that shape our society, our culture and our reality.

References

Dykes, N.; Evert, S.; Mahlberg, M.; Pipersperki, A. (2025). FlexiConc 0.1.18 https://pypi.org/project/FlexiConc/

Mahlberg, M., Stockwell, P., Wiegand, V. and Lentin, J. (2020) CLiC 2.1. Corpus Linguistics in Context, available at: <u>clic.bham.ac.uk</u> & <u>clic-fiction.com</u>

Short Bio:

Michaela Mahlberg is Professor of Digital Humanities at FAU Erlangen Nürnberg, Germany and honorary Professor of Corpus Linguistics at the University of Birmingham, UK. Her research focuses on the relationship of language and reality. She is editor of the *International Journal of Corpus Linguistics*. With her team, she has developed the CLiC web app for the digital reading of fiction, and she is the host of the "Life and Language" podcast.

Purposes and patterns in academic texts Randi Reppen (Northern Arizona University (NAU))

Writing academic texts is a task that requires control of many different aspects of language. The writer has to control much more than simply sentence structure at a local level, they must also be able to control patterns of language to construct texts to accomplish different rhetorical purposes. This presentation will look at some of the patterns of language use that are found in L1 and L2 English writing across a range of different task types. We'll explore some of the resources that can be used to identify different patterns of language use and how those patterns function to accomplish different rhetorical goals. Many examples will be presented and implications for instruction for both English L1 and L2 academic writers will be provided.

Short Bio:

Randi Reppen is Professor Emerita of Applied Linguistics and TESL at Northern Arizona University (NAU). She has a keen interest in using corpus research to inform language teaching and to develop better language teaching materials. Randi has given presentations in over 25 countries and directed NAU's Program in Intensive English for 11 years. She is the lead author of the multi-level ELT textbook series *Grammar and Beyond with Academic Writing* (2020) and her recent publications have appeared in the *Journal of English for Academic Purposes*, *English Language and Linguistics*, and the *International Journal of Learner Corpus Linguistics*. Randi enjoys many outdoor activities, especially, biking, hiking, and Nordic skiing.

WORKSHOPS

Publishing during, as part of and from the doctorate Achilleas Kostoulas & Richard Fay (University of Thessaly & The University of Manchester)

For many doctoral students, publishing is a rite of passage from their studies into early researcherhood. This is not just a question of gaining academic legitimacy, but also —for many early career researchers— a process of academic 'becoming', i.e., gradually developing a distinctive authorial voice and researcher identity. However, this formative, and potentially very rewarding process, is often fraught with feelings of uncertainty, self-doubt, and institutionalized or tacit pressure to publish. This workshop aims to address some of these challenges by familiarizing participants with the publication process. In the first part of the workshop, the facilitators will discuss the peer review process and discuss common issues that reviewers identify, drawing on examples from their editorial experience. The second half of the workshop will foreground the lived experiences of early career researchers making their first steps into publishing. Working in groups, participants will engage with narratives that describe early publication experiences, and extract insights that are relevant to their own contexts. By the end of the workshop, it is expected that participants will have a clearer understanding of the practicalities involved in academic publishing, as well as stronger confidence from learning about their peers' experiences.

Short Bio:

Achilleas Kostoulas is an applied linguist and language teacher educator at the University of Thessaly in Greece. He holds a PhD and an MA in TESOL from The University of Manchester, and a BA in English Studies from the University of Athens. He has published extensively about various aspects of language teaching and learning, often seen through a Complex Systems Theory perspective, and has served as editor to the Q1 journal Studies in Second Language Learning and Teaching. Some of his recent publications include the research monograph The Intentional Dynamics of TESOL (2021; with Juup Stelma) and the edited volume Doctoral Study and Getting Published (2025; with Richard Fay). More information about his work is available at www.achilleaskostoulas.com

Richard Fay is a Reader in Education (TESOL and Intercultural Education) at The University of Manchester, where he received his PhD in Education (2004). He is a critical intercultural educator working across disciplines, and his research (often collaborative, interdisciplinary, and multilingual in character) has focused topics including of language education, language teacher education, music education, researcher education, and global mental health. He is the co-editor, with Achilleas Kostoulas, of the recently published *Doctoral Study and Getting Published* (2025; Emerald Publishing).

Getting started with Open Science *Elen Le Foll (University of Cologne)*

In this informal coffee + chat session, I will provide a brief introduction to different aspects of Open Science from sharing materials, data, and code, to creating Open Educational Resources (OERs), and doing replications in Learner Corpus Research (LCR). We will cover the basics of pre-registration, FAIR data, reproducibility, and open access publishing. I will show concrete examples from my PhD project and from other published LCR studies. Together, we will discuss tips, incentives, and barriers to doing Open Science as early-career researchers in LCR.

Short Bio:

Elen Le Foll is a post-doctoral researcher and lecturer in linguistics at the Department of Romance Studies at the University of Cologne. She has a strong interest in quantitative corpus linguistics methods and applications of corpus research to language teaching and learning. She is co-project leader of a project on the role of gender as a prominence feature within the Collaborative Research Center "Prominence in Language". As a keen educator, she enjoys teaching about quantitative methods, R, statistics, data visualisation, critical data literacy, and Open Science practices for linguistics and language education research.

Networking for researchers: A strategic approach *John Kluempers*

With some careful, but not time-consuming, consideration, you can make visits to professional meetings more rewarding. With that in mind, we will see what strategies you can use to return satisfied from your next event.

- Understand what professional meetings are about
- Make yourself visible
- The researcher is a curious person by nature
- Cultural aspects of networking
- Small talk

This talk will also be participatory.

Short Bio:

For 15 years, **John Kluempers**, Ph.D. has helped researchers in Germany, Switzerland and Austria improve their careers. He uses his experience as a radio journalist at Deutsche Welle, Germany's international broadcaster to assist hundreds of Ph.D. candidates and post-docs take their presentations to the next level. Giving presentations in all their forms, e.g., conference talks and poster presentations, is a step. A step to getting in contact with fellow professionals at the various meetings scholars attend. Even if you are not presenting, you want to make new contacts. And that is a further specialty of John: helping junior researchers make the most of their conference visits. He got his doctorate in linguistics at University of California, Los Angeles (UCLA), lived in Berlin before moving to Tokyo, Japan where he taught. John has sought different teaching environments and loves the challenge of working with intelligent and ambitious researchers in all fields.

The lucky mindset - How to attract serendipitous opportunities in your career *Anne Schreiter (GSO*- Guidance, Skills & Opportunities for Researchers)*

Chance plays a crucial role in building a successful research career. While excellent research is the foundation, not all outstanding researchers secure academic positions - even when that is their goal. There is no guaranteed formula for becoming a professor, but you can create conditions that make serendipitous opportunities more likely.

In this talk, I will explore how communication and key aspects of leadership can help you position yourself for these fortunate coincidences. I will also offer insights into developing an authentic career strategy that works for academia & other sectors. The talk will be interactive, and there will be time for your questions.

Short Bio:

Anne Schreiter is Executive Director of GSO* - Guidance, Skills & Opportunities for Researchers, an independent nonprofit in Berlin offering tailored Career Guidance, targeted Skill Development, and Opportunities through innovative programs for PhD-level researchers. She studied Social and Business Communications and Sociology in Berlin and Chinese Language in Nanjing, China. After receiving her PhD in Organizational Sociology from the University of St. Gallen, she spent a year as postdoctoral researcher at UC Berkeley, USA. She then transitioned from academia to the nonprofit sector.

PRESENTATIONS

Combining corpus analysis and machine learning to predict the CEFR level of Estonian learner texts Kais Allkivi (Tallinn University School of Digital Technologies)

This study addressed the challenge of developing transparent machine learning models for language testing. Learner corpora with human-graded texts provide empirical material for identifying language-specific features of the CEFR levels and building text assessment tools. However, there is a lack of studies explicitly combining these two aspects. The novelty of presented research lies in careful selection of linguistic variables used for classification, only including such that can be associated with increasing complexity and accuracy in writing, rather than task-based factors. This allows to interpret the assessment results, enabling the learner and teacher to understand what is being graded and how it relates to writing proficiency.

In Estonia, automated language assessment has recently gained more relevance due to the planned digitisation of state exams, including Estonian as L1 and L2. The study focused on classifying Estonian L2 proficiency examination writings according to their level (A2–C1). Four categories of linguistic features were applied for level-to-level comparative analysis and training prediction models: lexical features – measures of lexical diversity and sophistication; grammatical features – frequencies of parts-of-speech and grammatical categories of nominals and verbs; surface features – text complexity measures related to word, sentence, and text length; error features – frequencies of spelling and grammatical errors as detected by correction tools.

First, the training set was analysed to detect reliable distinguishers of the proficiency levels. Then, classification models were trained for each feature category. The different features were also combined in a unified model. Using only pre-selected relevant features reduced variation in classifying different text types (e.g., argumentative writings, informal and semiformal letters). The best classifiers achieved test accuracy of around 90%. Their generalizability was also tested on a separate exam dataset. The models have been integrated to the writing evaluation module of an Estonian language learning environment.

Keywords: automated writing assessment, L2 proficiency, supervised machine learning, text classification, CEFR

Learner factors in Data-Driven Learning: A meta-analysis of population and treatment variables Reem Alojaimi (Lancaster University, King Saud University)

Data-driven learning (DDL) is an approach to language teaching providing corpora access to language learners (Johns, 1991). Its effectiveness in (second) language acquisition has been demonstrated in the literature (Boulton & Vyatkina, 2021; O'Keeffe, 2021; Pérez-Paredes, 2022) and synthesised by recent meta-analyses (Boulton & Cobb, 2017; Cobb & Boulton, 2015; Lee et al., 2019; Mizumoto & Chujo, 2015; Ueno & Takeuchi, 2023; Yoon & Lee, 2024), attracting constant research attention in the field. This study expands on previous syntheses by investigating the effect of learner-related factors (population variables) and intervention-related factors (treatment variables) in the classroom context and their influence on effect size. It employs two methods: a simple meta-analysis and a multiple meta-regression analysis using a multi-model approach. 97 studies met the inclusion criteria, and 125 unique between-subjects samples with over 6,000 participants were included.

The results show that DDL is an effective teaching approach with a large effect (g=0.95). Among the population variables, region and sample size, in order, predict a larger effect size. On the other hand, institution type and students' language proficiency are not important predictors. Among the treatment variables, corpus type, linguistic target of intervention and DDL interaction type (i.e. direct vs indirect DDL) are, in order, the most important predictors. Conversely, although experiment duration is not considered an important predictor of the effectiveness of DDL, further analysis shows that longer interventions could lead to more substantial learning outcomes. Further analysis also revealed the versatility of DDL: Despite its application being largely exclusive to research and universities (Timmis &Templeton, 2023, p.420), the findings prove that it is also promising in high schools as well. Overall, the meta-analysis suggests that medium- and long-term interventions with direct access to public corpora in medium-sized groups of learners are the most effective conditions in DDL for better learning gains. Implications for future studies and meta-analyses will be discussed.

Keywords: data-driven learning, meta-analysis, multi-model inference

References

- Boulton, A., & Cobb, T. (2017). Corpus Use in Language Learning: A Meta-Analysis. *Language Learning*, 67(2), 348–393. https://doi.org/10.1111/lang.12224
- Boulton, A., & Vyatkina, N. (2021). *Thirty years of data-driven learning: Taking stock and charting new directions over time*. http://hdl.handle.net/10125/73450
- Cobb, T., & Boulton, A. (2015). Classroom applications of corpus analysis. In D. B. & R. Reppen (Ed.), *Cambridge Handbook of Corpus Linguistics* (pp. 478–497). Cambridge University Press. https://doi.org/10.1017/CBO9781139764377.027
- Johns, T. (1991). "Should you be persuaded": Two samples of data-driven learning materials. *ELR Journal*, *4*, 1–16.

- Lee, H., Warschauer, M., & Lee, J. H. (2019). The Effects of Corpus Use on Second Language Vocabulary Learning: A Multilevel Meta-analysis. *Applied Linguistics*, *40*(5), 721–753. https://doi.org/10.1093/applin/amy012
- Mizumoto A., & Chujo K. (2015). A Meta-analysis of Data-driven Learning Approach in the Japanese EFL Classroom. *English Corpus Studies=英語コーパス研究*, 22, 1–18.
- O'Keeffe, A. (2021). Data-driven learning a call for a broader research gaze. *Language Teaching*, 54(2), 259–272. https://doi.org/10.1017/S0261444820000245
- Pérez-Paredes, P. (2022). A systematic review of the uses and spread of corpora and datadriven learning in CALL research during 2011–2015. *Computer Assisted Language Learning*, 35(1–2), 36–61. https://doi.org/10.1080/09588221.2019.1667832
- Ueno, S., & Takeuchi, O. (2023). Effective corpus use in second language learning: A metaanalytic approach. *Applied Corpus Linguistics*, *3*(3), 100076. https://doi.org/10.1016/j.acorp.2023.100076
- Yoon, K.-H., & Lee, D. J. (2024). A Meta-Analysis of Data-Driven Learning (DDL) in EFL/ESL Settings. *Brain, Digital, & Learning, 14*(2), 283–304. https://doi.org/10.31216/BDL.20240017

Uncovering the patterns beneath: Detecting suspicious essays in learner corpora by the analysis of keystrokes Ahood Al Sawar & Nicolas Ballier (Université Paris Cité)

The emergence of ChatGPT has made it necessary to detect learner texts which are Algenerated. We report a current experiment on a learner corpus collected on-line with keylogging data. Our study provides a novel methodological contribution to learner corpus research (LCR) by concentrating on the writing process and analysing typing behaviours, unlike traditional LCR essays, where the focus is on the final text. We employed the King's College London & University Paris Cité Keys (KUPA-KEYS) dataset, which consists of keystroke data collected from 1,006 participants writing essays in English. Our findings suggest practical applications for enhancing corpus annotation and academic integrity.

By comparing the key press latency for the copy-text task and the essay-writing task with an Anderson-Darling test, (Velentzas et al., 2024) showed that 20% of the texts had a high p-value for this test, suggesting that "the cognitive effort of the participant for the essay-writing task is not much different than for the copy-text task". We attempted to manually annotate the suspicious texts using a visualisation method to differentiate between regular essays and suspicious essays.

We compared the writing patterns of the copy task and of the writing task with the visualisation of two longitudinal representations of keylogging, with time as an X dimension and number of characters as a Y dimension. The existence of textual bursts and pauses was used for our identification of suspicious essays, and when in doubt, we used a replay tool to detect pauses in typing patterns and observe textual bursts. Analysing the "pattern beneath" the text allows us to reveal hidden insights into the writing process. Our replaying tool uncovers unusual pausing patterns potentially indicating suspicious behaviour, like when long pauses are followed by long textual bursts or very long textual bursts with no pause nor repair. Using these criteria, the first annotator found 68 suspicious texts, and inter-rater agreement is on its way for the second annotator.

Keywords: keystroke logging, learner corpus, suspicious writing detection, textual bursts, pauses patterns, writing process

References

Velentzas, G., Caines, A., Borgo, R., Pacquetet, E., Hamilton, C., Arnold, T., ... & Yannakoudakis, H. (2024, May). Logging keystrokes in writing by English learners. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)* (pp. 10725-10746).

Longitudinal development of syntactic complexity in Chinese ESL MA students *Liwen Bing (University of Birmingham)*

In L2 writing development, syntactic complexity (SC) is a vital construct, as its development is essential to a second language learner's overall progress in mastering the target language (Ortega, 2003). Previous research views SC as a multidimensional construct (Bulté & Housen, 2014; Norris & Ortega, 2009). However, the interpretation of measures and the dimensions they represent remain controversial. In addition, current research on SC in L2 writing development largely relies on language assessment data (e.g., TOEFL) or focuses on EFL contexts, with limited studies using disciplinary corpora in ESL settings.

To address these gaps, this study examined the longitudinal development of SC measures in disciplinary assignments written by Chinese ESL MA students over one academic year at a UK university. The dataset comprised 124 academic assignments produced by 62 participants in the field of language education, with one assignment collected per semester from each student. These texts formed a self-built corpus consisting of a comparable genre, written under untimed conditions.

Exploratory factor analysis was conducted on 14 SC measures computed using the L2 Syntactic Complexity Analyzer (Lu, 2010) to identify underlying constructs. Thirteen of the measures loaded onto two factors, while T-units per sentence (T/S) was excluded due to low loading. The first factor, clausal sophistication, included six measures such as clauses per T-unit (C/T) and dependent clauses per T-unit (DC/T). The second factor, phrasal sophistication, comprised seven measures, including mean length of clause (MLC) and complex nominals per clause (CN/C). Together, the two factors explained over 80% of the variance. Mixed-effects regression models were used to predict factor scores over time, controlling for individual differences. Results showed significant increases in both clausal and phrasal sophistication. Further analysis of individual measures confirmed growth in dependent clauses and complex nominals, possibly reflecting participants' current proficiency and developmental stage in academic writing.

Keywords: L2 writing development, learner corpus, syntactic complexity

References

- Bulté, B., & Housen, A. (2014). Conceptualizing and measuring short-term changes in L2 writing complexity. Journal of Second Language Writing, 26, 42–65.
- Lu, X. (2010). Automatic analysis of syntactic complexity in second language writing. International Journal of Corpus Linguistics, 15(4), 474-496.
- 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.
- 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.

Capturing underlying patterns in the acquisition of attribution strategies in L2 German Aylin Braunewell (Justus Liebig University Giessen)

This paper investigates how newly migrated learners of L2 German convey attributional content using strategies alternative to noun group-internal devices. Although noun group expansion is central to the "language of schooling" (Schleppegrell 2001) and has been studied in L2 German (e.g. Siekmeyer 2013), less is known about how learners express the same informational content without embedded modifiers.

The study draws on a two-year longitudinal corpus of 15 students in German preparatory classes (cf. Schlauch et al. in prep.), which prepare newly arrived children and adolescents for regular classes. It consists of oral and written retellings of six identical picture stories in varied communicative settings to capture register-specific development (cf. Wiese 2020; Lüdeling 2022).

After an initial focus on precursor structures of relative clauses, I adopted a functional, form-open perspective: each expanded noun group (e.g. das Kind aus dem roten Team 'the child from the red team') is annotated for its specifying information, here: the child's membership in the red team. These annotations yield content units, such as "team membership" in the example above. These units then were traced across retellings of the same picture story to identify instances where similar information is conveyed by alternative means—for example, by repetition and sequencing (Das rote Team ist traurig. Der Junge ist traurig. 'The red team is sad. The boy is sad.') rather than by employing formal attribution. This inductive analysis addresses two questions: Which alternative attribution strategies do learners employ, and do setting or modality (oral vs. written) affect their choice? In a subsequent quantitative phase, I will test whether the relative frequency of these strategies shifts with increasing proficiency.

By changing from a narrow, form-centric model to a broader functional analysis, this research clarifies how referential specificity is achieved in early L2 German and contributes to models of academic language development.

Keywords: L2 German, precursor structures, relative clause, attribution, newly immigrated students

References

Lüdeling, A., Alexiadou, A., Adli, A., Donhauser, K., Dreyer, M., Egg, M., Feulner, A. H., Gagarina, N., Hock, W., Jannedy, S., Kammerzell, F., Knoeferle, P., Krause, T., Krifka, M., Kutscher, S., Lütke, B., McFadden, T., Meyer, R., Mooshammer, C., ... Zeige, L. E. (2022). Register: Language Users' Knowledge of Situational-Functional Variation. *Register Aspects of Language in Situation*, 1, 1–58. https://doi.org/10.18452/24901

Schlauch, J., Braunewell, A., Gamper, J. (in prep.). *Das Seiteneinsteiger:innenkorpus SeiKo*. Schleppegrell, M. J. (2001). Linguistic Features of the Language of Schooling. *Linguistics and Education*, *12*(4), 431–459. https://doi.org/dt6q9c

- Siekmeyer, A. (2013). Sprachlicher Ausbau in gesprochenen und geschriebenen Texten. Zum Gebrauch komplexer Nominalphrasen als Merkmale literater Strukturen bei Jugendlichen mit Deutsch als Erst- und Zweitsprache in verschiedenen Schulformen. Dissertation. Universität des Saarlandes. https://publikationen.sulb.uni-saarland.de/handle/20.500.11880/23682
- Wiese, H. (2020). Language Situations: A Method for Capturing Variation within Speakers' Repertoires. In Y. Asahi (Ed.), *Proceedings of Methods XVI: Papers from the sixteenth international: Vol. v.59* (pp. 105–117). Peter Lang.

Pauses in L2 writing: A study of keystroke logging data from the Process Corpus of English in Education Amanda Chng Yi, Anna-Katharina Scholz & Christina Sanchez-Stockhammer (Chemnitz University of Technology)

While corpus linguistics has enabled large-scale analyses of written language, the underlying cognitive mechanisms of writing remain difficult to capture. This challenge is further amplified in second language (L2) writing, where additional cognitive demands influence the writing process. One way to investigate these underlying mechanisms is through the study of pauses, i.e. moments of hesitation that provide insight into cognitive processing during writing. While current works in the field observe pauses within discourse boundaries, little is known about the linear process during which these pauses take place.

This study investigates the pausing behaviour of L2 writers through a relatively novel visualization approach using line graphs that show pauses along a time axis (Wengelin 2006, Luuk Van Waes et al. 2016, Bécotte-Boutin et al. 2019). Drawing on keystroke logging data from the German part of the Process Corpus of English in Education (PROCEED, French component, Gilquin 2022) collected from undergraduate students at Chemnitz University of Technology, the research examines pause frequency, duration, and pause location to explore the extent to which line graphs effectively represent L2 writer's pausing behaviour. A pause threshold of 2000ms was applied following studies of similar nature (Spelman Miller 2000, Wengelin 2006, Elola & Mikulski 2016). Qualitative analysis of the data reveals patterns consistent with previous studies, while the line graphs suggest underlying trends in the pausing behaviour of L2 writers. For instance, recursion was a strategy present in most of the participants. Additionally, a peak in pause duration at the end of the writing process could be observed in more than half of the dataset. When viewed qualitatively, these peaks point to a pattern of proofreading. The relatively small sample size (n=13) remains a limitation.

The findings contribute to existing research on L2 writing processes by demonstrating the potential of observing pauses from an in-depth qualitative perspective. Nonetheless, the present study underscores the need for further research into writing processes, particularly regarding pause durations before lexical retrieval events versus revision, proofreading characteristics, and recursive writing.

Keywords: learner corpus, second-language writing, pausing patterns

References

Bécotte-Boutin, Hélène-Sarah, Gilles Caporossi, Alain Hertz & Christophe Leblay. 2019. Writing and rewriting: The coloured numerical visualization of keystroke logging. In Eva Lindgren & Kirk Sullivan (eds.), *Observing writing: Insights from keystroke logging and handwriting*, 96–120. Leiden: Brill.

Elola, Idoia & Ariana M. Mikulski. 2016. Similar and/or different writing processes? A study of Spanish foreign language and heritage language learners. *Hispania* 99(1). 87–102. Gilquin, Gaëtanelle. 2022. The Process Corpus of English in Education: Going beyond the

- written text. Research in Corpus Linguistics 10(1). 31-44.
- Spelman Miller, Kristyan. 2000. Academic writers on-line: Investigating pausing in the production of text. *Language Teaching Research* 4(2). 123–148.
- Van Waes, Luuk, Mariëlle Leijten, Eva Lindgren & Åsa Wengelin. 2016. Keystroke logging in writing research: Analyzing online writing processes. In Charles A. MacArthur, Steve Graham & Jill Fitzgerald (eds.), *Handbook of writing research*, 410–426. London: The Guilford Press.
- Wengelin, Åsa. 2006. Examining pauses in writing: Theory, methods, and empirical data. In Kirk P. H. Sullivan & Eva Lindgren (eds.), Computer keystroke logging and writing: Methods and applications, 107–130. Amsterdam: Elsevier.

Connecting the dots: Elements, structure, and coherence in deaf learners' written narratives Patrice Clarke (The University of the West Indies)

Studies investigating narrative competence suggest it improves with maturity and have routinely focused on hearing populations and discrete linguistic components rather than on pragmatic and discourse skills (Khan et al. 2016; Yoshinaga-Itano and Downey 1986). The limited research on deaf learners' narrative skills indicates that they generally compose short narratives which contain few cohesive devices, show an absence of or non-adherence to discourse rules, and exhibit difficulties with referencing, organising, and relating ideas (Marschark, Mouradian, and Halas 1994; Arfé and Boscolo 2008; Soares, Garcia de Goulart, and Chiari 2010).

This paper explores the development, quality, and structure (Applebee 1978; Heilmann et al. 2010) of 209 single-picture elicited narratives written by deaf learners aged 13 to 21 who are acquiring English as an additional language. Learners were enrolled in specialised schools for the deaf that adopted a sign bilingual approach involving Jamaican Sign Language and English. Data was subjected to a corpus analysis guided by a functional framework.

The qualitative findings revealed that learners performed at similar levels (Williams and Mayer 2015), mainly producing heaps or sequences and only a few stories that describe routines (40.62%) or narrate positive events (31.25%). Overall, younger learners produced more stories (27) than their older peers (2). The quantitative findings showed that narrative scores were low, with texts scoring 5.32 SD below the subcorpus mean (13.51) on the NSS. This level of underperformance points to disturbed coherence, developmental pragmatic and discourse skills and the absence or underdevelopment of core story elements.

Interestingly, character development (2.51) and referencing (2.74) were the narrative indices with the highest mean scores, and each had a significant main effect (prompt) on narrative performance (p < 0.001). Additionally, referential strategies used by learners create ambiguity and make it challenging to identify and track referents within each text (Arfé and Perondi 2008). Though it is common to use nominal strategies in writing, learners' reliance on nominal devices, namely, (un)determined full noun phrases to (re)introduce characters, disrupts thematic continuity and coherence. These findings underline the importance of linguistic and pragmatic skills in establishing and maintaining textual coherence. The paper concludes by discussing the study's limitations and its implications for test designs and the use of picture stimuli in writing assessment.

Keywords: deaf learners, English as an additional language, narrative macrostructure, picture-elicited narratives, coherence

References

Applebee, Arthur. *The child's concept of story: Ages two to seventeen.* Chicago: University of Chicago Press. 1978.

- Arfé, Barbara & Pietro Boscolo. 2006. Causal coherence in deaf and hearing students' written narratives. *Discourse Processes* 42 (3). https://doi.org/10.1207/s15326950dp4203_2
- Arfé, Barbara & Irene Perondi. 2008. Deaf and hearing students' referential strategies in writing: What referential cohesion tells us about deaf students' literacy development. First Language 28(4). https://doi.org/10.1177/0142723708091043
- Heilmann, John, Jon F. Miller, Ann Nockerts & Claudia Dunaway. 2010. Properties of the narrative scoring scheme using narrative retells in young school-age children. American Journal of Speech-Language Pathology 19. https://doi.org/10.1044/1058-0360(2009/08-0024)
- Khan, Kiren, Mihaiela Gugiu, Laura Justice, Ryan Bowles, Lori Skibbe & Shayne Piasta. 2016.

 Age-related progressions in story structure in young children's narratives. *Journal of Speech, Language, and Hearing Research* 59(6). https://doi.org/10.1044/2016_JSLHR-L-15-0275
- Marschark, Marc, Vera Mouradian, and Margaret Halas. 1994. Discourse rules in the language productions of deaf and hearing children. *Journal of Experimental Child Psychology* 57(1). https://doi.org/10.1006/jecp.1994.1005
- Soares, Alexandra, Bárbara Garcia de Goulart, and Brasilia Chiari. 2010. Narrative competence among hearing-impaired and normal-hearing children: analytical cross-sectional study. *Sao Paulo Medical Journal* 128. doi:10.1590/s1516-31802010000500008.
- Williams, Cheri and Connie Mayer. 2015. Writing in young deaf children. *Review of Educational Research* 85(4). https://doi.org/10.3102/0034654314564882
- Yoshinaga-Itano, Christine & Doris M. Downey. 1986. A hearing-impaired child's acquisition of schemata: Something's missing, *Topics in Language Disorders* 7(1). doi/10.1097/00011363-198612000-00007

The impact of age on the acquisition of scalar implicatures, tense and aspect in L2 Spanish spoken production by native English students Nebojsa Damnjanovic (University of Belgrade)

L2 learners differ in many ways. While L1 acquisition almost guaranteed, L2 learners reach their intended goal through different paths, which gives a distinctly individual tone to their language learning experience. One of the more prominent factors is age. By the same token, the same holds true for their pragmatic competence. The research shall investigate how language learners of different age groups acquire and interpret scalar implicatures (SIs), and whether tense-aspect acquisition supplements or hinders their pragmatic inference. Papafragou et al. (2016) mentioned that "implicatures are components of speaker meaning that constitute an aspect of what is meant in a speaker's utterance without being part of what is said", while Paradis et al. (2008) stated that "without question, developing the ability to use morphology to articulate aspect is one of the most significant challenges for L2 learners" (as cited in Thane [2018]). The topic of interest is the development of SIs in native English speakers learning Spanish as a second language (L2), since "Spanish is widely learned as an L2, and there is an actively developing research literature on its acquisition" (Mitchell et al. 2008). It will focus on their interaction with tense-aspect acquisition across three age groups: beginners (A2, age 13-14/15), intermediate students (B1, age 17-18), and undergraduates (C1, 4th year). Horn scales (lexical sets structured by the level of informativity [the ordering relation of entailment]) and the Aspect Hypothesis will be employed as theoretical frameworks. SPLLOC1 and SPLLOC2 (Spanish Learner Language Oral Corpus - 60 instructed learners, L1-English) corpora shall be utilized, with a small native speaker control sample. Ten students of every age and proficiency group shall be randomly selected for further quantitative analysis (ANOVA, regression, and partial correlation) from each individual task. One of questions of interest shall be do perfective and telic verbs align more often with stronger, monotone implicatures, and inperfective and atelic verbs with weak, non-monotone ones. One of the expected findings is that beginners shall predominantly rely on L1 transfer, because "learners depend on the context of the language they are producing and rely on their interaction with the interlocutor and transferred expressions from the L1" (Bardovi-Harlig, 2013).

Keywords: scalar implicatures, tense-aspect acquisition, L2 Spanish, quantifier scope, pragmatic inference

References

Bardovi-Harlig, K. (2013). Tense and aspect in second language acquisition. In The handbook of second language acquisition (pp. 235–252). Wiley. https://doi.org/10.1002/9781118584347.ch14

Chierchia, G. (2017). Scalar implicatures and their interface with grammar. Annual Review of Linguistics, 3, 245–264. https://doi.org/10.1146/annurev-linguistics-011516-033846

Mitchell, Rosamond & Domínguez, Laura & Arche, Maria J & Myles, Florence & Marsden,

Emma. (2008). SPLLOC: A new database for Spanish second language acquisition research. EUROSLA Yearbook. 8. 287-304. 10.1075/eurosla.8.15smit.

Papafragou, Anna & Skordos, Dimitrios. (2016). Scalar Implicature.

Thane, P. D. (2018). The present state of the Aspect Hypothesis: A critical perspective. Eurasian Journal of Applied Linguistics, 4(2), 261-273. doi: 10.32601/ejal.464173

Modal verb-argument construction development in Asian L2 English learners: A cross-register, corpus-based analysis Dilay Candan & Daniel Dixon (Georgia State University)

Verb-argument constructions (VACs) are fundamental to capturing how learners acquire and use meaning-bearing grammatical patterns in a second language, as they encode relationships between actions and participants as well as learners' ability to combine lexical and syntactic knowledge in contextually appropriate way (Goldberg, 2003). Grounded in usage-based acquisition theory (Ellis & Wulff, 2015), this study aims to extend Römer-Barron's (2024) research on modal VACs by examining their development among L1 Korean, Japanese, and Chinese learners of English, using both spoken and written data from the International Corpus Network of Asian Learners of English (ICNALE). The study aims to investigate developmental patterns targeted in previous research with Indo-European L1s and to explore how modality (spoken vs. written) and register influence VAC use. Further, results will also be used to assess the cross-linguistic generalizability from previous research.

Normalized frequencies (per 1,000 words) of three modal VAC types—nsubj-modal-v, nsubj-modal-v-dobj, and nsubj-modal-v-xcomp-were systematically analyzed using corpus linguistic tools and methods, including the natural language processing tool spaCy for dependency parsing which allows for register-based subcorpus comparisons. Proficiency levels (A2-B2, CEFR) were controlled for, and results were compared within ICNALE and, cautiously, with EFCAMDAT findings. The analysis accounted for differences in task type and register by comparing argumentative essays, spoken monologues, and dialogues, and by referencing frequency patterns in native-speaker corpora where possible. Findings indicate that modal VAC use increases with proficiency, supporting usage-based models, but overall frequencies are lower than in German and Spanish learner corpora, likely due to differences in corpus design, register, and cultural-pragmatic norms regarding modality. For example, interview data suggest that learners from Asian backgrounds may use modal constructions less frequently in spoken interaction due to cultural preferences for indirectness or deference. Written data consistently showed higher modal VAC density than spoken data, reflecting genre-driven expectations for explicit stance-taking. These results highlight that modal VAC frequency is shaped not only by learner competence but also by register, task, and cultural factors, and that over- or underuse may reflect pragmatic adaptation rather than developmental lag.

This study provides empirical evidence for the dynamic, context-sensitive nature of constructional development in L2 English and underscores the need for pragmatics-focused, corpus-informed pedagogy tailored to diverse learner backgrounds and communicative contexts.

Keywords: learner corpus research, modal verb-argument constructions, second language acquisition, corpus-based pedagogy

References

- Ellis, N. C. & Wulff, S. (2015). Usage-based approaches to SLA. In B. VanPatten & J. Williams (Eds.), Second language acquisition research series: Theories in second language acquisition (pp. 75-94). New York: Routledge.
- Goldberg, A.E. (2003). Constructions: A new theoretical approach to language. Trends in Cognitive Science, 7, 219-224.
- Römer-Barron, U. (2024). How do constructions with modal verbs develop in second language learners of English?. Journal of Second Language Studies, 7(2), 199-226.

Analyzing collocational errors among Russian learners of Persian: Evidence from Learner Corpus Research Mehrdad Vasheghani Farahani& Zahra Haghighi Naseri (Leipzig University& Iranian Ministry of Education)

This study explores collocational errors among Russian-speaking learners of Persian—a learner population rarely examined in second language acquisition research. The corpus comprises 72 short essays (100–120 words each) written on general topics such as daily routines, education, and personal experiences, resulting in a learner corpus of 11,340 tokens. The texts were manually annotated using CATMA, following Selinker's interlanguage framework, and analyzed with Sketch Engine to extract collocational patterns. Errors were identified by comparing learner output with native-speaker usage, considering gradient acceptability rather than a strict binary of correct vs. incorrect. Preliminary findings reveal recurring lexical mismatches, particularly in light verb constructions and adjective–noun pairings, often influenced by learners' first language. These results offer insights into developmental interlanguage features and highlight the pedagogical need to address collocational competence explicitly. The study contributes to learner corpus research, second language lexical acquisition, and the design of corpus-informed curricula for Persian language instruction.

Keywords: learner corpus, collocation errors, Russian learners, first and second language

IceLC: A new corpus and new methods for analyzing linguistic features in Icelandic L2 writing Isidora Glišić (University of Iceland)

We introduce IceLC, an Icelandic learner corpus in development designed to support research on second language (L2) writing in Icelandic. The corpus is annotated using SVALA, a specialized tool for learner corpus annotation developed by Språkbanken (Volodina & Matsson, 2019) that ensures anonymity and supports aligned rewrites and error corrections. SVALA enables the construction of parallel corpora by systematically linking original learner texts to their corrected versions at both the word and sentence level.

The data for IceLC was compiled during a 2024 data collection campaign, as part of the initial development of standardized assessment for Icelandic as a second language. The preparatory work involved adapting the Common European Framework of Reference (CEFR) descriptors to Icelandic and designing writing prompts. Learners participated in pilot writing assessments at various CEFR levels, producing texts that include narratives, arguments, and personal or societal descriptions. This variety offers rich insights into vocabulary development and language use across different stages of L2 acquisition. A statistical analysis will be conducted on IceLC to identify grammatical and lexical patterns across proficiency levels. These findings will be compared with data from IceL2EC, an Icelandic L2 error corpus (Ingason et al., 2022), to examine how error patterns and language usage evolve with learner proficiency.

The presentation will also address the challenges of processing learner language with NLP tools originally designed for native Icelandic, as interlanguage characteristics often disrupt tokenization, lemmatization, and part-of-speech tagging. However, the integration of SVALA allows for systematic error classification, enhancing the corpus's value for developing automated writing assessment tools. We underscore the importance of standardized learner corpora for both advancing computational methods and informing teaching, assessment, and curriculum development in Icelandic L2 education.

Keywords: learner corpus, Icelandic as a second language, SVALA, CEFR, skill level assessment

References

- Council of Europe. (2018). Common European Framework of Reference for Languages:

 Learning, Teaching, Assessment. Companion Volume with New Descriptors.

 Strasbourg: Council of Europe Publishing. [https://rm.coe.int/common- European framework-of-reference-for-languages-learning-teaching/16809ea0d4]
- Gilquin, G. (2020). Learner Corpora. In M. Paquot & S.T. Gries (Eds.), *A Practical Handbook of Corpus Linguistics*. Cham, Switzerland: Springer, pp. 283–304.
- Granger, S. (2017). Learner Corpora in Foreign Language Education. In *Learner Corpora in Foreign Language Education*, pp. 427–440. Cham: Springer International Publishing. Hvanndal, G., Þorláksdóttir, H., Jónsdóttir, H., and Þorvaldsdóttir, S. (2025). Þróun

- stöðumats í íslensku sem öðru máli. *Mannamál: vefrit um íslensku og önnur mál.* Aðgengilegt á https://mannamal.is/throun-stodumats-i-islensku-sem-odru- mali/
- Ingason, A., Stefánsdóttir, L., Arnardóttir, Þ., Xu, X., Glišić, I. & Guðmundsdóttir, D. (2022). The Icelandic L2 Error Corpus (IceL2EC) version 1.2. CLARIN-IS.
- Ingólfsdóttir, S., Loftsson, H., Daðason, J., and Bjarnadóttir, K. 2019. Nefnir: A high accuracy lemmatizer for Icelandic. In *Proceedings of the 22nd Nordic Conference on Computational Linguistics*, pages 310–315, Turku, Finland. Linköping University Electronic Press.
- Skills and Labour Market Integration of Immigrants and their Children in Iceland, Working Together for Integration. 2024. OECD Publishing, Paris.
- Steingrímsson, S., Kárason, Ö. and Loftsson, H. 2019. Augmenting a BiLSTM Tagger with a Morphological Lexicon and a Lexical Category Identification Step. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2019)*, pages 1161–1168, Varna, Bulgaria. INCOMA Ltd.
- Volodina, E. & Matsson, A. (Dan Rosén and Mats Wirén) (2019). SVALA: an Annotation Tool for Learner Corpora generating parallel texts. *Learner Corpus Research Conference-2019*. Poland, Warszawa.
- Westhoff, G. (2007). Challenges and opportunities of the CEFR for reimagining foreign language pedagogy. *The Modern Language Journal*, 91(4): 676–679. Aðgengilegt á http://www.jstor.org/stable/4626098

Turkish Sign Language Learner Corpus: TİD_{L2}-Corpus Yasemin Güçlütürk (Ankara University)

Recent developments in second language acquisition (SLA) research have increasingly focused on multimodal and bimodal contexts, highlighting signed languages as legitimate and complex L2 systems (Cormier et al., 2012; Ortega & Ozyurek, 2020). Learner corpora offer powerful empirical tools for uncovering language acquisition patterns and informing pedagogy (Granger, 1998, 2008; Myles, 2005; Alonso-Ramos, 2016). While over 200 spoken L2 learner corpora have been developed (see CECL, 2023), signed language corpora for M2L2 learners are scarce (Schönström, 2015; Mesch & Schönström, 2023). As yet no such resource currently exists for Turkish Sign Language (TİD).

To address this gap, this study presents the design of the TİD_{L2}-Corpus, the first learner corpus of TİD created for hearing adults learning TİD as M2L2. The corpus includes learner productions collected through both naturalistic and elicited tasks. In addition to free conversations and guided dialogues, participants are engaged in structured tasks such as picture descriptions, video narrations, and storytelling. Unique to this corpus is the inclusion of specially designed visual and video stimuli aimed at eliciting classifier constructions. These stimuli were created using AI-supported tools (e.g., Canva, Adobe Firefly) and include widely validated materials from international research (such as, images from the Benjamin Bruening Scope Fieldwork Project, images from Zeshan's Pictorial Contrasts Test, a redrawn version of Meyer's Frog Story (1969), and Pixar Shorts for video-based tasks (Slobin, 1996)). Annotation is conducted using ELAN with over 90 tiers, covering lexical, phonological, syntactic, semantic, and error-related features (Johnston, 2010), and the corpus will be uploaded to the Swedish National Data Service (SND) to promote transparency and reusability.

TİD_{L2}-Corpus offers a uniquely structured dataset that combines natural and elicited productions with classifier-focused visual materials. Its design enables the observation of interlanguage development in a modality-specific context, paving the way for innovative language teaching tools, assessment frameworks, and corpus-based language planning. By foregrounding the visual and grammatical richness of sign language learning, this study sets a foundation for future research in multimodal, inclusive language acquisition.

Keywords: learner corpus, Turkish Sign Language, M2D2 acquisition, corpus linguistics classifier

References

- Alonso-Ramos, M. (2016). Learner corpora and second language acquisition. John Benjamins.
- Callies, M., & Götz, S. (2015). *Learner corpora in language testing and assessment*. John Benjamins.
- Cormier, K., Schembri, A., & Woll, B. (2012). Pronouns and pointing in sign languages. *Lingua*, 122(3), 345–362.

- Ellis, R. (2004). The study of second language acquisition (2nd ed.). Oxford University Press.
- Granger, S. (1998). The computer learner corpus: A versatile new source of data for SLA research. In S. Granger (Ed.), *Learner English on computer* (pp. 3–18). Longman.
- Granger, S. (2008). Learner corpora. In A. Lüdeling & M. Kytö (Eds.), *Corpus linguistics: An international handbook* (Vol. 1, pp. 259–274). Mouton de Gruyter.
- Granger, S. (2015b). Contrastive interlanguage analysis: A reappraisal. *International Journal of Learner Corpus Research*, 1(1), 49–66.
- Johnston, T. (2010). From archive to corpus: Transcription and annotation in the creation of signed language corpora. *International Journal of Corpus Linguistics*, 15(1), 104–129.
- Lillo-Martin, D., Quadros, R. M. de, & Chen Pichler, D. (2022). Second language acquisition of sign languages. In M. Marschark & P. E. Spencer (Eds.), *The Oxford Handbook of Deaf Studies in Language* (2nd ed., pp. 303–322). Oxford University Press.
- Mesch, J., & Schönström, M. (2023). Developing a learner corpus for Swedish Sign Language. *International Journal of Learner Corpus Research*, 9(1), 42–65.
- Myles, F. (2005). Interlanguage corpora and second language acquisition research. Second Language Research, 21(4), 373–391.
- Ortega, L., & Ozyurek, A. (2020). Sign languages and SLA: A natural laboratory for embodied language learning. *Studies in Second Language Acquisition*, *42*(1), 161–175.
- Ortega, L., & Quadros, R. M. de. (2023). Toward a more inclusive SLA: Bimodal bilingualism and sign language acquisition. *Language Learning*, 73(S1), 179–205.
- Schönström, M. (2015). Sign language learner corpora: Opportunities and challenges. In L. Borin & A. Saxena (Eds.), *Approaches to measuring linguistic differences* (pp. 341–358). De Gruyter Mouton.
- Seidlhofer, B. (2002). The shape of things to come? Some basic questions about English as a lingua franca. In M. Hewings (Ed.), *Academic writing in context* (pp. 269–280). Birmingham University Press.
- Slobin, D. I. (1996). From "thought and language" to "thinking for speaking". In J. Gumperz & S. Levinson (Eds.), *Rethinking linguistic relativity* (pp. 70–96). Cambridge University Press.

Al-based writing evaluation system for L1-specific graph description task Darya Kharlamova (National Research University Higher School of Economics)

This ongoing research aims to develop a new automated writing evaluation (AWE) system optimized for L1 Russian learners of English completing IELTS-like graph description tasks. The tool can assist English as a foreign language (EFL) teachers and students with providing feedback on writing tasks (Wang et al. 2020:2). We address shortcomings in existing tools, including tendency for erroneous corrections, challenging teacher-student communication (Benali 2021), and inability to adjust systems to different tasks or types of learners by creating a learner-specific and task-specific system. Additionally, we introduce a comprehensive and user-friendly error annotation scheme inspired by Bryant (2019).

We start with the new error annotation scheme for REALEC corpus (Vinogradova & Lyashevskaya 2022) of learner English texts by native Russian speaking students. We group and restructure the existing error tags into a simpler two-level scheme: correction level (3 tags) and linguistic level (11 tags), with the full tagging scheme resulting from the combination of the two levels, following Bryant (2019). The introduced taxonomy addresses imbalanced error distribution in REALEC. Next, we fine-tune error detector (encoder) and error corrector (decoder) models following Yuan et al. (2021): the former identifies and classifies erroneous spans while the latter provides the corrections. We compare the performance of this system on our corpus against prompting decoder-only language models, cf. (Fan 2023). The resulting most successful pipeline is further reused to generate additional tasks based on student error patterns, inspired by Vinogradova & Login (2021). The effectiveness of both approaches (encoder-decoder-based and prompt-based) is assessed by monitoring performance changes in volunteering students.

Expected results include the openly accessible and scalable AWE system for EFL learners and teachers, and release of restructured REALEC benchmark. We also plan to contribute into the discussion on the effectiveness of AWE design choice under learner-specific and task-specific conditions.

Keywords: English as foreign language, applied machine learning, automated writing evaluation, grammatical error correction, language model fine-tuning

- Benali, A. (2021). The Impact of Using Automated Writing Feedback in ESL/EFL Classroom Contexts. English Language Teaching 14(12):189. doi: 10.5539/elt.v14n12p189.
- Bryant, Ch. (2019). Automatic Annotation of Error Types for Grammatical Error Correction (Publication No 938). [Ph.D. thesis, University of Cambridge, Churchill College]. https://www.cl.cam.ac.uk/techreports/UCAM-CL-TR-938.pdf
- Fan, Y, Jiang, F., Li, P., & Li, H. (2023). GrammarGPT: Exploring Open-Source LLMs for Native Chinese Grammatical Error Correction with Supervised FineTuning. https://arxiv.org/abs/2307.13923

- Vinogradova, O, & Login, N. (2017). The Design of Tests with Multiple Choice Questions Automatically Generated from Essays in a Learner Corpus. SSRN Electronic Journal. doi: 10.2139/ssrn.3087215.
- Vinogradova, O, & Lyashevskaya, O. (2022). Review of Practices of Collecting and Annotating Texts in the Learner Corpus REALEC. Pp. 77–88 in Text, Speech, and Dialogue. Vol. 13502, Lecture Notes in Computer Science, edited by P. Sojka, A. Horák, I. Kopeček, and K. Pala. Cham: Springer International Publishing. https://link.springer.com/chapter/10.1007/978-3-031-16270-1_7
- Wang, E. L., Matsumura, L. C., Correnti, R., Litman, D., Zhang, H., Howe, E., Magooda, A., & Quintana, R. (2020). eRevis(Ing): Students' Revision of Text Evidence Use in an Automated Writing Evaluation System. Assessing Writing 44:100449. doi: 10.1016/j.asw.2020.100449.
- Yuan, Z., Taslimipoor, S., Davis, Ch., & Bryant, Ch. (2021). Multi-Class Grammatical Error Detection for Correction: A Tale of Two Systems. Pp. 8722–36 in Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.687

Profiling L2 learner proficiency for AI-supported writing feedback: A corpus-based study using ICNALE Hyunhwa Kim (Georgia State University)

This study explores how natural language processing (NLP) and generative AI can support second language (L2) writing instruction by profiling learner proficiency using corpus-based methods. Focusing on argumentative essays from the International Corpus Network of Asian Learners of English (ICNALE), the research extracts a range of lexico-grammatical features such as part-of-speech (POS) tag sequences, syntactic n-grams, and dependency patterns—using the spaCy NLP toolkit. These features are analyzed quantitatively across learners' CEFR proficiency levels and L1 backgrounds to uncover systematic patterns of linguistic development. Grounded in corpus-based approaches to academic writing, and drawing on Biber and Gray's (2016) model of grammatical complexity, the study operationalizes L2 proficiency through measurable structural variation in learner texts. By identifying linguistic features that reliably differentiate proficiency levels, the research creates detailed learner profiles that serve as a foundation for adaptive feedback mechanisms. These findings directly inform the design of an Al-assisted feedback system that personalizes writing suggestions based on learners' proficiency profiles. Rather than replacing teacher input, the proposed framework envisions a hybrid feedback model in which Al-generated guidance complements instructor intervention. This ensures both scalability and pedagogical integrity in writing instruction. This research contributes to the growing field of learner corpus analysis by demonstrating how computational tools can bridge the gap between descriptive linguistic patterns and pedagogically meaningful interventions. It offers insights into the integration of learner corpora and AI in applied linguistics, highlighting their combined potential to inform writing assessment, curriculum design, and the future of technology-enhanced language education.

Keywords: learner corpus research, natural language processing, generative AI, L2 writing, writing assessment

References

Biber, D., & Gray, B. (2016). Grammatical complexity in academic English: Linguistic change in writing. Cambridge University Press.

Toward dual-level modeling of conjunction: Developing a position-aware framework for L1/L2 writing *Jose Lema-Alarcon (Universidad de las Fuerzas Armadas ESPE, University of Exeter)*

There is growing interest in using automated tools to identify cohesive lexicogrammatical features in texts produced by L1 and L2 writers. This work-in-progress presentation introduces a position-sensitive algorithm designed to detect conjunctions functioning as cohesive devices across structural (intra-clausal) and non-structural (inter-sentential and paragraph-level) positions. The approach extends Halliday's framework by incorporating a positional dimension into the classification of conjunctions by discourse function.

The algorithm has been tested on three corpora: EFCAMDAT and IELTS (L2 writing), and the GIG corpus (L1 writing). Preliminary findings show that conjunction usage patterns vary systematically with proficiency level and textual position. Statistical analysis and Random Forest classification support these observations, highlighting consistent contrasts in how L1 and L2 writers deploy causal, adversative, and additive conjunctions. This model achieved promising classification accuracy across text levels and contributes to a broader characterization of cohesion in L1 and L2 learner writing.

To facilitate pedagogical and research applications, a prototype Streamlit interface has been created. This web-based tool displays the model's output and enables users to analyze new texts by examining conjunction use at different levels. It is designed to enhance transparency, replicability, and practical insight into cohesion through conjunction analysis. In addition, a standalone downloadable tool for the automatic analysis of conjunction (TAACONJ) will also be presented. TAACONJ is intended for local analysis of larger corpora and allows researchers to extract position-sensitive conjunction indices without relying on an internet connection. Together, these tools support the evaluation of cohesion patterns in learner writing and facilitate large-scale corpus analysis.

Keywords: cohesion, conjunctive devices, sentence position, learner corpus, NLP, non-structural cohesion, Python

Development of noun phrase complexity in L2 spoken and written Chinese *Yilei Li* (University of Arizona)

Noun phrases (NP) are typical features of formal writing (Biber & Gray, 2010), and NP complexity has proven powerful in predicting writing scores (Crossley & McNamara, 2014; Lu & Wu, 2022). It is also expected to increase as learners continue L2 study, compared to the clausal development at the beginning (Biber, Gray, & Poonpon, 2011; Parkinson & Musgrave, 2014). However, NP development in L2 Chinese remains under investigation. Furthermore, how NP complexity in Chinese writing is different from learners' oral production is still unclear. Therefore, this study aims to find the developmental patterns and register differences in NP complexity among L2 Chinese learners in an American university. Data are extracted from the Chinese Learner Corpus in MACAWS (Multilingual Academic Corpus of Assignments: Writing & Speech), covering elementary to advanced learning stages (equivalent to four years/levels of Chinese courses in college). Written data includes narratives and argumentative writings, totaling 286 texts. Spoken data includes presentations and role play, totaling 175 texts. To examine NP complexity in a fine-grained way, all NPs will be identified first and categorized into four types [lexical modifier + noun; phrasal modifier + noun; clausal modifier + noun; multiple modifiers + noun] based on the properties of premodifiers. Then, the Kruskal-Wallis Test will be performed to determine whether significant differences exist among the four levels for each type of NP, and the Mann-Whitney U Test will be run to compare register differences at each level. I expect to depict the progression of L2 Chinese NP complexity in writing and when learners start to develop register awareness in terms of NP use. Findings could contribute to the understanding of how complexity is constructed in Chinese NPs as learners progress in their learning journey.

Keywords: noun phrase complexity, Chinese learner corpus, second language Chinese, linguistic development

Study of the various types of tasks in the compilation of the longitudinal spoken multi-L1 (Slavic) learner corpus of L2 Italian: pros and cons Kristýna Lorenzová (Masaryk University)

The study investigates the advantages and disadvantages of various tasks used to design a longitudinal spoken corpus of university students of Italian as a second language (L2), with levels from A1 to C1 (according to CEFR levels, Council of Europe 2020), whose mother tongue are Slavic languages (Czech, Slovak, Russian and Ukranian). These tasks include semi-structured interviews, role-plays, linguistic autobiographies, personal narrations, approximation tasks, focus groups and coffee talks (Migliorini/Raina 2001; Pallotti et al. 2010; Krueger/Casey 2015).

This protocol of data collection is administered once a year during the three-year bachelor program. The first data collection session is completed, and the second is in progress. I will examine the strengths and limitations of each task, along with their detailed description, based on the criteria regarding task complexity (e.g. time available), task condition (e.g. numbers of participants) and task difficulty (e.g. participant characteristics) on the target language production (Robinson and Gilabert 2007).

Learning languages represents a complex process that occurs over time; therefore, it is believed that any statement about learning is meaningful only if made from a longitudinal perspective (Ortega/Iberri-Shea 2005; Robinson/Ellis 2008). Although still limited, longitudinal learner corpora are a valuable resource for linguistic research, providing insights into interlanguage development (Granger et al. 2015). To optimize data collection and ensure comprehensive coverage of linguistic phenomena, it's crucial to understand the pros and cons of various oral tasks included.

A combination of tasks, directed towards the meaning rather than the linguistic form (Ellis 2003), is expected to provide varied interlanguage samples, offering a broader perspective than typically seen in acquisition research (Pallotti et al. 2010; Granger et al. 2015). Personal narrations and linguistic autobiographies are likely to offer insights into learners' experiences and linguistic awareness (Bonvino et al. 2023). Approximation tasks show strategies for dealing with linguistic gaps, while semi-structured interviews may provide controlled yet flexible data collection. Role-plays highlight interactive skills, while focus groups encourage more natural communication.

Keywords: longitudinal spoken corpus, L2 Italian, L1 Slavic, corpus compilation, data collection

References

Bonvino, E., Cortés Velásquez, D., De Meo, A., & Fiorenza, E. (2023). *Agire in L2: Processi e strumenti nella linguistica educativa*. Milano: Hoepli.

Council of Europe. (2020). Common European Framework of Reference for Languages: Learning, teaching, assessment – Companion volume. Council of Europe Publishing.

https://www.coe.int/en/web/common-european-framework-reference-languages

- Granger, S., Gilquin, G., & Meunier, F. (Eds.). (2015). *The Cambridge Handbook of Learner Corpus Research*. Cambridge: Cambridge University Press.
- Krueger, R. A., & Casey, M. A. (2015). Focus groups: A practical guide for applied research (5th ed.). Sage Publications.
- Migliorini, L., & Rania, N. (2001). I focus group: Uno strumento per la ricerca qualitativa. *Animazione sociale*, 2, 82–88.
- Ortega, L., & Iberri-Shea, G. (2005). Longitudinal research in second language acquisition: Recent trends and future directions. *Annual Review of Applied Linguistics*, 25, 26-45.
- Pallotti, G., Ferrari, S., Nuzzo, E., & Bettoni, C. (2010). Una procedura sistematica per osservare la variabilità dell'interlingua. *Studi Italiani di Linguistica Teorica e Applicata*, 39(2), 215–241.
- Robinson, P.; Ellis Nick C. (2008). *Handbook of Cognitive Linguistics and Second Language Acquistion*, New York, Routledge 270 Madison Ave.
- Robinson, P. & Gilabert, R. (2007). Task complexity, the Cognition Hypothesis and second language learning and performance. *International Review of Applied Linguistics in Language Teaching*, 45(3), 161-176. https://doi.org/10.1515/iral.2007.007

Applying learner corpus tools in the classroom: A data-driven approach to collocation development in EFL writing Bahareh Malmir (Boğaziçi University)

This work-in-progress paper reports on the application of learner corpus insights in a preparatory program at a Turkish university, where AntConc software was integrated into a five-session classroom workshop aimed at enhancing EFL students' lexical awareness and collocational range. Participants were B1–B2 level learners who engaged in guided corpusbased activities to explore academic collocations related to themes such as education and social media.

To evaluate the effect of the intervention, students completed a short essay task on the same topic both before and after the workshop. A comparison of these two writing samples revealed noticeable improvement in the use of lexical items introduced during the sessions. These gains, although limited to the short term, suggest that repeated exposure to such patterns could foster longer-term retention if reinforced over time.

The pedagogical value of AntConc was further supported by student reflections and post-session surveys, which highlighted the benefits of modeling, guided search strategies, and collaborative tasks. During lessons, informal discussions often compared AntConc to generative AI tools. Students noted that while AI can produce word lists or examples, it lacks the transparency of pattern discovery that AntConc offers. Moreover, effective use of AI requires advanced prompting literacy, and does not provide the keyword-in-context (KWIC) insights that are central to grammatical and collocational understanding. By contrast, AntConc promotes active noticing and discovery learning, particularly when supported with topic-specific corpora that offer richer, more relevant lexical input.

The study contributes to ongoing discussions on how data-driven learning can support vocabulary development, and proposes a practical model for integrating corpus tools into classroom-based instruction.

Keywords: Data-Driven Learning (DDL), learner corpus, AntConc, collocation awareness, vocabulary development, corpus-based pedagogy, lexical resources, EFL Writing

Behind the scenes of linguistic accuracy: Investigating the impact of automated writing evaluation on L2 academic theses *Katharina Maschke* (Chemnitz University of Technology)

As part of the PhD project "The Impact of Automated Writing Evaluation on Linguistic Accuracy in Academic Theses: A Comparative Analysis", this study aims to examine whether and how the use of Automated Writing Evaluation (AWE) impacts the linguistic accuracy of academic theses. Despite a growing emphasis on communication and meaning negotiation in foreign language learning (Wang et al., 2022, pp. 1–2), accuracy is still essential in conveying meaning, especially in higher-complexity language contexts such as academic writing, where concepts of elevated complexity must be communicated precisely and effectively (Ferris, 2006, p. 81; Hyland & Hyland, 2006, p. 2; Liao, 2016, p. 308).

One source of feedback on accuracy in L2 academic settings is automated writing evaluation (AWE), which provides "instantaneous feedback on spelling and grammar" often integrated into word processors or as separate software (McCarthy et al., 2022, p. 1). Existing research into the use of AWE in academia has primarily focused on applications inside the classroom and their impact on learning L2 writing, leaving the effect AWE has on the quality of high- stakes academic writing outside of the classroom underexplored. To address this research gap, two surveys are distributed. The first survey gathers student perspectives on the perception and popularity of AWE tools. The second survey collects L2 English theses accompanied by relevant metadata for corpus compilation and analysis. The theses will be coded via a hybrid automatic/manual annotation process, and lexicogrammatical error patterns will be analyzed using multiple regression analysis. By comparing texts written with and without the help of AWE, the study aims to identify statistically significant differences in error pattern frequency. Preliminary results from both data collection instruments will be presented at the conference to offer insights into both the perceived and actual impact of AWE on the writing accuracy of advanced L2 English academic authors.

Keywords: academic writing, L2 writing, automated writing evaluation, linguistic accuracy

- Ferris, D. (2006). Does error feedback help student writers? New evidence on the short and long-term effects of written error correction. In K. Hyland & F. Hyland (Eds.), Feedback in second language writing: Contexts and issues (1st ed., pp. 81–104). Cambridge University Press. https://doi.org/10.1017/CBO9781139524742
- Hyland, K., & Hyland, F. (2006). Feedback in second language writing: Contexts and issues (1st ed.). Cambridge University Press. https://doi.org/10.1017/CBO9781139524742
- Liao, H.-C. (2016). Using automated writing evaluation to reduce grammar errors in writing. *ELT Journal*, 70(3), 308–319. https://doi.org/10.1093/elt/ccv058
- McCarthy, K. S., Roscoe, R. D., Allen, L. K., Likens, A. D., & McNamara, D. S. (2022).

 Automated writing evaluation: Does spelling and grammar feedback support high-quality writing and revision? *Assessing Writing*, 52, 100608. https://doi.org/10.1016/j.asw.2022.100608

Wang, M., Wang, H., & Shi, Y. (2022). The role of English as a foreign language learners' grit and foreign language anxiety in their willingness to communicate: Theoretical perspectives. *Frontiers in Psychology*, 13, 1002562. https://doi.org/10.3389/fpsyg.2022.1002562

Investigating second language pragmatic competence across proficiency level and first language, using a multimodal corpus Gerard O'Hanlon (Mary Immaculate College)

This research focuses on how English language learners use videoconferencing for spoken interaction, charting what features of requests are evident across CEFR levels.

This presentation presents initial findings from a work-in-progress study. It employs a multimodal corpus pragmatics approach (Adolphs and Chen 2021) to investigate how request sequences (Felix-Brasdefer 2021) align with politeness theory (Culpeper et al 2017) and nonverbal communication (e.g. Rohrer et al 2022) in an audiovisual dataset.

The corpus consists of the spoken interactions of 12 pairs of English language learners. Participants were recruited at three CEFR levels (B1, B2 and C2) from four L1 backgrounds (Brazilian Portuguese, Japanese, Spanish and German). All participants interacted with partners from the same L1 at the same CEFR level. The spoken data was elicited through open roleplays (Felix-Brasdefer 2018) and delivered and recorded via *Zoom*, resulting in a video dataset of approximately 5 hours, 45,000 tokens and 108 request roleplays. ELAN was used for transcription and annotation.

Corpus pragmatics involves manually annotating the roleplays' core requests. Request type is also annotated (direct/indirect) to gauge variation. Given that requests are impositive, dispreferred and often occur deeper into conversations due to their face-threatening nature (Al-Gahtani and Roever 2012), their onset times are also recorded. Finally, multimodality (Jewitt et al 2016) permits exploration of speech acts beyond orthographic transcription, specifically embodied actions (e.g. gesture or gaze).

Initial findings show uniformity occurring regarding request type (indirect) and form (modal verbs) across all CEFR levels. Other features, such as the use of pragmatic modifiers (downtoning, hedging) increase in use and sophistication at higher proficiency levels.

There is wide variability pertaining to request onset times across levels (i.e. the speaker's perception of request imposition on the interlocutor). This highlights the dynamics of request co-construction across pairings and the contextual of speech act production. Finally, multimodal embodiment is flexible, variable and unique to each request, requiring further qualitative insight.

The focus on pragmatics and multimodal communication makes for a logical intersection in furthering insights into the competencies of these interactions in digital discourses.

Keywords: corpus pragmatics, multimodality, online audiovisual interaction, CEFR, English language learner spoken language

References

Adolphs, S. and Chen, Y. (2021) 'Corpus Pragmatics', in M. Haugh, D.Z. Kádár, and M.

- Terkourafi (eds.) The Cambridge Handbook of Sociopragmatics. Cambridge: Cambridge University Press (Cambridge Handbooks in Language and Linguistics), pp. 639–662.
- Al-Gahtani, S. and Roever, C., 2012. Proficiency and sequential organization of L2 requests. Applied Linguistics, 33(1), pp.42-65.
- Culpeper, J., Terkourafi, M. (2017). Pragmatic Approaches (Im)politeness. In: Culpeper, J., Haugh, M., Kádár, D. (eds) The Palgrave Handbook of Linguistic (Im)politeness. Palgrave Macmillan, London.
- Félix-Brasdefer, J.C., 2018. Role plays. In A.H. Jucker, K.P. Schneider and W. Bublitz (eds), Methods in Pragmatics . Berlin: Mouton de Gruyter, pp. 305–331.
- Félix-Brasdefer, J.C., 2021. Pragmatic competence and speech-act research in second language pragmatics. New directions in second language pragmatics, pp.11-26.
- Jewitt, C., Bezemer, J. and O'Halloran, K., 2016. Introducing multimodality. Routledge.
- Rohrer, P. L., Vilà-Giménez, I., Florit-Pons, J., Gurrado, G., Esteve-Gibert, N., Ren-Mitchell, A., Shattuck-Hufnagel, & Prieto, P. (2023). The MultiModal MultiDimensional (M3D) labeling system.

Do subpart frequency and phrase frequency affect articulatory durations of multi-word expressions produced by L2 English learners? Yi Qi & Anna Siyanova-Chanturia (Victoria University of Wellington)

Whether or not multi-word expressions (MWE) are stored and retrieved holistically in the mental lexicon has been a central question discussed in the literature. Following Wray's (2002) seminal proposal on the holistic storage of MWEs, studies have shown that MWE processing is not entirely holistic. For instance, Arnon and Cohen Priva (2014) showed that during the production of trigrams, phrase frequency effects increased, while word frequency effects decreased, but did not disappear. However, studies in this area have predominantly focused on L1 speakers and/or have relied on elicited data. L2 phrase frequency effects during language production have so far received far less attention, with no study to date looking at naturalistically produced (rather than laboratory elicited) data.

The present study uses the International Corpus Network of Asian Learners of English (ICNALE) to investigate whether and how phrase frequency and subpart frequency influence the articulatory durations of verb-preposition-noun (VPN, e.g., ask for help) sequences produced by L2 English learners in spontaneous speech. The learners' proficiency levels range from beginning to upper-intermediate. The corpus consists of interviews involving picture descriptions and role-plays conducted by a language instructor and the learners. Subpart frequency here refers to the frequency of a lexical combination within the sequence. For example, the subpart frequencies of the MWE "ask for help" include unigram frequency (i.e., ask, for, help), first bigram (i.e., ask for) and second bigram frequencies (i.e., for help). The spoken data with word-level timestamps were transcribed and were subsequently assigned part-of-speech tags. All VPN sequences were extracted, and the articulatory durations of the three-word sequence were calculated. Mixed-effect models were used to analyse the data. The present study provides important evidence as to L2 learners' sensitivity to subpart and phrase frequency in language production and further contributes to the holistic storage debate.

Keywords: multi-word expressions/MWEs, spoken learner corpus, frequency effect, language processing, second language/L2

References

Arnon, I., & Cohen Priva, U. (2014). Time and again: The changing effect of word and multiword frequency on phonetic duration for highly frequent sequences. *The Mental Lexicon*, 9(3), 377–400. https://doi.org/10.1075/ml.9.3.01arn

Wray, A. (2002). Formulaic language and the lexicon. Cambridge University Press.

Introducing an index for analysis of grammatical diversity in writing Christian Holmberg Sjöling & Taehyeong Kim (Luleå University of Technology & Northern Arizona University)

Much research using an automated and quantitative approach to investigate grammatical complexity features in writing do so with a frequency-based approach. In this paper, we build on the frequency-based approach to tap into a complementary dimension of grammatical complexity that we term grammatical diversity. That is, while frequencies of individual grammatical features offer valuable insights, they may mask how diverse grammatical features are used by writers with different proficiency levels and across different registers.

To tap into such a construct of grammatical diversity, we have written a Python script that builds on tagging initially carried out with the Lexicogrammatical Tagger (Biber et al., 2025). The script counts the number of unique grammatical complexity features for a moving five-sentence window of a text (e.g., sentence 1–5, 2–6, 3–7 and so on) (see also Zenker & Kyle, 2021). This approach captures local variability and patterns of grammatical complexity features used in a text and accounts for the finite number of grammatical features that writers eventually must repeat. The measure is applied in a pilot study where a sample of 7,702 L2 texts from the EF-Cambridge Open Language Database (EFCAMDAT) (Geertzen et al., 2014; Huang et al., 2014) is analysed. Firstly, diagnostics were run to establish window size, then, a minimally sufficient approach from Staples et al. (2023) was used to determine if there was a difference of grammatical diversity between different proficiency levels in the corpus. The findings show a steady developmental increase (i.e., non-overlapping Cis and 10% increase) across all proficiency levels for both phrasal and clausal diversity. The computation of the measurement and the findings are discussed further.

Keywords: grammatical complexity, grammatical diversity, learner writing

- Biber, D., Dirkx, J., Dussan, H., Egbert, J., Kyle, K., Paddock, A., Reppen, R., Sung, H., & Walker., R. (2025). *Lexicogrammatical Tagger*.
- Geertzen, J., Alexopoulou, T., & Korhonen, A. (2014). Automatic linguistic annotation of large scale L2 databases: The EF-Cambridge Open Language Database (EFCamDat). In R.T.
- Millar, K.I. Martin, C.M. Eddington, A. Henery, N.M. Miguel, & A. Tseng (Eds.), Selected proceedings of the 2012 Second Language Research Forum (pp. 240–254). Somerville, MA: Cascadilla Proceedings Project.
- Huang, Y., Geertzen, J., Baker, R., Korhonen, A., & Alexopoulou, T. (2017). *The EF Cambridge Open Language Database (EFCAMDAT): Information for users* (pp. 1–18). Retrieved from https://ef-lab.mmll.cam.ac.uk/EFCAMDAT.html
- Staples, S., Gray, B., Biber, D., & Egbert, J. (2023). Writing Trajectories of Grammatical Complexity at the University: Comparing L1 and L2 English Writers in BAWE. *Applied Linguistics*, 44(1), 46–71. https://doi.org/10.1093/applin/amac047
- Zenker, F., & Kyle, K. (2021). Investigating minimum text lengths for lexical diversity indices. *Assessing Writing*, *47*, 100505.

The influence of grammatical complexity features on grade in Swedish high-stakes EFL exams Christian Holmberg Sjöling (Luleå University of Technology)

This paper examines the effect of eight grammatical complexity features (e.g., verb + that complement clauses, finite adverbial clauses, attributive adjectives, pre-modifying nouns) on grading during high-stakes EFL exams. This is done by analysing a learner corpus consisting of 142 graded example texts and 175 teacher graded texts written during the National Tests of English in Sweden (these tests are similar to the *International English Language Testing System*, IELTS).

The graded example texts are graded by expert raters and are a part of the assessment instructions that should guide teachers in the assessment of their own students' texts. Teachers should also take course specific grading criteria into consideration when carrying out the assessment, and these criteria mention that some form of grammatical structures should influence grade, but it is not clear which features nor how.

The following research questions are answered: 1. To what extent does the frequency of grammatical complexity features characteristic of spoken and informational registers influence grade? 2. To what extent is the influence of grammatical complexity features in teachers' grading in line with the expert raters? To answer these questions, the material was initially tagged with the *Lexicogrammatical Tagger* (Biber et al., 2025). Then, statistical analysis was carried out using a minimally sufficient approach adapted from Staples et al. (2023). The results show that only one feature characteristic of informational writing, *Ofgenitive* phrases (e.g., *Do you know the name of this flower?*), has a clear positive influence (i.e., non-overlapping CIs) on grade in assessment carried out by expert raters and teachers. However, one feature characteristic of spoken registers, verb + *that* complement clauses, shows a negative correlation with grade in texts assessed by expert raters, but this is not reflected in teachers' grading. These findings are further discussed in relation to language testing, test construction, and classroom teaching.

Keywords: grammatical complexity, learner writing, language testing, high-stakes exams

References

Biber, D., Dirkx, J., Dussan, H., Egbert, J., Kyle, K., Paddock, A., Reppen, R., Sung, H., & Walker.,R. (2025). *Lexicogrammatical Tagger*. https://lcr-ads-lab.github.io/LxGrTagger-Documentation/ [visited 2025-03-30].

Staples, S., Gray, B., Biber, D., & Egbert, J. (2023). Writing Trajectories of Grammatical Complexity at the University: Comparing L1 and L2 English Writers in BAWE. *Applied Linguistics*, 44(1), 46–71. https://doi.org/10.1093/applin/amac047

Communicative strategies in oral and written Romanian: A case study *Isabella Şinca* (University of Bucharest)

This doctoral research investigates the communicative strategies used by learners of Romanian as a foreign language when facing lexical gaps in oral and written discourse. Building on prior studies (Tarone 1981; Dörnyei & Scott 1997; Putri 2013) and continuing previous MA-level work on written discourse (Şinca 2024a, 2024b), the study has three main goals: (1) to identify the types of communicative strategies used in both oral and written discourse; (2) to analyze the differences between oral and written strategies; (3) to explore how learners' L1s influence strategy choice.

Using Dörnyei & Scott's (1997) taxonomy, the study analyzes data from students with diverse linguistic backgrounds (e.g., Arabic, Turkmen, Bulgarian) learning Romanian at the University of Bucharest. The corpus comprises classroom-based written texts and oral transcriptions from the LECOR corpus (Barbu et al. 2023).

Preliminary findings (Şinca 2024a, 2024b) from written data show frequent use of approximation (e.g., *vremea României*, instead of *clima României* 'Romania's climate') and code-switching (e.g., *puțin pepper*, instead of *puțin piper* 'a little pepper'). Strategy use appears to differ based on learners' L1s, suggesting that linguistic and possibly cultural backgrounds shape preferences. Differences also emerge between oral and written discourse, with some strategies being specific to oral or written discourse.

By examining both spoken and written communication, this research contributes a broader perspective to the field, extending beyond earlier studies focused solely on oral production (Vasiu 2020). The findings can inform teaching practices and help develop pedagogical materials tailored to learners' strategic needs in both forms of expression.

Keywords: Romanian learners, communicative strategies, spoken vs. written discourse

References

Barbu, Ana Maria, Elena Irimia, Carmen Mîrzea Vasile, Vasile Păiș, 2023, "Designing the LECOR Learner Corpus for Romanian", in Galia Angelova, Maria Kunilovskaya and Ruslan Mitkov Varna (ed.), 2023, Deep Learning for Natural Language Processing Methods and Applications (Proceedings of International Conference Recent Advances in Natural Language Processing, RANLP 2023, Varna, 4–6 September, 2023), p. 143-152 (online ISBN: 978-954-452-092-2, e-book site: www.acl-bg.org, series online: ISSN 2603-2813.2023, INCOMA Ltd. Shoumen, Bulgaria).

Dörnyei, Zoltán, Mary Lee Scott, 1997, "Communication strategies in a second language: Definitions and taxonomies", *Language learning*, 47.1: 173-210.

Kennedy, Sara, 2022, "Second Language Speaking Strategies", in Derwing, Tracey M., Murray J. Munro, Ron I. Thomson (ed.), *The Routledge Handbook of Second Language Acquisition and Speak*, New York, Routledge Taylor & Francis, 261-272.

- Mîrzea Vasile, Carmen, Ana-Maria Barbu, Valentina Cojocaru, Mihaela Cristescu, Elena Irimia, Simona Neagu, Vasile Păiș, Isabella Șinca, Monica Vasileanu, 2025, "Aspects regarding the use of the Learner Corpus of Romanian (Lecor)", Synergy the Journal of the Department of Modern Languages and Business Communication, Faculty of International Business & Economics, the Bucharest University of Economic Studies, in the process of being published.
- Putri, Lidya Ayuni, 2013, "Communication strategies in English as a second language (ESL) context", *Advances in Language and Literary Studies* 4.1: 129-133.
- Şinca, Isabella, 2024a, "Strategii de comunicare la nivelul A1 în producții scrise de nativi arabi, turkmeni, albanezi și francezi. Studiu de caz", in E. Platon, C. Bocoș, D. Roman, L. Vasiu (ed.), *Discurs polifonic în româna ca limbă străină (RLS). Actele Conferinței Internaționale*, Cluj-Napoca, 4th edition/ 2023, Nr. 4/2024, 189-203.
- Şinca, Isabella, 2024b, *Strategii de comunicare în româna ca limbă nenativă în texte scrise* (nivelul A1). Câteva studii de caz, dissertation, Faculty of Letters, University of Bucharest, unpublished.
- Tarone, Elaine, 1981, "Some thoughts on the notion of communication strategy", *TESOL Quarterly* 15.3: 285-295.
- Vasiu, Lavinia-Iunia, 2020, Achiziția limbii române ca L2. Interlimba la nivelul A1, Editura Universitară Clujeană, Cluj-Napoca.

Analysis of metaphor production in Finnish L1 learner English *Renata Turunen (University of Inland Norway)*

This research aims to provide an empirical quantitative portrait of metaphor production in Finnish L1 English, contributing to the previous work on the ubiquity of metaphor in language and describing metaphor use of the previously unexplored learner group. The investigation is performed by comparing learner produced written texts containing deliberately elicited metaphor (40 collected essays, 14,336 lexical units) to learner produced written texts where metaphor was not elicited (27 essays from ICLE corpus, 14,526 lexical units). Metaphors were identified with metaphor identification procedure MIPVU (Steen et al., 2010), which involves determining the metaphorical status of every lexical unit. Metaphor density was calculated by counting the ratio of metaphor-related words to the total number of lexical units (Littlemore et al., 2014). The results showed no statistically significant difference in metaphor density between the two datasets (p=0.2599). The overall metaphorical density for the whole sample of Finnish learner essays is 16.80%.

Metaphor density fluctuates throughout a text, showing varied frequencies of metaphors. Segments of discourse with concentrated metaphor usage are referred to as 'metaphor clusters', which are especially valuable for deeper analysis because these clusters often play significant roles in communication (Nacey, 2020, p. 295). They are used to explain and summarize ideas; to manage discourse, for example in discussion openings and endings (Nacey, 2022). The current research outlines a qualitative investigation of metaphor clusters in the learner essays.

Keywords: Finnish learner English, metaphor production, metaphor density, metaphor clusters

- Littlemore, J., Krennmayr, T., Turner, J., & Turner, S. M. (2014). An Investigation into Metaphor Use at Different Levels of Second Language Writing. Applied linguistics, 35(2), 117-144. https://doi.org/10.1093/applin/amt004
- Nacey, S. (2020). Figurative production in a computer mediated discussion forum. In J. Barnden & A. D. Gargett (Eds.), Producing figurative expression: theoretical, experimental and practical perspectives (Vol. Volume 10, pp. 363-388). John Benjamins Publishing Company.
- Nacey, S. (2022). Development of metaphorical production in learner language: A longitudinal perspective. Nordic Journal of Language Teaching and Learning, 10(2), 272-297. https://doi.org/10.46364/njltl.v10i2.975
- Steen, G. J., Dorst, A. G., Herrmann, J. B., Kaal, A., Krennmayr, T., & Pasma, T. (2010). MIPVU. In. The Netherlands: John Benjamins Publishing Company.

Validating automated measures of phraseological competence in L2 speaking: Evidence from human judgments *Tingting Wang (Nanjing University)*

Phraseological competence is crucial for language acquisition, processing, and fluency (Ellis *et al.*, 2008; Paquot *et al.*, 2020) but remains challenging for L2 learners (Laufer & Waldman, 2011; Paquot & Granger, 2012). While existing studies have developed automated indices to operationalize this construct and demonstrated their predictive validity by linking them to L2 proficiency, few have directly addressed whether the indices adequately capture the construct itself (Paquot & Naets, 2025). This raises concerns about their construct validity – the extent to which a measure accurately represents the theoretical construct it is intended to assess. Establishing construct validity requires evidence of alignment between automated measures and alternative methods, with human ratings serving as a critical benchmark in applied linguistics (Crossley *et al.*, 2013). However, few studies have evaluated the alignment between human judgments and automated measures of phraseological competence, particularly in L2 oral production, highlighting the need for further validation research.

Grounded in an argument-based framework (Kane, 2006), this study explores the following research questions to provide convergent evidence for automated measures of phraseological competence in L2 speaking: **RQ1.** To what extent do automated measures of phraseological competence align with human ratings? **RQ2.** What features of phraseological competence do human raters focus on, and how well do automated indices capture them?

A mixed-methods approach was adopted, consisting of two phases: (1) a corpus linguistic analysis using automated measures and (2) an experiment with human raters. In the first phase, oral performances from 98 test-takers of the TEM 8-Oral (Test for English Majors-Band 8) were analyzed using automated indices. Phraseological competence was operationalized along three dimensions–accuracy, diversity, and sophistication—and measured with respect to two key phenomena in phraseology: co-occurrence and recurrence (Granger & Paquot, 2008). Co-occurrence was examined through six grammatical relations (adjectival modifier (amod), direct object (dobj), adverbial modifier (advmod), adjectival complement (acomp), nominal subject (nsubj), and prepositional modifier (prep)). Recurrence was explored through the analysis of three-word lexical bundles. In the second phase, 30 human raters were recruited to evaluate the same performances using a comparative judgment method. They provided both ratings of phraseological competence and qualitative comments explaining their decisions. These human judgments served as external benchmarks for automated measures and offered insights into the features that influenced rater evaluations.

Based on a pilot study in which two raters evaluated 30 texts and provided comments on their decisions, we anticipate the following findings: 1) Automated measures of phraseological competence will demonstrate varying degrees of alignment with human ratings, with the sophistication measures for lexical bundles expected to show the strongest

correlation and explanatory power; 2) Automated indices will effectively capture certain features prioritized by human raters, particularly phraseological diversity and sophistication. However, they may be less effective in adequately capturing some nuanced qualitative aspects that are equally emphasized in human evaluations of phraseological competence, such as idiomaticity and contextual appropriateness. The findings of this study will provide convergent validity evidence for the automated measures of phraseological competence and highlight areas where computational measures require refinement, particularly in capturing qualitative features emphasized by human raters. By bridging automated analysis and human judgment, the study aims to inform the development of more robust assessment tools and offer pedagogical insights for fostering phraseological competence in L2 learners.

Keywords: phraseological competence, automated measures validity, comparative judgment, L2 speaking

- Crossley, S., Salsbury, T., & McNamara, D. S. (2013). Validating lexical measures using human scores of lexical proficiency. In S. Jarvis & M. Daller (Eds.), *Studies in bilingualism* (Vol. 47, pp. 105-134). John Benjamins Publishing Company. https://doi.org/10.1075/sibil.47.06ch4
- Ellis, N. C., Simpson-Vlach, R., & Maynard, C. (2008). Formulaic language in native and second language speakers: Psycholinguistics, corpus linguistics, and TESOL. *TESOL Quarterly*, 42(3), 375–396. https://doi.org/10.1002/j.1545-7249.2008.tb00137.x
- Granger, S., & Paquot, M. (2008). Disentangling the phraseological web. In S. Granger & F. Meunier (Eds.), *Phraseology: An Interdisciplinary Perspective* (pp. 27–49). Amsterdam & Philadelphia: John Benjamins. https://doi.org/10.1075/z.139.07gra
- Kane, M. T. (2006). Validation. In R. Brennen (Ed.), *Educational measurement* (4th ed., pp. 17–64). Praeger and Greenwood Publishing.
- Laufer, B., & Waldman, T. (2011). Verb-noun collocations in second language writing: A corpus analysis of learners' English. *Language Learning*, 61(2), 647–672. https://doi.org/10.1111/j.1467-9922.2010.00621.x
- Paquot, M., & Granger, S. (2012). Formulaic language in learner corpora. *Annual Review of Applied Linguistics*, 32, 130–149. https://doi.org/10.1017/S0267190512000098
- Paquot, M., & Naets, H. (2025). Phraseological sophistication as a multidimensional construct: Exploring the relationship between association, register specificity and frequency of word combinations. In T. Larsson & D. Biber (Eds.), *Cumulative knowledge building and replication in Learner Corpus Research*. International Journal of Learner Corpus Research, 11(1). https://doi.org/10.1075/ijlcr.23033.paq
- Paquot, M., Gries, S. Th., & Yoder, M. (2020). Measuring lexicogrammar. In P. Winke & T. Brunfaut (Eds.), *The Routledge handbook of second language acquisition and language testing* (pp. 223–232). Routledge. https://doi.org/10.4324/9781351034784-25

POSTER PRESENTATIONS

LC-meta: A core metadata schema for L2 data *Jennifer-Carmen Frey*¹, *Alexander König*², *Hubert Naets*³, *Egon W. Stemle*¹ *and Magali Paquot*³ (¹*Institute for Applied Linguistics, Eurac Research,* ²*CLARIN ERIC,* ³*UCLouvain*)

The Core Metadata Schema for L2 data consists in a comprehensive set of variables that encapsulate crucial information about L2 data. It is organized into several sections that describe specific aspects of a learner corpus. These include administrative details (e.g. authors or license), corpus design, text-related variables, learner-related variables, in-built annotation (e.g. details about manual or automatic annotation), information about annotators or transcribers (e.g. native language or language repertoire) and task-related details (e.g. instructions, time constraints) (Paquot et al., 2024).

While the schema started out as the result of extensive collaboration between learner corpus compilers at the Centre for English Corpus Linguistics (UCLouvain, Belgium) and Eurac Research (Bolzano, Italy), and a research data infrastructure expert and member of CLARIN's metadata taskforce, further development and user involvement is driven by the LC-meta working group organized under the aegis of the learner corpus association (LCA). In this poster session, we invite young researchers to discuss their questions and feedback on the current metadata schema v2.0 with members of the working group in a personal question and answer session.

Keywords: learner corpus, metadata, schema, study comparability, findability

- Frey, J.-C., König, A., Stemle, E. & M. Paquot (2023). A core metadata schema for L2 data.

 Paper presented at the 32nd Conference of the European Second Language
 Association (EUROSLA), 30 August 2 September 2023, University of Birmingham,
 UK.
- König, A., Frey J.-C., Stemle, E., Glaznieks, A. & M. Paquot (2022). Towards standardizing LCR metadata. Paper presented at Learner Corpus Research 6, 22-24 September 2022, University of Padua, Italy.
- Paquot, M., König, A., Stemle, E. & J.-C. Frey (2023). Core Metadata Schema for Learner Corpora, https://doi.org/10.14428/DVN/4CDX3P

Separating agreement and assignment of gender marking in L2 Spanish writing Gabriela Sanchez (University of Texas at Arlington)

Existing research on patterns of accuracy amongst L2 learners of grammatical gender suggests that the masculine gender appears to be acquired more easily, shown by higher accuracy rates assigning gender to masculine nouns (Sabourin et al., 2004; Alarcón, 2011; Gudmestad et al., 2012), higher percentage of accuracy in grammaticality judgment tasks for masculine, and overall higher assignment of nouns to masculine, even incorrectly (Kirova & Camacho, 2021). These and other studies have concluded that learners default to assigning unfamiliar nouns as masculine. The current study seeks to shed light on these results by determining whether noun phrase type (assignment only vs. agreement and assignment) and noun canonicity dictate accuracy in grammatical gender use by L2 learners of Spanish.

In this investigation, written texts of L1 English Spanish learners are examined from the *Corpus Escrito del Español L2* (CEDEL2; Lozano 2021), a learner corpus, to track accuracy of gender assignment and agreement productions. By separating the data into proficiency levels for comparison, this study aims to increase understanding of error production and patterns as well as the sequence of acquisition regarding grammatical gender, specifically, noun agreement and assignment. Nouns will be separated by canonicity: canonical, ending in "-o" like *banco* and noncanonical, ending in "-e" or a consonant like *puente*, *lapiz*, *sol*. The focus of this study is masculine and feminine determiners in singular and plural forms with special attention paid to gender agreement and assignment. This includes: *el/la*, *los/las*, *varios/as*, *unos/as*, *muchos/as*, *todos/as*, thus focusing on twelve noun phrase types. Crucially, the current study's results might clarify whether certain types of nouns are acquired first, in this case focusing on canonicity, whether gender assignment and agreement are learned concomitantly, and whether certain types of noun phrases are more likely to be accurate than others.

A preliminary analysis using a small subset of the data focused exclusively on definite determiners, finding that learners exhibited higher accuracy rates using masculine agreement and assignment (90%-98%) compared to feminine agreement and assignment (78%-100%), as previous studies have shown. As expected, the advanced proficiency learners had little to no errors, while lower proficiency learners had the highest error rate for all determiner forms (22%, 18%). Interestingly, the error rate for the masculine singular determiner 'el' shows a gradual decline as the learner proficiency increases, a pattern not paralleled by the other determiners. This result could suggest acquisition of the masculine gender before the feminine, easier acquisition of certain determiner-noun pairs than others, or it could have something to do with the proportion of canonical nouns in each category. The full analysis examines each incorrect grammatical gender production by conducting separate analyses on types of nouns (canonical, noncanonical), and noun phrase types, i.e., assignment only (determiner + noun), or assignment and agreement (determiner + noun + adjective, etc.), hypothesizing significantly higher accuracy rates for determiner + noun phrase types (assignment without agreement) using masculine canonical nouns.

Keywords: Spanish, grammatical gender, L2 writing

Complexity development of Intermediate German learners of English: A longitudinal corpus analysis *Philine Metzger (Philipps University Marburg)*

Complexity is a topic of wide range, overall separated into various levels, namely lexis (Linnarud 1986), morphology (Brezina & Pallotti 2019), syntax (Lee 2004) and phraseology (Paquot 2019). It is generally assumed that syntactical complexity of a language correlates positive with overall language competence and development. (Bulté 2008 & Paquot, Naets, Gries 2021) For example, length of T-Units, clauses per T-Unit or dependent clauses per T-Unit. (Lee 2004: 108) Previous studies already treated the topic, whereas their corpora were much smaller. (Kyle, Crossley, Verspoor 2021).

The present project therefore poses the following research question: How does the quantitative and qualitative linguistic complexity in written English of intermediate learners of a German high school develop between 9th and 12th grade? To answer these research question, the "Marburg Corpus of Intermediate Learner English" (MILE; Kreyer 2015) gives the opportunity to analyze data of 90 students from 9th to 12th grade. Within more than 500,000 words in total, the MILE corpus for enables the conduction of a truly longitudinal analysis of a large number of intermediate learners of English over four years. Also, metadata such as gender and age was made available. This corpus will be analyzed using the "Tool for Automatic Analysis of Syntactic Sophistication and Complexity" (TAASSC; Kyle 2016) which includes Lu's 14 parameters of complexity.

The method that will be used replicates previous studies connected to syntactic complexity (Crossley, McNamera 2012) and is called the "bottom-up" principle, which means the parameters will first be extracted (TAASSC) and will then be analyzed via mixed effect regression modeling in R. The project analyses the complexity development considering Lu's 14 parameters of measuring complexity, namely the main categories mean length of production unit, sentence complexity ratio, subordinations, coordinations and particular structures including subcategories. (Lu 2010: 479). After that, quantitative (TAASSC with Lu's 14 parameters) and qualitative (teacher ratings such as grammar, lexis, variation and sentence structure) will be compared to one another to gain knowledge in both fields of syntactic development. Results are expected to show the impact of given metadata on complexity development and will be discussed in the light of their language-pedagogical. Finally, the project results in an individual and multifactorial learner curve.

Keywords: syntactic complexity, corpus based study, longitudinal study, Intermediate learners

References

Brezina, Vaclaw; Palotti, Gabriele. 2019. Morphological Complexity in Written L2 Texts. Second Language Research. 35 (1): 99-119

Bulté, Bram; Housen, Alex; Pierrad, Michel; Van Daele, Siska. 2008. Investigating Lexical Proficiency Development over Time – The Case of Dutch-Speaking Learners of French

- in Brussels. Journal of French Language Studies 18 (3): 277-298.
- Crossley, Scott; McNamera, Danielle. 2012. Predicting Second Language Writing
 Proficiency: The Roles of Cohesion and Linguistic Sophistication. Journal of
 Research in Reading 35 (2): 115-135.
- Kreyer, Rolf. 2015. The Marburg Corpus of Intermediate Learner English (MILE). In Learner Corpora in Language Testing and Assessment, Marcus Callies & Sandra Götz, eds. Amsterdam: John Benjamins. 13-34.
- Kyle, K. (2016). Measuring syntactic development in L2 writing: Fine grained indices of syntactic complexity and usage-based indices of syntactic sophistication (Doctoral Dissertation)
- Lee, Jiyoujng (2004). Syntactic Complexity, Clausal Complexity, and Phrasal Complexity in L2 Writing: The Effects of Task Complexity and Task Closure. The Journal of Asia TEFL. South Korea: 108-124
- Linnarud, Moira. 1986. Lexis in Composition: A Performance Analysis of Swedish Learners Written English. Malmö: C.W.K. Gleerup.
- Paguot, Magali. 2019. The Phraseological Dimension in Interlanguage Complexity Research. Second Language Research 35 (1): 121-145.
- Paquot, Magali; Naets, Hubert; Gries, Stefan. 2021. Using Syntactic Co-occurrences to Trace Phraseological Complexity Development in Learner Writing: Verb + Objekt Structures in LONGDALE. In: LeBruyn, Bert Simonne Walter; Paquot, Magali (hrsg): Learner Corpus Research Meets Second Language Acquisition. Cambridge: Cambridge University Press.

Comparative analysis of the modal verb 'could' among Indonesian and Japanese EFL learners Nida Ghaziyah (Kanazawa University)

Modal verbs, such as could, play essential roles in expressing modality across academic and formal discourse. However, their acquisition by English as a Foreign Language (EFL) learners remains complex due to the semantic flexibility of modal expressions and the influence of learners' first languages (L1s). This study investigates the usage patterns of could among Indonesian and Japanese EFL learners, drawing on written data from the International Corpus Network of Asian Learners of English (ICNALE). Notably, the differing equivalents of could in Indonesian (bisa, dapat) and Japanese (dekiru, kamoshirenai) may contribute to divergent learner behaviors, shaped by their typologically distinct L1s (Shibatani, 1990; Sneddon, 1996; Maynard, 1993; Oktavianti, 2019). The current stage of analysis focuses on frequency as a foundational step for broader exploration. Frequency analysis was selected for two main reasons. First, it provides an initial quantitative overview that supports further comparison across L1 groups and proficiency levels. By normalizing frequency per 10,000 tokens and applying proportion tests, the study examines CEFR-level differences. Results indicate that Japanese learners tend to use could more frequently at the B1_1 level, though this statistical significance may not directly reflect functional contrasts. Second, frequency findings serve as a basis for guiding qualitative investigation. A functional classification of could was developed based on pre-established categories drawn from previous research, including Ability, Hypothetical Reasoning, and Politeness (Coates, 1983; Palmer, 2001; and Biber et al., 1999), who outline dynamic, epistemic, and speech-act related functions. This study seeks to connect quantitative patterns with functional usage to explore how learners' L1 backgrounds shape their use of could in written English. The analysis provides a foundation for future investigations into deeper patterns such as collocational behavior and syntactic dependencies.

Keywords: EFL learners, could, modal verbs, ICNALE