# Syntactic and lexical complexity in written L2 English

A comparative study of Finnish and Hungarian university-level students

Iina Hesso

Master's Thesis

Degree Programme in Language Learning and Teaching, Department of English

School of Languages and Translation Studies

Faculty of Humanities

University of Turku

April 2024

Master's Thesis

**Degree Programme in Language Learning and Teaching, Department of English:**
**Iina Hesso**
**Syntactic and lexical complexity in written L2 English – A comparative study of Finnish and Hungarian university-level students**
**Number of pages**: 59 pages, 6 appendices

This thesis aims to compare *syntactic* and *lexical complexity* in essays written by Finnish and Hungarian university-level English language students and examine how syntactic structures and lexical elements manifest in observed high-complexity and low-complexity essays. Thirty essays by Finnish students and thirty essays by Hungarian students from the ACALEX corpus were utilised as the primary data for the analysis. Syntactic complexity and lexical complexity were measured by using web-based interfaces of automatic analyser tools *The L2 Syntactic Complexity Analyzer* and *The Lexical Complexity Analyzer*. The results were further analysed in Microsoft Excel and IBM SPSS Statistics 27. Finally, a qualitative extreme case analysis was conducted based on the statistical analysis' results.

In the quantitative analysis, it was discovered that Finnish students, on average, display higher syntactic and lexical complexity in their essays than Hungarian students. Finnish students, on average, tend to produce longer syntactic structures, use more complex grammar, and write more lexically dense language with more variation and sophistication in word choices than Hungarian students. The qualitative analysis revealed that high complexity may manifest in the choice of long syntactic structures and varied sophisticated lexical elements. Low complexity is characterised by short syntactic structures, repetition of words, and a lack of sophistication in vocabulary. The results of the thesis, thus, suggest that phenomena such as written learner language, L1 background, and linguistic complexity in L2 can affect each other in various ways. Therefore, the findings indicate that these phenomena should be considered in foreign language instruction for effective language learning.

**Key words**: syntactic complexity, lexical complexity, linguistic complexity, learner language, L2 proficiency

# Table of contents

# Tables and Figures

# Abbreviations

CAF = complexity, accuracy, and fluency

EFL = English as a foreign language

ESL = English as a second language

L1 = first language

L2 = second language

L2SCA = The L2 Syntactic Complexity Analyzer

LC = lexical complexity

LCA = The Lexical Complexity Analyzer

NCC = National Core Curriculum

SC = syntactic complexity

SLA = second language acquisition

TTR = type-token ratio

# 1   Introduction

It is widely acknowledged in *second language acquisition* (henceforth SLA) research that more proficient *second language* (henceforth L2) writers are simultaneously "more fluent, accurate, and complex in their writing than less proficient writers" (Wolfe-Quintero, Inagaki, and Kim 1998, 4). There are various ways to measure *L2 proficiency*. However, the connection between L2 proficiency and *complexity*, *accuracy*, and *fluency* (henceforth CAF) is well-established and, therefore, the CAF framework remains a prominent approach. Furthermore, CAF measures are not only used to investigate L2 proficiency but also *L2 development* (Wolfe-Quintero, Inagaki, and Kim 1998). Complexity, one of the elements of CAF, relates to how elaborated the language produced by the learner is (Ellis and Barkhuizen 2005, 139). This thesis is interested in the relationship between linguistic complexity, possible variables that affect it, for example, different language backgrounds and socio-cultural factors, and in the manifestation of syntactic structures and lexical elements in essays of varying levels of complexity. Therefore, the present mixed methods study compares *syntactic* and *lexical* complexity (henceforth SC and LC, respectively) in essays written by Finnish and Hungarian university-level English language students.

L2 complexity has been researched considerably in the past few decades. For example, De Clerq and Housen (2017) researched the relationship between SC development and general L2 development. They found that SC increases when the learner's language develops. Research by Kuiken and Vedder (2019, 193) suggests that there is variation in SC "across proficiency levels, across languages, and between L2 and L1". Likewise, LC has been researched before. For example, De Clerq (2015) mapped out the development of LC in L2 English and L2 French. Kim, Crossley, and Kyle (2018) further explored the connection between LC, lexical proficiency, development, and writing quality. Additionally, a study by Abdi Tabari, Lu, and Wang (2023) investigated the impact of complexity of tasks on LC.

Because of the multitude of studies that have been conducted in the past decades, there already exists a rich understanding of how complexity and its dimensions manifest in learner language. Nevertheless, the existing studies tend to focus on either syntactic or lexical complexity and seldom include both aspects in the same analysis. Consequently, they cannot necessarily paint a descriptive and holistic picture of L2 complexity as the research is limited in terms of different types of complexities. Neither can the existing research efficiently comment on the relationship between syntactic and lexical complexity. It is important to note

that the same sample of learner language can vary in terms of different types of complexities. Additionally, the possible effect of the *first language* (henceforth L1) of the learner and socio-cultural factors on L2 complexity have still not been studied thoroughly as not many studies focus on them exclusively but merely comment on the phenomena superficially. Hence, it is essential to focus the analysis on these areas which have previously been neglected by research. This ensures that the present study will bring forth new and important knowledge that can help us understand the variability in L2 complexity.

This study compares Finnish and Hungarian university-level English language students in terms of SC and LC by analysing essays written by the two groups. University students of English were determined suitable for the analysis because their level of English as an L2 was thought to be appropriate for investigating complexity. The decision to analyse SC and LC was motivated not only by the lack of earlier research featuring both types of complexity but also by the fact that these two elements would provide a wider understanding of linguistic complexity in L2 proficiency. What makes the comparison in this study especially remarkable are the two groups that are compared in the analysis. The two groups are Finnish and Hungarian university-level English language students whose L1s are Finnish and Hungarian, respectively. These two L1s will be the starting point for the comparison. Finnish and Hungarian are languages that have some similarities as they are related but also differ greatly as the relation is distant. This creates a fascinating and complex setting for the analysis of the thesis.

The analysis of the essays was conducted with the help of web-based linguistic complexity analysers: *The L2 Syntactic Complexity Analyzer* (henceforth L2SCA) (Lu 2010, Lu 2011, Ai and Lu 2013, and Lu and Ai 2015) and *The Lexical Complexity Analyzer* (henceforth LCA) (Ai and Lu 2010 and Lu 2012). The results given by the analyser tools were further analysed quantitatively. A quantitative analysis was chosen as the methodological starting point because it allows for a large amount of data to be analysed easily, which may help the generalisability of the analysis. The primary material for the analysis comes from the ACALEX corpus. The thesis aims to investigate the differences and/or similarities in linguistic complexity between the two groups to explore what factors explain the observed differences and/or similarities, and examine the manifestation of syntactic structures and lexical elements in the essays. Therefore, the thesis will answer the following research questions:

1A. How do English texts written by Finnish and Hungarian university students compare in terms of syntactic complexity?

1B. How do English texts written by Finnish and Hungarian university students compare in terms of lexical complexity?

2. How do syntactic structures and lexical elements manifest in observed high-complexity and low-complexity essays?

Thus, the goal of this master's thesis is to compare the L2 complexity of two groups, to find potential factors that affect the similarities and/or differences between the two groups, and to create an understanding of the structures and elements in essays of varying levels of complexity. The research will shed light on the L2 complexity with a more holistic approach and theorise the potential effects of language background and socio-cultural factors. As the thesis will bring forth more knowledge about written learner language, the results of the study may be beneficial for professionals in the teaching field. New knowledge about how learner language manifests in writing may help teachers understand the challenges faced by L2 learners and, consequently, be able to instruct L2 learners better. Furthermore, the results of this study may help to guide the development of language teaching curricula in the direction where L2 learners will be provided with more apt instruction. This will, of course, be beneficial on a grander scale as this may contribute to SLA as a whole.

This thesis consists of six sections, some of which also have subsections. The following section, Section 2, focuses on the theoretical background of the thesis, and, thus, aims to define the necessary terms and concepts in the analysis and discusses prior research in the field. Section 3 introduces the materials and methods used in this thesis. Section 4 explores the results of the study. Section 5 discusses the results in terms of how they relate to previous studies and any ideas and questions which may arise from the analysis. Finally, Section 6 will conclude everything discussed in the thesis.

## 2   Theoretical background

In this section, a comprehensive review of the relevant literature on the themes of the thesis will be provided and the necessary terminology will be defined. First, in Subsection 2.1, the key term of the thesis, linguistic complexity, will be defined and its connection to L2 proficiency and development will be discussed in general. Then, both SC and LC will be defined, and they will be discussed in more detail with a special focus on earlier research on the phenomena. In Subsection 2.2, the effect of cross-linguistic influence on linguistic complexity and related studies will be discussed briefly. Finally, in Subsection 2.3, Finland and Hungary will be discussed from two different viewpoints: Finnish and Hungarian as L1s, and English as an L2 in both countries.

### 2.1   Linguistic complexity in L2 proficiency and development

In this subsection, the main objective is to define linguistic complexity in SLA and to discuss how it relates to L2 proficiency and development. However, complexity is a challenging phenomenon to define as there is no one commonly accepted definition of it (Bulté and Housen 2012, 22). This lack of standardisation poses challenges because one must be vigilant when reading about the phenomena. Additionally, Ehret and Szmrecsanyi (2019, 24) observe that terms such as 'complexity' and 'complex' are often used to refer to different phenomena like "complexity of linguistic structures but also second language (L2) acquisition difficulty". This may easily lead to further confusion around the definition. That is why it is crucial to define these terms clearly to avoid any confusion. Dahl (2004, 25–26) makes a distinction between *relative* and *absolute* complexity. On one hand, relative complexity concerns how language learners perceive the systems and features of a language. On the other hand, absolute complexity concerns the objective reality of the complexity in a language. This thesis takes an absolute approach to complexity.

Fortunately, many efforts to define the multiplex term that is complexity have been made. For example, Bulté and Housen (2014, 46) define complexity as "absolute, objective, and essentially quantitative property of language units, features, and (sub)systems thereof in terms of (i) the number and the nature of discrete parts that the unit/feature/system consists of and (ii) the number and the nature of the interconnections between the parts". Furthermore, Ellis and Barkhuizen (2005, 139) define complexity more simply as "the extent to which learners produce elaborated language". Therefore, it is imperative to clarify that when these terms,

'complexity' and 'complex', are used in this thesis, they refer to absolute linguistic complexity. They do not refer to relative complexity or L2 acquisition difficulty, which highlight different phenomena in SLA.

Complexity is one of the three dimensions that form CAF, which is a framework used to analyse both L2 proficiency and development (Wolfe-Quintero, Inagaki, and Kim 1998). CAF is a relatively new framework. Skehan (1996, 46) first suggested separating learner goals into three main areas: accuracy, complexity, and fluency. CAF can also be thought of as dimensions of performance, where CAF is first divided into two: fluency and form. Form can further be divided into accuracy and complexity (Skehan and Foster 2001, 190). CAF has been criticised as an approach to measure L2 performance. For example, Pallotti (2009) makes a notion that meaning-wise nonsensical sentences such as "colorless green ideas sleep furiously on the justification where phonemes like to plead vessels for diminishing our temperature" could achieve high scores using CAF measures. Thus, Pallotti suggests that adequacy be also considered in CAF research.

When discussing complexity, Skehan (1996) also discusses a process he calls *restructuring*. Skehan sometimes even refers to complexity as 'complexity-restructuring' or opts to use the term 'restructuring' to refer to the phenomenon as a whole (ibid.). Restructuring means that when learners realise the limitations of the language they produce and that it may need some modifications, "they are more likely not simply to use more complex language but also to attempt to pressure their own language systems" (Skehan and Foster 2001, 191). Skehan and Foster (2001, 190) note that learners may vary in terms of how willing they are to take risks and try new things when producing learner language. Additionally, learners may prioritise different dimensions of CAF over the other elements. This is known as the *Trade-off Hypothesis* (Skehan 2009, 511). Thus, according to the hypothesis, if a learner prioritises complexity it may negatively affect the other dimensions, accuracy and fluency, leading to errors or lessened fluency. According to Mehnert (1998), in oral language, learners tend to prioritise accuracy and fluency over complexity. Like other aspects of learner language, linguistic complexity develops over time when language skills progress. It is hypothesised that language produced by beginners is usually simpler in terms of linguistic complexity than that of more advanced learners (Kisselev et al. 2022, 798).

Linguistic complexity is a well-researched phenomenon, and many of its aspects have been studied comprehensively. For example, the link between the nature of the task the learners

perform when producing language and linguistic complexity has been researched. For example, Lee (2021) researched the effect of genre on linguistic complexity and found that argumentative essays tend to be more linguistically complex than descriptive essays. The relationship between task complexity and linguistic complexity is a notable area in L2 complexity research. Robinson (2001) claims that higher task complexity leads to higher linguistic complexity. This means that pushing learners to complete more difficult tasks will result in the learners producing more elaborated language. This is known as the *Cognition Hypothesis*. More recent research has challenged the hypothesis. For example, Allaw (2021) found evidence that task difficulty may contribute to linguistic complexity but acknowledged that other factors may also have as significant an effect. Kuiken and Vedder (2011) found that task complexity affects mainly accuracy, another dimension of CAF. Furthermore, in another study by Kuiken and Vedder (2012), no correlation between task difficulty and linguistic complexity was found.

Different types of linguistic complexities have been considered in research and different ways of classification have been suggested. For example, Ellis and Barkhuizen (2005, 152) list five different types of complexity: interactional, propositional, functional, grammatical, and lexical. Another suggested classification has four types of complexity: lexical, morphological, syntactic, and phonological (Bulté and Housen 2012). This thesis focuses on two types of complexity: syntactic and lexical. Between SC and LC, the former has been researched more thoroughly (De Clerq and Housen 2017, 317). Earlier research on both types of complexities explored in this thesis will be discussed more in the following subsections.

### 2.1.1 Syntactic complexity

In this subsection, SC will be defined and discussed considering existing literature. As with discussing the phenomenon of complexity in general, defining SC also has its challenges. It is important to highlight the ambiguity around the terms 'syntactic complexity' and '*grammatical complexity*'. Grammatical complexity means that the learner has "a wide variety of both basic and sophisticated structures" available to them and that they can access these structures quickly (Wolfe-Quintero, Inagaki, and Kim 1998, 69). 'Syntactic complexity' is a closely related term, which refers to the complexity of the syntactic structures in the language used. These two terms are very frequently used interchangeably. Often, for example, in research, the term 'grammatical complexity' may be used, but the measures used in the methodology are purely syntactic in nature. Some measures are not necessarily purely

syntactic but can be used to measure grammatical complexity. For example, Wolfe-Quintero, Inagaki, and Kim (1998, 80) describe how counting types and numbers of pronouns and articles can be used to measure grammatical complexity. However, the term 'syntactic complexity' is generally more commonly used to refer to the phenomenon. For the sake of this thesis, the term 'syntactic complexity' will be used to advance clarity and consistency and as the measures used to calculate the type of complexity are syntactic in nature in the analysis of this thesis.

SC may manifest in different ways in writing and for that reason, there are also myriad different ways to measure SC. However, some indices are more common than others. SC is most often measured by the amount of subordination. The idea is that the more subordination, the more syntactically complex learner language is (Ellis and Barkhuizen 2005, 139–140). Thus, a high number of dependent clauses in writing is linked to higher SC. With this logic a sentence such as "I like books because I like reading" would be considered more syntactically complex than a sentence such as "I like books and I like reading". The two sentences are almost identical but the way they connect the clauses is an important difference. The first sentence displays subordination and the second sentence displays coordination. However, the amount of coordination can also be used to measure SC even if it is not used as frequently. SC also manifests in the length of different structures produced. When examining the length of structures and measuring SC, the general logic is that the longer the structure, the more complex it is. For example, a sentence such as "I often read books before bed" would be considered more syntactically complex than a sentence such as "I read books" because of the length of production.

One important element used to analyse SC is the *T-unit* (*minimally terminable unit*). A T-unit contains "one main clause with all the subordinate clauses attached to it" (Hunt 1965, 20). Norris and Ortega (2009, 560) argue that the T-unit may be ideal for measuring intermediate or advanced written data. It has been observed that longer T-units tend to be more common in texts written by more proficient L2 learners (Kyle and Crossley 2018, 334). The T-unit serves as the base for many SC indices, for example, the mean length of T-unit (MLT), verb phrases per T-unit (VP/T), clauses per T-unit (C/T), dependent clauses per T-unit (DC/T), T-units per sentence (T/S), complex T-unit ratio (CT/T), coordinate phrases per T-unit (CP/T), and complex nominals per T-unit (CN/T). It has been concluded that mean length of T-unit, mean length of clause, clauses per T-unit and dependent clauses per clause may be the most

accurate measures for SC (Ortega 2003, 493). SC measures will be discussed in more detail in the methodology section of this thesis concerning the SC analyser tool used.

SC in learner language has been researched considerably, and the majority of studies that focus on linguistic complexity tend to focus on SC over other types of complexities: "[t]he most commonly used measures of complexity are grammatical in nature" (Ellis and Barkhuizen 2005, 154). Differences between learners and native speakers have been mapped out in several studies, for example, Kuiken and Vedder 2019, Lan et al. 2022, and Zhang 2022. All three studies found differences in the complexity of the grammar between texts written by L1 and L2 speakers. Bulté and Housen (2014) mapped out the development of linguistic complexity in L2 writing. In their study, SC development manifests as "a significant increase in the length of linguistic units at all levels of syntactic organization" (Bulté and Housen 2014, 53). This suggests that more proficient learners write longer phrases, clauses, sentences, and T-units. Furthermore, Norris and Ortega (2003, 563) argue that higher phrasal elaboration is linked to a higher level of language development. Additionally, more specific and detailed aspects of syntactic complexity have been researched. For example, Sarte and Gnevsheva (2022), Díez-Bedmar and Pérez-Paredes (2020), and Liu and Li (2016) have shed light on noun phrase syntactic complexity. The studies argue that more proficient writers use more noun modifiers. Sarte and Gnevsheva (2022, 10) note that the topic of an essay being written may affect noun phrase complexity.

## 2.1.2 Lexical complexity

In this subsection, the definition of LC and the relevant literature on the phenomenon will be discussed. Contrary to SC, the definition of LC is more straightforward and less debated although similarly multiplex. Therefore, LC means that the learner has "a wide variety of both basic and sophisticated words" available to them and that they can access these words quickly (Wolfe-Quintero, Inagaki, and Kim 1998, 101). Hence, LC is not simply equal to the usage of advanced vocabulary but one also needs to know basic vocabulary. Furthermore, LC can be used to describe the effectiveness of the L2 learner's communication skills in both spoken and written form (Lu 2012, 252). This type of complexity manifests in writing "primarily in terms of the range (lexical variation) and size (lexical sophistication) of a second language writer's productive vocabulary" (Wolfe-Quintero, Inagaki, and Kim 1998, 101).  It is relevant to mention that LC differs from phenomena such as *vocabulary breadth*, which refers to "the number of words a language learner knows", and *vocabulary depth*, which refers to "how well

these words are known" (Harkio and Pietilä 2016, 1079). Thus, vocabulary depth is more clearly separated from LC as it does not consider the complexity of the learner's vocabulary but what type of knowledge the learner has about individual words. Vocabulary breadth is closer to LC as aspects of LC also relate to vocabulary size. However, vocabulary breadth is much less multiplex as a phenomenon as it is more focused on the literal number of words rather than whether the learner knows a wide variety of different types of words.

LC can be measured in various ways. One commonly used approach is to consider LC as consisting of three main dimensions: *lexical density*, *lexical diversity*, and *lexical sophistication*. Bulté and Housen (2012, 28) suggest a fourth dimension: *lexical compositionality*. Nevertheless, LC is generally seen as a multidimensional phenomenon. Lexical density refers to "the ratio of the number of lexical (as opposed to grammatical) words to the total number of words in a text" (Lu 2012, 191). *Lexical words*, which are also known as *content* words, are words that have semantic meanings, and *grammatical* words, which are also known as *function words*, are used to express grammar and sentence structure (Biber, Conrad and Leech [2002] 2019, 15–16). Lexical diversity is often measured by examining *type-token ratios* (henceforth TTR). TTR ratio is measured "by dividing the number of different words (word types) by the total number of words (word tokens)" (Cunningham and Haley 2020, 711). It is important to note that lexical diversity has been criticised. In a study by Yoon and Polio (2017, 288), it was found that the lexical diversity of the same L2 writer may vary across genres. Lexical sophistication, which is also known as *lexical rareness*, measures "the proportion of relatively unusual or advanced words in the learner's text" (Read, 2000, 203). Lexical compositionality considers "the number of formal and semantic components of lexical items" such as morphemes (Bulté and Housen 2012, 28). LC measures will be discussed in more detail in the methodology section of this thesis concerning the LC analyser tool used.

The multidimensional nature of LC has been researched notably although not as much as SC. It has been suggested that LC increases steadily as the L2 learner's proficiency level increases (e.g., Kisselev et al. 2022, Kim, Crossley, and Kyle 2018, and Alexopoulou et al. 2017). Crossley and McNamara (2012, 120) examined a corpus by using automated indices of cohesion, lexical sophistication, syntactic complexity, and conceptual knowledge, and found that lexical sophistication is generally the strongest predictor in detecting the L1 of a writer. Qin and Uccelli (2020, 10) suggest that the LC of a writer may vary across genres as they found that *English as a foreign language* (henceforth EFL) "learners' academic texts were

more lexically diverse than their colloquial texts". Additionally, Rahayu, Utomo, and Setyowati (2021, 259) made an important discovery when the participants of their study achieved higher scores in LC in their L2 compared to their L1. The researchers suggest that this may be because the use of dictionaries or thesauri may be more common in L2 writing than in L1 writing ending in higher LC in L2.

## 2.2  Cross-linguistic influence on linguistic complexity

This subsection aims to briefly discuss the existing literature on the relationship between different L1 backgrounds and linguistic complexity. The cross-linguistic aspect of linguistic complexity has not been researched significantly. However, some research has been made on differences between linguistic complexity in different L2s. For example, De Clercq and Housen (2017) studied how SC manifests in L2 English and French produced by L1 speakers of Dutch and found evidence of differences in SC. Furthermore, Bernardini and Granfeldt (2019) mapped out variation in Swedish L1 learners of English, French, and Italian and learnt that variation between the target languages affects linguistic complexity in learner performance. The results suggest that different L1 and L2 combinations may contribute to linguistic complexity in different ways (Bernardini and Grandeldt 2019, 226). This means that learners of an L2 may have differences in how linguistic complexity manifests because of their differing L1 backgrounds and that learners with the same L1 will have similarities in this regard.

Furthermore, Ehret and Szmrecsanyi (2019) studied learners from different L1 backgrounds and suggest that L1 is a good predictor for L2 complexity. Similarly, Lu and Ai (2015, 26) suggest that there may be differences in the development of SC between learners with different L1 backgrounds. This suggests that the L1 background of an L2 learner may help explain how SC manifests in their learner language. Phuoc and Barrot (2022) suggest similar results to Lu and Ai. Crossley and McNamara (2012) studied whether they could detect the L1 of a learner by using automated tools which analysed cohesion, lexical sophistication, syntactic complexity, and conceptual knowledge in the learners' texts. They found that automated tools could detect the L1 with varying degrees. They also discovered that Finnish learners of English show high lexical sophistication but average SC (Crossley and McNamara 2012, 121).

Ortega (2003) explored the impact of instructional setting and proficiency on SC and found that *English as a second language* (henceforth ESL) learners tend to produce learner language

with higher SC than their EFL counterparts. Furthermore, Barrot and Gabinete (2021) and Zhang and Kang (2022) have researched differences between ESL and EFL learners using CAF measures. The former study found that both the skill level and the L1 of the learner influence writing proficiency. Most importantly the results indicate that "learners of different L1 background may not develop the same way in various measures of CAF despite their similar proficiency level" (Barrot and Gabinete 2021, 227). The latter study suggests that language learning contexts, English language education, and teaching strategies may influence the use of syntactic features in writing.

The difference between ESL and EFL learners is not necessarily always a matter of L1 but it is highly related to culture and the approach towards English learning. There may be differences in instruction based on whether English is considered a second or a foreign language. However, it is crucial to note that the difference between ESL and EFL is somewhat ambiguous as these terms are frequently used interchangeably. Thus, it is also not always sensible to group all ESL into one group and all EFL learners into another group.

## 2.3  Finland and Hungary

The primary purpose of this subsection is to introduce Finland and Hungary. The majority of the focus will be on the national languages and the educational systems of the two countries as these themes are crucial to the analysis of this thesis.

### 2.3.1  Finnish and Hungarian as L1s

Finnish and Hungarian, which are mainly spoken in Finland and Hungary, respectively, share linguistic heritage and both belong to the Finno-Ugric language group (KOTUS n.d.). This means that Finnish and Hungarian have developed into their own languages from the same proto-language. Thus, it is then evident that these languages have some similarities linguistically. For example, in both languages primary word stress is on the first syllable, both languages have a similar vowel harmony system, and both languages are agglutinative in nature (Abondolo and Valijärvi 2023, 17, 20, 28–32). These are just some examples of the similarities between Finnish and Hungarian but they also illustrate how the two languages differ from English, as English has more variation in word stress, has no vowel harmony, and is analytic in nature.

However, there is a significant amount of linguistic distance between Finnish and Hungarian as the two languages are not closely related. Within the Finno-Ugric language group, Finnish belongs to the Finnic branch (KOTUS n.d.) while Hungarian belongs to the Ugric branch (Abondolo and Valijärvi 2023, 11). Within the same language group, Finnish is more similar to the languages in the same branch, for example, Estonian and Karelian, than to Hungarian (KOTUS n.d.). Likewise, Hungarian is more similar to languages in the Ugric branch, Khanty and Mansi (Abondolo and Valijärvi 2023, 5). Finnish and Hungarian are not mutually intelligible, and similarities between the two languages, although evident, are not necessarily easy to detect.

### 2.3.2 English as an L2 in Finland and Hungary

This subsection aims to explore Finland and Hungary from the viewpoint of English as an L2 in both countries. The educational systems of the two countries regarding foreign language teaching will also be discussed in this subsection. It is important to discuss the impact of socio-cultural factors on L2 proficiency and development as language learning does not happen in a vacuum, but it is affected by a multitude of different variables. Finland and Hungary possess some cultural similarities when it comes to the English language. For example, neither Finland nor Hungary lists English as an official language. Instead, Finland's official languages are Finnish and Swedish (KOTUS n.d.) and Hungary's official language is Hungarian (EU n.d.). This means that, to be precise, English is taught as a foreign language rather than as a second language in both countries. Furthermore, neither Finland nor Hungary has a significant minority group of native English speakers (European Commission 2012). Thus, if we consider the countries with Kachru's (1985) *three concentric circles of English* in mind, we find that both countries belong in *the third or expanding circle*. Additionally, it is appropriate to mention that both Finland and Hungary are member states of the Council of Europe and therefore both countries are encouraged to use *The Common European Framework of Reference for Languages* (CEFR) as a tool in language instruction.

The two countries' *National Core Curriculums* (henceforth NCC) have similar regulations about when the teaching of the first foreign language should start. The current Finnish NCC states that teaching should start already during the first grade in primary school (OPH 2019). This is a relatively new regulation and before the reform, teaching used to start during the third grade (OPH 2014). The current Hungarian NCC states that the teaching of the first foreign language should start no later than during fourth grade (Igazságügyi Minisztérium

2020). Neither NCC states that English has to be the first foreign language to be studied. However, during the 2021–2022 school year in Finland, 91% of first graders studied English as a foreign language (SUKOL n.d.). In comparison, during the same school year in Hungary, 79% of the time a student was learning a foreign language, they were learning English (KSH n.d.). Thus, it is safe to say that English is the most frequent foreign language learnt in schools in both countries. This may be due to the position of English as a global lingua franca.

Naturally, there are differences between the two countries as well. According to the 2012 Eurobarometer (European Commission 2012), seventy per cent of Finns speak English well enough to have a conversation, while only twenty per cent of Hungarians can do the same. This suggests that the overall proficiency level of Finns is higher than that of the Hungarians. Additionally, EF (2022) reported that out of 111 countries, Finnish speakers of English rank seventh and Hungarian speakers of English rank 18th in proficiency and, thus, the English speakers from the countries have "very high proficiency" and "high proficiency", respectively. Still, this report has its flaws as it is based on an online test, which is available to anyone, but which attracts mostly people who are interested in language learning and thus, the report is not effective in describing the two groups as a whole. A study by Thi, Van Do, and Nikolov (2023) compared the syntactic complexity of Hungarian learners of English to that of Myanmar learners' of English. According to the study, Hungarian learners show higher syntactic complexity when compared to their Myanmar counterparts (Thi, Van Do and Nikolov 2022, 137). Nevertheless, it is suggested that, overall, Finns may have higher proficiency in English than Hungarians. Even so, there is no significant research on whether there would be any difference when comparing Finns and Hungarians of the same academic level.

# 3   The present study

The goal of this section is to introduce the present study by examining the materials and methods used. In Subsection 3.1, the primary material of the thesis, the corpus, will be discussed. Subsequently, in Subsection 3.2, the discussion shifts to the automatic analysers that were used in the analysis. Finally, in Subsection 3.3, the rest of the analysis process will be described. However, first, let us briefly examine the research questions and discuss the hypotheses that arise. Research questions 1A and 1B relate to comparing Finnish and Hungarian university-level English language students by quantitatively comparing the groups' performance in SC and LC. Research question 2 brings a more qualitative angle to the present study by shedding light on what syntactic structures and lexical elements look like in texts with different levels of complexity.

The initial hypothesis is that the two groups will differ as there is some existing research which suggests that there may be a link between linguistic complexity and the L1 background of the learner (e.g., Bernardini and Grandeldt 2019, Lu and Ai 2015, and Phuoc and Barrot 2022). The difference may be that one group shows higher SC or LC. However, how significant that difference would be, is complicated to hypothesise as there is not enough prior research on the topic to make an educated estimation. The two groups have not previously been compared this way, so it is also difficult to estimate what kinds of differences there could be. One study suggests that Finnish learners display perhaps higher than average LC but perform averagely in terms of SC (Crossley and McNamara 2012, 121). The hypothesis will be further discussed in the discussion.

## 3.1   Corpus

Using corpora in SLA research is very common because of its many advantages. The convenience of corpora stems from the fact that corpora are existing large data sets which were created specifically for research purposes. Using large data sets in research may help the generalisability of a study and using an already existing data set saves time in the research process. The primary material used for the present study is the ACALEX corpus. The ACALEX corpus is part of *the Academic Lexis in L2 Speech and Writing* project. The corpus contains material from students from three European universities: the University of Turku, Åbo Akademi University, and the University of Szeged. The former two universities are based in southwestern Finland in Turku, and the latter university is based in southern Hungary

in Szeged. The data of the corpus were collected during three years from 2016 to 2018. In total, there are 175 samples of learner language of which 106 samples were produced by major and minor students of English. The corpus consists of both written data and spoken data, and some personal information about some of the participants. The identities of the learners are well-protected in the corpus as each learner is referred to only by an identification code such as ACALEX-1101. Thus, the participants are not identifiable to a person using the corpus. The ACALEX corpus has not been widely used in research as it is not available for free use online. However, it has been used previously, for example, in earlier master's theses (e.g., Mäkynen 2023 and Peltomäki 2018).

The ACALEX corpus presents many research opportunities. The specific part of the corpus used in the present study was perceived as meaningful for research and thus, the idea for the present study was born from the corpus. The primary data used for this thesis are the thirty essays written by Hungarian university-level English language students and 32 essays by Finnish students of English from the University of Turku found in the corpus. The texts by Hungarian students of English were written in 2016 at the University of Szeged for an English course. The texts by Finnish students were written in 2017 at the University of Turku for a course part of basic studies in the English language. Thus, all participants of the study are either major or minor students of English. To have equal numbers for the comparison, two texts by Finnish students were removed from the analysis. This was taken as an opportunity to remove the two shortest texts which were considered potential outliers for the analysis. The topic of all the essays is "My Professional Future the Way I See it at the Moment". Thus, in all the essays, the students describe the potential jobs and professions they may do in the future concerning their current studies. Advantageously, the essays do not vary topic-wise, which makes the comparison between the two groups of essays more sensible. It has been suggested that genre and the choice of topic may have an influence on linguistic complexity (e.g., Sarte and Gnevsheva 2022, Lee 2021, Yoon and Polio 2017, and Qin and Uccelli 2020). Thus, it is beneficial that all the essays are of the same genre and topic.

## 3.2   Automatic analysers

Automatic analyser tools are the main method used in the analysis of the present study. The decision to use automatic analysers was driven by the fact that it allows a large amount of data to be analysed conveniently and quickly. Additionally, it ensures that the subanalyses be free from human error and mistakes. Most existing studies use automatic analysers, but the

choice of the automatic analyser tool varies from study to study with some studies using multiple analyser tools. The two most cited developers of both syntactic and lexical automatic analysers are Xiaofei Lu and Kristopher Kyle. The two chosen tools for the analysis of the present study are L2SCA and LCA, which were developed by Lu. However, analysis tools TAASSC (*Tool for the Automatic Analysis of Syntactic Sophistication and Complexity*), TAALES (*Tool for the Automatic Analysis of Lexical Sophistication*), and TAALED (*Tool for the Automatic Analysis of Lexical Diversity*) by Kyle were also considered for the analysis. The three tools are a little newer compared to the tools by Lu and feature a wider range of more detailed measures. However, for the present study, the programs by Lu were thought to suffice as they lean on well-established theories, research on their accuracies has been promising, and they are overall simpler and more user-friendly. Furthermore, Kyle's tools have not yet been evaluated as widely in terms of reliability.

Thus, both analyser tools used in the present study were developed by Lu. It was deemed beneficial to use analysers essentially from the same source as this would create a balance between the two subanalyses, SC analysis and LC analysis. The two analyser tools are often grouped together and sometimes even referred to as simply 'SynLex'. Thus, it is not uncommon for studies to use both as a basis for their analysis. For example, Rahayu, Utomo, and Setyowati (2021) used the two analyser tools when comparing SC and LC in L1 and L2 writing. The indices used in the analyser tools are largely based on classical measures most of which were reviewed by Wolfe-Quintero, Inagaki, and Kim (1998). For this thesis, the web-based interfaces of the analysers were used.[1] The decision to use the web-based analysers instead of the downloadable programs was influenced by the fact that the web-based tools do not require *POS tagging* (POS = *part of speech*) or *lemmatisation* by the user unlike the downloadable programs. Lemmatisation refers to a process where a lemma which consists "of a headword and some of its inflected and reduced (n't) forms" (Nation 2001, 7) is reduced to its basic form. All required processes such as POS tagging and lemmatisation are integrated into the web-based interfaces of the programs and the user just needs to enter raw text into the analysers.

---

[1] The wed-based interfaces for L2SCA and LCA are available for free use online at https://aihaiyang.com/software/.

### 3.2.1 The L2 Syntactic Complexity Analyzer

This subsection describes the computational tool used to measure SC in the analysis of the present study and briefly highlights some studies that have utilised the tool. L2SCA is an automatic analyser tool that uses 14 measures to analyse SC in written text. The analyser tool was created because, according to Lu (2010, 476), most existing computational tools used measures that were primarily suitable for researching L1 acquisition. This lack of computational tools focusing on SC in SLA inspired the creation of L2SCA, which is meant specifically for analysing learner language. Most of the indices used in L2SCA are described by Wolfe-Quintero, Inagaki, and Kim (1998, 82–99). The indices are often described as 'classical' or 'traditional' measures in SC research.

Table 1. Description of L2SCA indices (Lu 2010, 479 and Kyle 2016, 54)

| Category | Index Abbreviation | Index Name | Index Description |
|---|---|---|---|
| length of production | MLC | mean length of clause | number of words per clause |
| | MLS | mean length of sentence | number of words per sentence |
| | MLT | mean length of T-unit | number of words per T-unit |
| sentence complexity ratio | C/S | clauses per sentence | number of clauses per sentence |
| amount of subordination | C/T | clauses per T-unit | number of clauses per T-unit |
| | CT/T | complex T-unit ratio | number of complex T-units per T-units |
| | DC/C | dependent clauses per clause | number of dependent clauses per clause |
| | DC/T | dependent clauses per T-unit | number of dependent clauses per T-unit |
| amount of coordination | CP/C | coordinate phrases per clause | number of coordinate phrases per clause |
| | CP/T | coordinate phrases per T-unit | number of coordinate phrases per T-unit |
| | T/S | T-units per sentence | number of T-units per sentence |
| relationship between syntactic structures and production units | CN/C | complex nominals per clause | number of complex nominals per clause |
| | CN/T | complex nominals per T-unit | number of complex nominals per T-unit |
| | VP/T | verb phrases per T-unit | number of verb phrases per sentence |

The indices used by the analyser are described in Table 1. Nine of the L2SCA measures are considered to be large-grained measures, which "provide analysis at the clausal and/or

sentence level", and five are considered fine-grained, which provide more detailed analysis (Kyle 2016, 48). The L2SCA measures are included in Kyle's competing SC analyser program, TAASSC, which features also fine-grained measures of SC unique to the program. According to Lu (2010, 479), the L2SCA indices can be categorised into five different categories. The first three indices in Table 1, MLS, MLT, and MLC, measure the "length of production at the clausal, sentential, or T-unit level" (ibid.). The second category consists of just one index, C/S, which measures the sentence complexity ratio. C/T, CT/T, DC/C, and DC/T are the four indices that comprise the third category, which measures the amount of subordination (ibid.). The second measure in the category in the table, CT/T, features a *complex T-unit*, which is a T-unit that contains at least one subordinate or embedded clause (Casanave 1994). The fourth category measures the amount of coordination and consists of three indices, CP/C, CP/T, and T/S (Lu 2010, 479). Finally, the fifth category consists of the last three indices in Table 1, CN/C, CN/T, and VP/T, and measures "the relationship between particular syntactic structures and larger production units" (ibid.). Two of the measures feature a *complex nominal*, which consists either of a noun plus adjective, possessive, relative clause, prepositional phrase, participle, or appositive, nominal clauses or gerunds and infinitives in subject position (Cooper 1976, 180 and Lu 2010, 483).

Even if L2SCA is a very common tool in learner language research, it is somewhat rare for SC studies to include all 14 measures in the analysis. However, Larsson and Kaatari (2020) utilised all 14 measures available in the analyser when analysing the relationship between SC and formality in writing. Furthermore, Lei, Wen, and Yang (2023) also included all 14 measures in their study. Most studies that use L2SCA do not use all the indices but choose to use the indices that are best fit for the study in question. For example, Yoon (2017), Kim and Crossley (2018) and Jiang, Bi, and Liu (2019) all chose seven measures for their studies, but the measures vary between the studies, although five measures are used in all three (MLC, MLS, MLT, CP/C, and CN/C).

The reliability of L2SCA has also been researched. Lu (2010, 487) reported that L2SCA "achieves a very high degree of reliability for identifying" units and structures in texts. However, later Lu (2011, 57) studied the reliability of L2SCA and reported that four of the 14 indices (C/T, CT/T, T/S, and VP/T) measure L2 development poorly as "they provide no useful information about a learner's proficiency level". Thus, he recommends using the remaining ten in future studies of complexity. The analysis of the present study adheres to this recommendation and uses the ten remaining indices in the analysis.

Contrary to prior research, Hwang and Polio (2023) suggest that choice of topic, genre, or text length may not have a significant influence on SC when analysed by computational tools. In their study, they tested different computational tools used to calculate SC, including L2SCA, and found that the tool can measure SC in L2 reliably with variability in the choice of topic, genre, and text length. However, they recommend using texts longer than one hundred words for analysing SC as they found that shorter texts could be more easily affected by choice of topic, genre, and text length.

### 3.2.2  The Lexical Complexity Analyzer

This subsection describes the computational tool used to measure LC in the analysis of the present study and briefly highlights some studies that have utilised the tool. LCA is an automatic analyser tool that uses 25 measures to analyse lexical complexity in written text. Around half of the indices are described by Wolfe-Quintero, Inagaki, and Kim (1998). Lu (2012, 191) developed the tool as a response to the lack of computational tools available to automate some of the measures and to make LC analysis less labour-intensive.

Table 2. Description of LCA indices (Lu 2012)

| Dimension | Index Abbreviation | Index Name | Index Description |
|---|---|---|---|
| Lexical Density | LD | Lexical density | Lexical words to total words ratio |
| Lexical Sophistication | LS1 | Lexical sophistication-I | Sophisticated lexical words to total words ratio |
| | LS2 | Lexical sophistication-II | Sophisticated types to total types ratio |
| | VS1 | Verb sophistication-I | Sophisticated verb types to total verbs ratio |
| | CVS1 | Corrected VS1 | |
| | VS2 | Verb sophistication-II | Squared VS1 |
| Lexical Variation | NDW | Number of different words | Number of types |
| | NDWZ-50 | NDW (first 50 words) | Number of types in the first 50 words of sample |
| | NDW-ER50 | NDW (expected random 50) | Mean number of types of 10 random 50-word samples |
| | NDW-ES50 | NDW (expected sequence 50) | Mean number of types of 10 random 50-word sequences |
| | TTR | Type/Token ratio | Types to total words ratio |
| | CTTR | Corrected TTR | |
| | MSTTR-50 | Mean Segmental TTR | Mean type to token ratio of all 50-word segments |

| Dimension | Index Abbreviation | Index Name | Index Description |
|---|---|---|---|
| | RTTR | Root TTR | Types to square root of total words ratio |
| | logTTR | Bilogarithmic TTR | Bilogarithmic TTR |
| | Uber | Uber Index | Squared logarithm of tokens to logarithm of token-type ratio |
| | VV1 | Verb variation-I | Ratio of verb types to number of verbs |
| | CVV1 | Corrected VV1 | |
| | SVV1 | Squared VV1 | Squared VV1 |
| | LV | Lexical word variation | Lexical word types to total lexical words ratio |
| | VV2 | Verb variation-II | Verb types to total lexical words ratio |
| | NV | Noun variation | Noun types to total lexical words ratio |
| | AdjV | Adjective variation | Adjective types to total lexical words ratio |
| | AdvV | Adverb variation | Adverb types to total lexical words ratio |
| | ModV | Modifier variation | Adjective and adverb types to total lexical words ratio |

Table 2 describes the indices used in LCA. According to Lu (2012), LC consists of three dimensions: lexical density, lexical sophistication, and lexical variation. LCA aims to measure each dimension. Lexical density is measured by just one index, LD, the first index in Table 2. The next five indices in the table comprise the measures for the next dimension, lexical sophistication. In these measures, a word is considered 'sophisticated' if it does not belong to the 2,000 most frequent words in English (Lu 2012, 192). The third and final dimension, lexical variation, is measured by the last 19 indices. There are three measures in Table 2, CVS1, CTTR, and CVV1, which are 'corrected' versions of other measures. These corrected versions aim to fix the discrepancy created by the variation of sample length. For example, in the case of observing TTRs, "[a]s the sample increases, the probability of introducing new words decreases because the vocabulary that characterizes individuals at any given time is considered finite" (Fergadiotis 2011, 13). This creates an issue where sample length may affect the TTR in an undesirable way and which CTTR aims to fix.

LCA is a common tool in learner language research. Many studies use just some of the indices available in LCA. For example, Seidinejad and Nafissi (2018) utilised the tool to investigate whether teaching creative thinking techniques could affect LC in EFL writing with a focus on LD, LS, and LV. Additionally, Nasseri and Thompson (2021) analysed dissertation abstracts and used ten indices from LCA (LD, NDWERZ, NDWESZ, LOGTTR, UBER, VV1, CVV1, VV2, LV, and NV) alongside other computational tools. Furthermore, Tsai

(2021), utilised eight indices while studying the effect of corpus consultation in a writing course.

Little research has been done on the reliability of LCA, which is likely partially due to that LC in general is much less researched than SC. However, Lu (2012, 204) mapped out how the indices correlate with rankings made by teachers when researching LC in spoken learner language. He determined that the indices that function the best for assessing spoken language are NDW, NDW-ER50, NDW-ES50, CTTR, RTTR, MSTTR, SVV1, and CVV1, which are all measures of lexical variation. Lu (2012, 198) suggests that there may be "a difference in the role lexical sophistication plays in spoken and written proficiency" and, thus, the indices might function differently when analysing written language. Therefore, in this analysis, all 25 measures are included as there is little research on the reliability of them when researching written learner language.

## 3.3    Procedure

This subsection aims to describe the steps taken during the analysis of the thesis. Preparation for the empirical part of the thesis started with proofreading all sixty essays used for the analysis, which was necessary as the primary analysis involved using automatic analysers. This was done because a simple writing mistake, for example, a typo or lack of punctuation, could affect the results of the analysis. For example, let us consider one typo in one of the essays, 'scool', which is, when observed in context, an incorrect spelling of the word 'school'. Now in the same essay, 'school' is also spelt correctly. To a human reader, this variation in spelling is easily understood as a typo and would not affect our perception of LC in the essay. However, the automatic analyser does not recognise 'scool' as a typo of 'school' but counts them as two different types and would likely interpret 'scool' as a rare word. This would result in an inaccurate calculation of LC. Likewise, incorrect punctuation would affect the results in an unwanted way. For example, if a writer by mistake does not add a full stop at the end of a declarative sentence and immediately starts a new sentence, the automatic analyser tool interprets these two sentences to be one long sentence. This would have an unwanted effect on the SC score. Thus, proofreading involved fixing typos and adding punctuation where needed. Editing was kept minimal to not affect the results per se but to merely ensure that the automatic analysers would interpret the essays correctly.

Most essays needed some editing. This most often meant that a couple of commas were added or deleted. However, some essays needed a lot more editing. The essays in question might

have been written without using automatic proofreading as they were missing punctuation and had consistent spelling errors. What is important to note is that there was a significant difference in the number of required edits between the two groups. On one hand, the essays by Finnish students needed, on average, five edits. On the other hand, the essays by Hungarian students needed, on average, twenty edits. This was due to consistent writing errors made repeatedly by the students.

While reading each essay, it was also determined whether the essay was closer to standard British English or standard American English as this is something LCA asks before running the analysis. It was determined by examining the spellings of words and word choices. It is important to note that learner language often features characteristics from multiple standards of the target language. Thus, categorising the essays in the present study into British English and American English is somewhat problematic, too rigid, and not an accurate description of learner language. However, this is something that LCA asks the user to do when running the analysis. As a test, five of the essays were run through LCA first with British English chosen as the standard for the analysis and then with American English as the standard. The results were the same for all five essays no matter which standard was chosen.

After proofreading, the essays were run through both web-based automatic analysers, L2SCA and LCA. The results of these analyses were put into Microsoft Excel for further analysis. In Microsoft Excel, some basic descriptive calculations were made first. These included calculating the means, medians, and standard deviations for the results given by the automatic analyser tools. Additionally, the minimum and maximum scores for each index were examined at this stage of the analysis. The rest of the calculations for the analysis were made with IBM SPSS Statistics 27, which functions more effectively for statistical analysis. "A Guide to Doing Statistics in Second Language Research Using SPSS and R" by Larson-Hall (2016) was used as a general guide for conducting the statistical analyses in the thesis.

To choose an appropriate statistical test to calculate the statistical significance for the subanalyses, the normality of the data had to be analysed. This was done using a Shapiro-Wilk test as it has been suggested that it is one of the most effective for small sample sizes (Ricci 2005, 20). The Shapiro-Wilk test revealed that the data is not normally distributed, which meant that a non-parametric test for two independent samples was needed. Thus, the Mann-Whitney U test was chosen as recommended by Larson-Hall (2016, 286). The test was then used to calculate statistical significance of the subanalyses. After conducting the tests,

the p-values were examined, and the statistical significance was determined for each index. A result was deemed statistically significant if the p-value was under .05. After all the calculations were made, comparisons between the groups were made.

After completing the analysis for research questions 1A and 1B, the results were used to prepare for research question 2. To answer research question 2, four essays were chosen for more detailed and qualitative analysis. To conduct an extreme case analysis, the essays with the highest and lowest scores in SC and LC were chosen for further analysis. The essays were chosen by comparing their scores to the other essays and by examining the consistency of the high or low scores in the indices. After determining which essays exhibit the highest and lowest SC and LC, the essays were carefully analysed further by examining the syntactic structures and lexical elements in them. Finally, two examples from each essay were chosen to illustrate how the syntactic structures and lexical elements manifest in the essays.

# 4   Results

In this section, the results of the analysis of the present study will be discussed. First, in Subsection 4.1, the results of the SC analysis will be highlighted by discussing the results with the categories of L2SCA indices in mind and research question 1A (How do English texts written by Finnish and Hungarian university students compare in terms of syntactic complexity?) will be answered. After, in Subsection 4.2, the results of the LC analysis will be examined. Finally, in Subsection 4.3, the results of research question 2 will be discussed with examples from the essays.

## 4.1   Syntactic complexity analysis

It is beneficial to examine the average number of syntactic structures counted by L2SCA first before looking at the results of the indices. Table 3 describes the syntactic structures used in the essays and compares Finnish university-level English language students to Hungarian university-level English language students.

Table 3. Description of the syntactic structures in the essays

| Index | Minimum | | Maximum | | Mean | | Standard Deviation | |
|-------|---------|-----|---------|-----|--------|--------|--------|-------|
|       | FIN     | HUN | FIN     | HUN | FIN    | HUN    | FIN    | HUN   |
| W     | 488     | 477 | 703     | 909 | 557.17 | 554.5  | 52.09  | 73.48 |
| S     | 19      | 20  | 37      | 62  | 25.63  | 30.7   | 4.94   | 8.51  |
| C     | 43      | 53  | 75      | 108 | 57.87  | 68.47  | 8.02   | 10.57 |
| T     | 21      | 25  | 43      | 72  | 32.07  | 38.73  | 6.14   | 9.28  |
| CT    | 11      | 13  | 26      | 27  | 18.4   | 20.37  | 3.88   | 3.49  |
| VP    | 69      | 61  | 112     | 148 | 86.03  | 91.2   | 12.66  | 14.11 |
| CP    | 3       | 6   | 34      | 25  | 12.83  | 12.03  | 5.86   | 4.19  |
| DC    | 14      | 17  | 43      | 41  | 25.83  | 27.97  | 7      | 6.04  |
| CN    | 48      | 37  | 82      | 92  | 61.17  | 54.9   | 7.77   | 9.49  |

FIN = Finnish university-level English language students, HUN = Hungarian university-level English language students, W = words, S = sentences, C = clauses, T = T-units, CT = complex t-units, VP = verb phrases, CP = coordinate phrases, DC = dependant clauses, CN = complex nominal.

As illustrated in Table 3, the mean word counts of the essays are relatively similar between the compared groups. However, one cannot ignore the fact that Hungarian students have a larger standard deviation in terms of the word count of the essays compared to Finnish students. What is remarkable is that, on average, the essays written by Hungarian students contain more of almost all the syntactic structures measured by L2SCA than the essays

written by Finnish students. Essays written by Hungarian students, on average, have more sentences (S), clauses (C), T-units (T), complex T-units (CT), verb phrases (VP), and dependent clauses (DC). Essays written by Finnish students, on average, have more total words (W), coordinate phrases (CP), and complex nominals (CN), although the differences between the groups' mean total word counts and numbers of coordinated phrases are relatively small.

Table 4 displays the results of the syntactic complexity analysis and compares Finnish university-level English language students to Hungarian university-level English language students.

Table 4. Syntactic complexity results

| Index | Minimum | | Maximum | | Mean | | Standard Deviation | | Statistical Significance |
|---|---|---|---|---|---|---|---|---|---|
| | FIN | HUN | FIN | HUN | FIN | HUN | FIN | HUN | |
| MLC | 8.10 | 6.95 | 12.33 | 10.47 | 9.74 | 8.17 | 1.17 | .81 | <.001* |
| MLS | 15.24 | 13.74 | 27.89 | 26.52 | 22.31 | 18.98 | 3.46 | 4.28 | .002* |
| MLT | 13.12 | 10.10 | 24.90 | 20.92 | 17.87 | 14.77 | 3.04 | 2.63 | <.001* |
| C/S | 1.54 | 1.55 | 2.95 | 3.55 | 2.31 | 2.33 | .41 | .53 | .923 |
| DC/C | .26 | .23 | .62 | .61 | .44 | .41 | .09 | .08 | .117 |
| DC/T | .35 | .31 | 1.55 | 1.46 | .82 | .76 | .30 | .26 | .487 |
| CP/C | .04 | .09 | .71 | .34 | .23 | .18 | .13 | .07 | .058 |
| CP/T | .09 | .14 | 1.06 | .66 | .41 | .32 | .20 | .13 | .053 |
| CN/C | .71 | .56 | 1.49 | 1.11 | 1.07 | .81 | .18 | .13 | <.001* |
| CN/T | 1.40 | .90 | 3.00 | 2.16 | 1.97 | 1.47 | .42 | .33 | <.001* |

FIN = Finnish university-level English language students, HUN = Hungarian university-level English language students. The statistically significant results are followed by asterisks.

The first three indices in Table 4, mean length of clause (MLC), mean length of sentence (MLS), and mean length of T-unit (MLT), make up the first category of L2SCA indices: the length of production. According to the Mann-Whitney U test, the results of all three indices are statistically significant. MLC and MLT have the p-value of <.001 and MLS has the p-value of .002. Examining the numerical data, one can observe that Finnish students, on average, achieve higher scores in each three indices in the category. This suggests that Finnish students tend to use longer clauses, sentences, and T-units in their essays when compared to Hungarian students. Figure 1 further illustrates the differences in the mean length of production indices between the two groups and visualises the lengths of the observed syntactic structures.
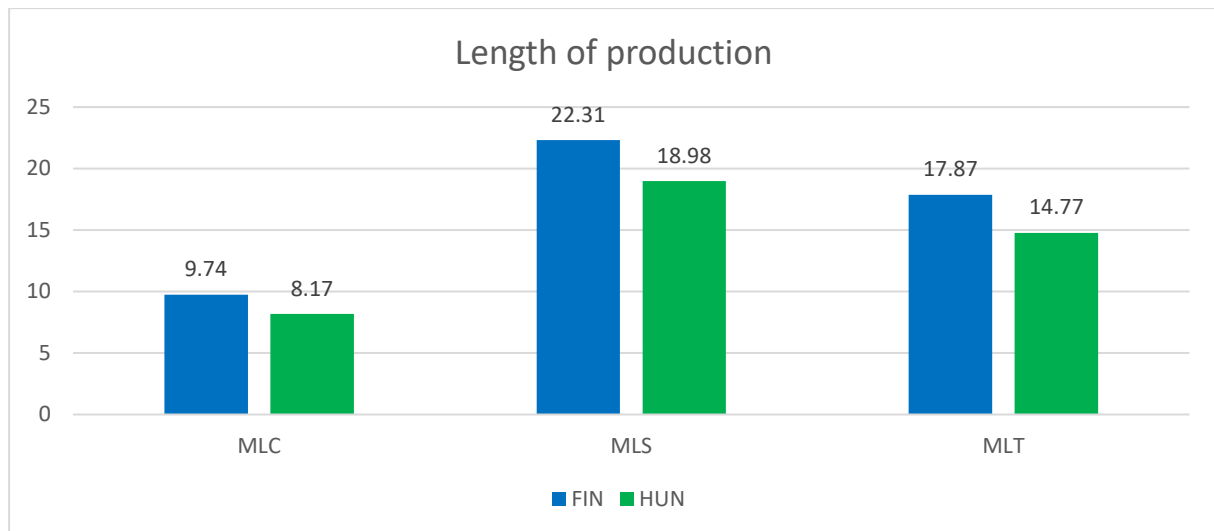
Figure 1. Length of production

In more detail, one can see in Figure 1 that, on average, Finnish students write 1.57 words longer clauses, 3.33 words longer sentences, and 3.1 words longer T-units than Hungarian students. It is important to note that while Hungarian students tend to produce more clauses, sentences, and T-units than Finnish students as is illustrated in Table 3, the produced structures tend to be shorter than those of Finnish students'. As longer length of production is considered more complex than shorter length of production, we can conclude that length of production is an aspect that separates the two groups in SC with Finnish students achieving higher SC according to the indices in this category.

The next index in Table 4, sentence complexity ratio (C/S), makes up its own category in the L2SCA indices. There is no statistically significant difference between the groups when it comes to C/S as the p-value is .923 according to the conducted Mann-Whitney U test. Furthermore, the difference between the mean C/S of Finnish students and Hungarian students is only .02, which is not enough to conclude that Hungarian students would, on average, write more clauses per sentence in a statistically significant way than Finnish students. Thus, we can conclude that the two groups have no significant difference in this category but instead have a similarity and that average sentences written by a Finnish student and a Hungarian student have the same number of clauses. However, one may note that this index is the only one in the entire analysis of the thesis, in which Hungarian students achieve a higher mean score than Finnish students.

The next two indices, dependent clauses per clause (DC/C) and dependent clauses per T-unit (DC/T), make up the category of the amount of subordination in SC indices. As Table 4

shows, Finnish students, on average, have slightly higher scores in both indices. The differences are small with .03 more dependent clauses per clause, and .06 more dependent clauses per T-unit. Furthermore, these differences are not statistically significant as the p-values are .117 and .487 according to the Mann-Whitney U test. Thus, it cannot be concluded that Finnish students use more subordination than Hungarian students. On the contrary, this is a similarity between the groups, and it can be concluded that the groups, on average, use a similar number of dependent clauses per clause and T-unit.

The results for the two indices in the category of amount of coordination are discussed next. As seen in Table 4, Finnish students, on average, have slightly higher scores in coordinate phrases per clause (CP/C) and coordinate phrases per T-unit (CP/T) indices than Hungarian students. The difference in CP/C is .05 and the difference in CP/T is .09. However, neither result is statistically significant as the CP/C index has a p-value of .058 and the CP/T index has a p-value of .053 according to the Mann-Whitney U test. Consequently, it cannot be suggested that there would be a statistically significant difference in the amount of coordination between the two groups. On the contrary, like in the amount of subordination, the groups display similarity in the category and, they, on average, use a similar number of coordinate phrases per clause and T-unit.

The final two indices which measure SC in the analysis are part of the final category of L2SCA indices: the relationship between syntactic structures and production units. As seen in Table 4, the results are statistically significant with p<.001 for both indices according to the Mann-Whitney U-test. Figure 2 illustrates the differences between the two groups.
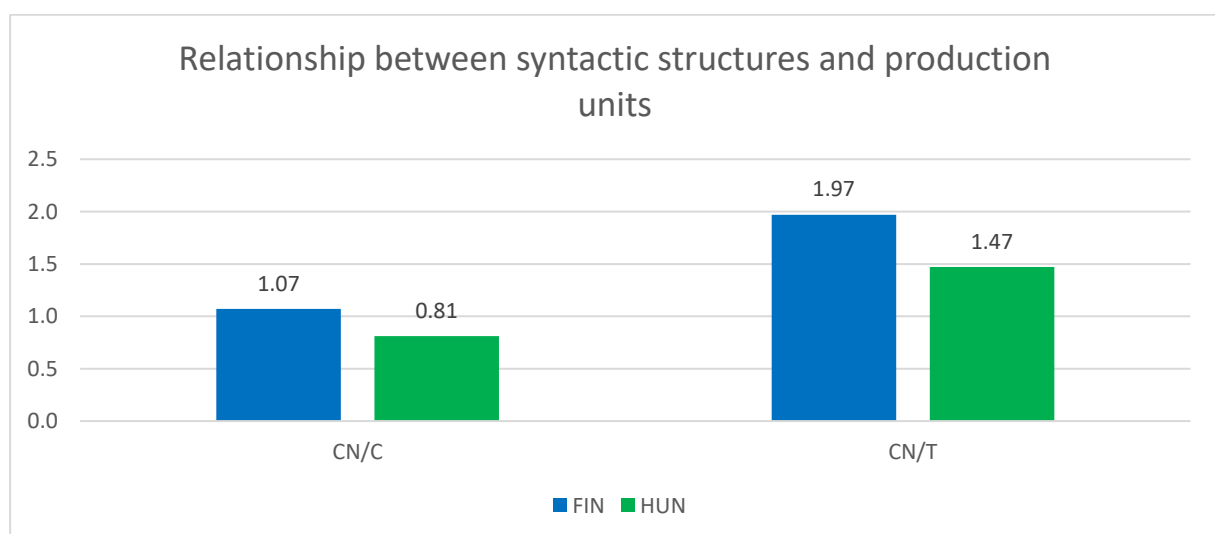


Figure 2. Relationship between syntactic structures and production units

As one can notice in Figure 2, Finnish students, on average, write .26 more complex nominals per clause (CN/C) and .5 more complex nominals per T-unit (CN/T) than Hungarian students. Thus, it can be suggested that Finnish students achieve higher SC compared to Hungarian students according to this category and that this category marks another difference between the groups.

## 4.2 Lexical complexity analysis

In this subsection, the results of the lexical complexity analysis concerning research question 1B (How do English texts written by Finnish and Hungarian university students compare in terms of lexical complexity?) will be discussed. The results of the LC analysis will be discussed considering the three dimensions of LC.

Table 5 illustrates the results for the first two dimensions of LC, lexical density and lexical sophistication, and compares Finnish university-level English language students to Hungarian university-level English language students.

Table 5. Lexical density and lexical sophistication results

| Index | Minimum | | Maximum | | Mean | | Standard Deviation | | Statistical Significance |
|---|---|---|---|---|---|---|---|---|---|
| | FIN | HUN | FIN | HUN | FIN | HUN | FIN | HUN | |
| LD | .44 | .43 | .54 | .49 | .49 | .46 | .02 | .02 | <.001* |
| LS1 | .09 | .08 | .39 | .25 | .17 | .15 | .06 | .04 | .088 |
| LS2 | .12 | .11 | .39 | .28 | .19 | .17 | .06 | .04 | .152 |
| VS1 | .00 | .01 | .51 | .19 | .09 | .06 | .09 | .04 | .026* |
| VS2 | .00 | .01 | 17.25 | 2.05 | 1.05 | .31 | 3.11 | .45 | .010* |
| CVS1 | .00 | .09 | 2.94 | 1.01 | .52 | .33 | .52 | .21 | .010* |

FIN = Finnish university-level English language students, HUN = Hungarian university-level English language students. The statistically significant results are followed by asterisks.

The first index in Table 5 concerns the first dimension of LCA indices, lexical density (LD). The p-value of the result is <.001 according to the Mann-Whitney U test, which means that the results of this dimension mark a difference between Finnish students and Hungarian students in LC. Figure 3 illustrates the differences in LD between the two groups.
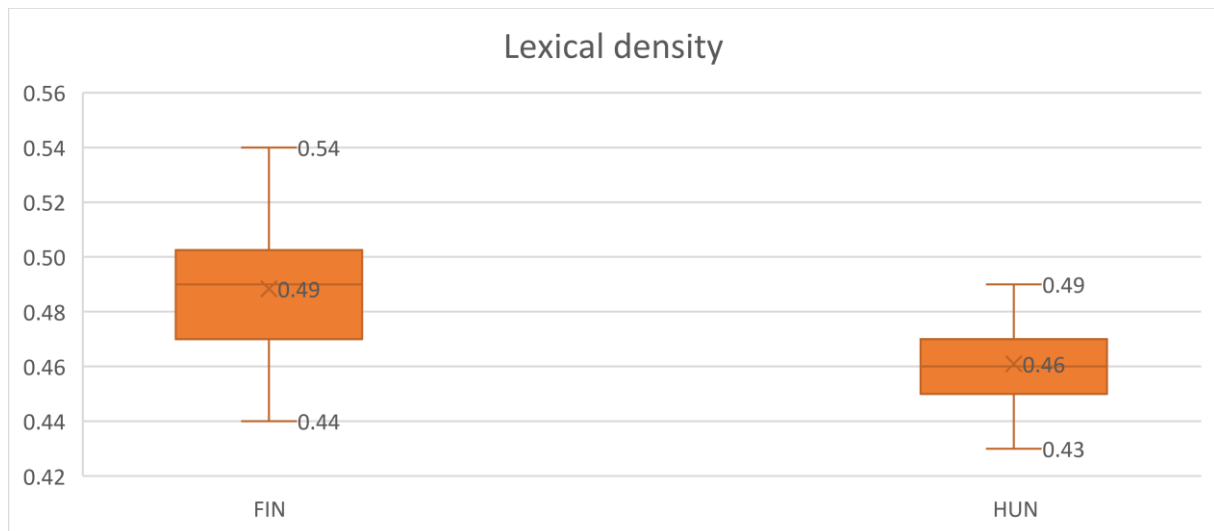
Figure 3. Lexical density

As illustrated in Figure 3, the mean value of LD is .49 for Finnish students and .46 for Hungarian students. This means that, on average, lexical words make up 49 per cent of all words in essays written by Finnish students and 46 per cent of all words in essays written by Hungarian students. Thus, for both groups, on average, just under half of all words in the essays are lexical words opposed to grammatical words. Furthermore, the results suggest that Finnish students, on average, write more lexically dense essays with a statistically significant difference between the groups. However, the difference between the two groups is relatively small at .03.

The rest of the indices in the table concern another dimension of LCA: lexical sophistication. As can be observed in Table 5, Finnish students, on average, achieve a higher mean score in all the indices in the dimension compared to Hungarian students. Lexical sophistication-I (LS1) and lexical sophistication-II (LS2) indices yield very similar results with a difference of .02 in both. On one hand, the mean LS1 is .17 for Finnish students and .15 for Hungarian students. This means that, on average, 17 and 15 per cent of the total words written by the groups are considered lexically sophisticated meaning that the words do not belong to the 2,000 most frequent words in English.

On the other hand, the mean LS2 is .19 for Finnish students and .17 for Hungarian students. Thus, on average, 19 and 17 per cent of the types written by the groups are considered sophisticated according to the index. However, the results for LS1 and LS2 are not considered statistically significant as the p-values are .088 and .152, respectively, according to the Mann-Whitney U test. Thus, a statistically significant difference between the groups cannot be

suggested. On the contrary, it can be suggested that, on average, the groups use a similar number of sophisticated lexical words or types when compared to the total word counts or total types in writing.

Verb sophistication was also measured with the indices in this category. It is very important to note the high variation in the Finnish students' results and thus, it is sensible to discuss the medians for some of the results as well. The mean for Finnish students' VS1 is .09. and the mean for Hungarian students' verb sophistication-I (VS1) is .06 , as seen in Table 5. The difference is .03. The mean corrected verb sophistication-1 (CVS1) for Finnish students is .52 and .33 for Hungarian students. Considering the results, it can be suggested that Finnish students, on average, display higher verb sophistication. The results for VS1 and CVS1 are both considered statistically significant as the p-values are .026 and .010, respectively, according to the Mann-Whitney U test.

The results for the verb sophistication-II (VS2) index are remarkable. The median for Hungarian students is .16 and the mean is .31. The median for Finnish students is .35 and the mean is 1.05. Thus, the mean for Finnish students is three times larger than the median. The significant differences in the medians and means in the Finnish students' results are likely due to the essay with the highest score of 17.25 in the index, which can be seen as an outlier in the data. The second-highest score in an essay written by a Finnish student is considerably smaller at 2.64. Thus, there is evidently a large gap between the essay with the highest score and the rest of the essays in this index. The essay in question will be further discussed in more detail in Subsection 4.3 concerning research question 2 as it is also the essay with the highest score in LC overall in the analysis.

The results for VS2 are statistically significant with a p-value of .010 according to the Mann-Whitney U test conducted. Thus, examining the results of the three indices, VS1, CVS1, and VS2, which measure verb sophistication, it can be suggested that there is a statistically significant difference between the groups in that Finnish students, on average, tend to use more sophisticated verbs when compared to Hungarian students. Thus, Finnish students display higher LC in this regard.

The final and largest dimension of LCA indices used in the LC subanalysis, lexical variation, consists of 19 indices. The results for the first ten of the 19 indices in the dimension are displayed in Table 6.

Table 6. Lexical variation results: NDW, TTR, and their transformations

| Index | Minimum | | Maximum | | Mean | | Standard Deviation | | Statistical Significance |
|---|---|---|---|---|---|---|---|---|---|
| | FIN | HUN | FIN | HUN | FIN | HUN | FIN | HUN | |
| NDW | 169 | 165.0 | 295 | 329 | 219.20 | 204.20 | 33.33 | 29.94 | .031* |
| NDWZ-50 | 33 | 27 | 41 | 41 | 37.30 | 35.70 | 2.32 | 3.19 | .046* |
| NDW-ER50 | 34.10 | 34.60 | 41.40 | 40.70 | 38.31 | 37.73 | 1.56 | 1.34 | .071 |
| NDW-ES50 | 33.70 | 32.50 | 40.80 | 39.10 | 37.75 | 36.92 | 1.76 | 1.44 | .049* |
| TTR | .32 | .33 | .50 | .44 | .39 | .37 | .04 | .03 | .008* |
| MSTTR | .64 | .68 | .81 | .78 | .75 | .74 | .04 | .02 | .031* |
| CTTR | 5.18 | 5.29 | 8.58 | 7.69 | 6.56 | 6.11 | .80 | .6 | .015* |
| RTTR | 7.32 | 7.48 | 12.13 | 10.87 | 9.27 | 8.64 | 1.13 | .82 | .015* |
| LOGTTR | .82 | .82 | .89 | .87 | .85 | .84 | .02 | .01 | .018* |
| UBER | 14.89 | 15.22 | 25.46 | 20.83 | 18.71 | 17.32 | 2.36 | 1.51 | .010* |

FIN = Finnish university-level English language students, HUN = Hungarian university-level English language students. The statistically significant results are followed by asterisks.

The results in Table 6 concern NDW, TTR, and their transformations. Thus, these results tell us of the overall lexical variation in the essays. On average, as seen in Table 6, Finnish students achieve a higher mean score with each index when compared to Hungarian students. The first index, the number of different words (NDW) shows that Finnish students, on average, write 15 more different types compared to Hungarian students in their essays. The result is statistically significant according to the Mann-Whitney U test with a p-value of .031. However, this result does not necessarily tell us much and it is not as important as the rest of the results in the table as it does not take the total word counts of the essays into account. That said, it gives us a hint as to what the TTRs might be as we know the mean total word counts between the groups are very similar with Finnish students, on average, writing 557 words and Hungarian students writing 555 words.

The indices which measure the NDW of the first 50 words (NDWZ-50), the NDW of the expected random 50 words (NDW-ER50), and the NDW of the expected sequence of 50 words (NDW-ES50), or the transformations of NDW, bring us much more insight. The indices' means for Finnish students are 37.30, 38.31, and 37.75. For Hungarian students, the means are 35.70, 37.73, and 36.92. The differences between the groups are then, on average, 1.6, .58, and .83 types. The results for NDW-ER50 are not statistically significant according to the Mann-Whitney U test with a p-value of .071. However, the p-value for the results of

NDWZ-50 is .046 and for the results of NDW-ES50, it is .049. Thus, the results for both indices are statistically significant. This means that according to the indices' results, Finnish students, on average, write 1.6 or .58 types more than Hungarian students. Therefore, it can be observed that Finnish students tend to use a larger number of different words in their essays when compared to Hungarian students.

All the results for the mean type-token ratio (TTR) and its transformations are statistically significant according to the Mann-Whitney U test. Figure 4 illustrates differences in TTR between the two groups.
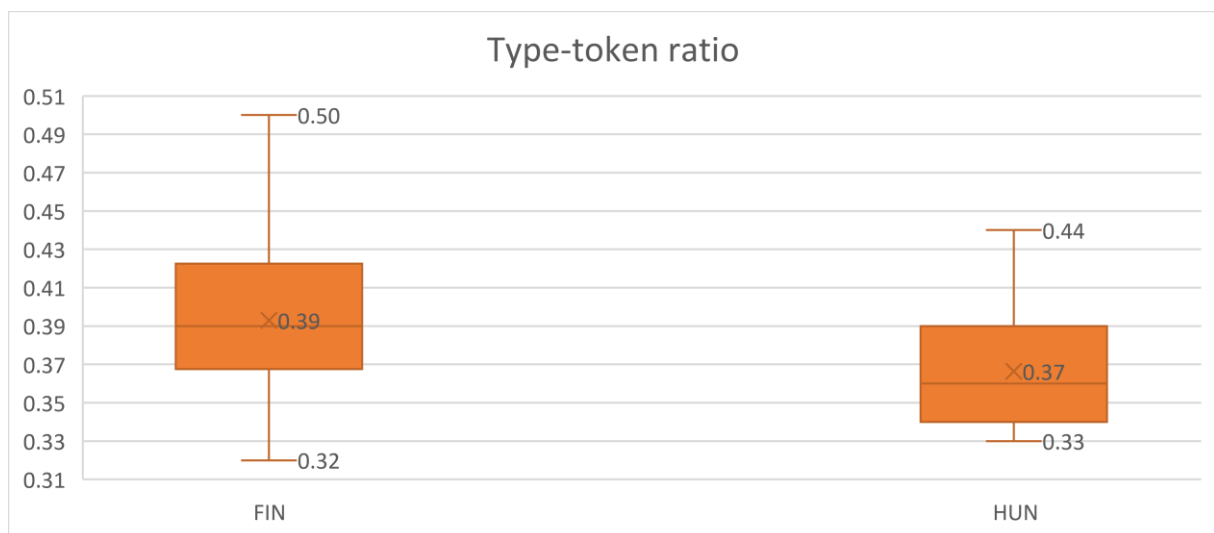


Figure 4. Type-token ratio

As seen in Figure 4, the mean TTR for Finnish students is .39 and .37 for Hungarian students. The difference between the groups is small at .02. However, the p-value is .008 making the result statistically significant. Thus, Finnish students, on average, write .02 types per token more in their essays than Hungarian students.

The scores of mean segmental TTR (MSTTR), corrected TTR (CTTR), root TTR (RTTR), bilogarithmic TTR (logTTR), and Uber Index (Uber) for Finnish students are .75, 6.56, 9.27, .85, and 18.71. The scores for Hungarian students are .74, 6.11, 8.64, .84, and 17.32. The differences between the groups are .45, .63, .01, and 1.39. Results for all four are statistically significant. MSTTR has a p-value of .031, CTTR and RTTR has a p-value of .015, logTTR has a p-value of .018, and Uber has a p-value of .010. These results further illustrate how Finnish students use more lexically varied language in their essays when compared to Hungarian students. Furthermore, Finnish students display higher LC in this regard.

Table 7 illustrates the rest of the results for the lexical variation dimension of LCA indices.

Table 7. Lexical variation results: variation of lexical elements

| Index | Minimum | | Maximum | | Mean | | Standard Deviation | | Statistical Significance |
|---|---|---|---|---|---|---|---|---|---|
| | FIN | HUN | FIN | HUN | FIN | HUN | FIN | HUN | |
| VV1 | .42 | .45 | .93 | .73 | .63 | .59 | .11 | .07 | .065 |
| SVV1 | 13.24 | 13.00 | 57.37 | 39.18 | 28.40 | 23.09 | 9.55 | 5.94 | .017* |
| CVV1 | 2.57 | 2.55 | 5.36 | 4.43 | 3.72 | 3.37 | .63 | .43 | .018* |
| VV2 | .11 | .11 | .22 | .19 | .16 | .15 | .03 | .02 | .126 |
| LV | .50 | .52 | .80 | .72 | .63 | .61 | .07 | .05 | .110 |
| NV | .53 | .50 | .85 | .76 | .64 | .62 | .07 | .07 | .264 |
| ADJV | .08 | .09 | .18 | .17 | .14 | .12 | .02 | .02 | .002* |
| ADVV | .06 | .06 | .16 | .11 | .10 | .09 | .02 | .02 | .002* |
| MODV | .17 | .16 | .30 | .28 | .24 | .21 | .03 | .03 | <.001* |

FIN = Finnish university-level English language students, HUN = Hungarian university-level English language students. The statistically significant results are followed by asterisks.

As displayed in Table 7, the mean verb variation-I (VV1) for Finnish students is .63 and .59 for Hungarian students. The difference between the groups is .04 and the p-value is .065 and thus, the results are not statistically significant according to the Mann-Whitney U test. The results for squared VV1 (SVV1) and corrected VV1 (CVV1) are statistically significant with p-values of .017 and .018. The mean SVV1 and CVV1 are 28.40 and 2.72 for Finnish students and 23.09 and 3.37 for Hungarian students. The differences are 5.31 and .35. The means for verb variation-II (VV2) are similar between the groups with .16 for Finnish students and .15 for Hungarian students. The p-value is .126. The results thus indicate that Finnish students, on average, may use a larger variety of verbs in their writing than Hungarian students, but the result is not necessarily clear.

Looking at lexical word variation (LV) and noun variation (NV), there is not a statistically significant difference between the groups with the mean values being .63 and .64 for Finnish students and .61 and .62 for Hungarian students. Additionally, according to the Mann-Whitney U test, the p-values exceed .05 as they are .110 and .264, respectively. However, it can be suggested that Finnish and Hungarian students display similarities in variation considering lexical words and nouns.

The mean adjective variation (ADJV) for Finnish students is .14 and .12 for Hungarian students. The difference is small at only .02. However, the Mann-Whitney U test assesses the

result to be statistically significant as the p-value is .002. This means that, on average, Finnish students have more variation in their use of adjectives than Hungarian students. The results for adverb variation further demonstrate a difference between Finnish students and Hungarian students. ADVV is .10 for Finnish students and .09 for Hungarian students. The p-value is .002. This again means that, on average, Finnish students have more variation in their use of adverbs than Hungarian students.

The final index of the dimension and the LC subanalysis, modifier variation (MODV), shows that, on average, adjective and adverb types make up 24 per cent of the lexical words written by Finnish students and 21 per cent of lexical words written by Hungarian students. The p-value is <.001 making the results statistically significant according to the Mann-Whitney U test. Observing the results of this dimension's indices, it can be suggested that the language written by Finnish students is, on average, more lexically varied than the language written by Hungarian students.

## 4.3 Syntactic structures and lexical elements in observed high-complexity and low-complexity essays

This subsection aims to answer research question 2 (How do syntactic structures and lexical elements manifest in observed high-complexity and low-complexity essays?). First, in subsection 4.3.1, the highest and lowest scores in SC will be examined. After, in subsection 4.3.2, the highest and lowest scores in LC will be discussed. All the examples discussed are from the essays from the ACALEX corpus.

### 4.3.1 Manifestation of syntactic structures

The highest overall score in SC according to the statistical analysis is in an essay written by a Finnish student, ACALEX-1102. They receive higher than average scores in five out of the ten indices measured when compared to other Finnish students. These five measures are notably also the five indices in which a statistically significant difference was found in the comparison between Finnish and Hungarian university-level English language students: MLC, MLS, MLT, CN/C, and CN/T. Furthermore, the essay's scores are the highest out of all essays in MLC and CN/C. In the remaining five indices, the essay's scores are average when compared to other Finnish students.

The mean length of a sentence (MLS) in the essay is 27.12 words, the mean length of a T-unit (MLT) is 23.38, and the mean length of a clause (MLC) is 12.33. On average, the essay contains 2.83 complex nominals per T-unit (CN/T) and 1.49 complex nominals per clause (CN/C). In the remaining five indices, the essay's scores are average when compared to other Finnish students. Thus, it can be concluded that the strength of the essay in SC manifests in the categories of the length of production and the relationship between syntactic structures and production units. Example 1 from ACALEX-1102's essay displays high SC in multiple ways.

> (1) "After continuously building upon my interest and skills, I went to high school, where I finally became confident enough in my skills to devote myself to making a future out of something that once was just an intriguing school subject".

The length of production in the essay is notably longer than that of an average essay written by a Finnish student. Example 1 is one of the longest sentences found in the essay and consists of forty words. There is also a long dependent relative clause of 27 words starting with "where". Even if the essay's scores are average in the indices measuring the amount of subordination it is important to note how example 1 also features a significant amount of subordination as there are multiple clauses dependent on the main clause "I went to high school". One should note that while it is evident that the sentence is long, it also contains a lot of information, and the use of different syntactic structures works in harmony with the semantic meaning. Example 2 is a much shorter sentence compared to the first example but shows high SC differently.

> (2) "Being able to cope with change is key, and that is what I will strive for".

ACALEX-1102' essay also shows frequent use of complex nominals. Example 2 features two complex nominals in one sentence. "Being able to cope with change" is a complex nominal functioning as a subject and "what I will strive for" is another complex nominal functioning as a subject predicative. Upon analysis, it is evident that the sentence showcases sophisticated use of language where it is not always the length of production implicating complexity but the intelligent usage of syntactic structures. The two example sentences reflect the overall SC of the essay in the way that even if the essay is also remarkable in terms of length of production, the essay displays SC in its full multiplexity using different types of syntactic structures.

The essay of the Hungarian participant ACALEX-1211 has consistently the lowest score in SC according to the statistical analysis. Compared to other Hungarian students, the essay's scores are lower than average in each of the indices with the lowest scores in four indices: MLT, DC/C, DC/T, and CN/T. However, for this analysis, let us focus on the indices which were found to be statistically significant in explaining the differences between Finnish and Hungarian university-level English language students and which were also found to be the strengths of the syntactically complex essay: MLC, MLS, MLT, CN/C, and CN/T. MLS in the essay is 14.67 words, which is significantly shorter compared to the high SC essay by ACALEX-1102 with 27.12 words. MLT in the essay is 10.10 and MLC is 7.51. The essay contains few complex nominals with its CN/T being .90 and CN/C being .67. Example 3 displays an especially short sentence from the essay.

(3) "I think my teacher was the best".

One notable feature in ACALEX-1211's essay is the short length of production. The sentence in example 3 consists of only six words. However, it does feature two clauses: "I think" as the main clause and "my teacher was the best" as a dependent clause. As the whole sentence is short, the clauses and T-units are short as well. Comparing this sentence with examples 1 and 2, which show high SC, is highly intriguing. All three are sentences with multiple clauses yet it is evident that examples 1 and 2 carry more meaning than example 3. This might be due to a significant difference in the number of sentences between the two essays even with relatively similar word counts. The essay with lower SC consists of 42 sentences and 616 words. On the contrary, the essay with higher SC, although longer in total words with 678 words, consists of only 25 sentences. Example 4 is a longer sentence from ACALEX-1211's essay with longer clauses.

(4) "Teachers have a lot of tasks though, but they have more free time in the summer, and they can spend more time with their families".

Another notable feature in the essay is the low amount of subordination. Example 4 displays a sentence with three main clauses separated by coordinators "but" and "and" without any dependent clauses. Without subordination, the sentence is somewhat crowded and the relationships between the clauses are unclear. Furthermore, by adding subordination and changing the order of the clauses the sentence would be more coherent: "Teachers have a lot of tasks though, but they can spend more time with their families because they have more free time in the summer". The two examples demonstrate how the lack of SC affects the essay.

The words and the information may be there, but the simplicity of the syntactic structures used does not add to the semantics of the essay and the reader may be left wanting more.

### 4.3.2  Manifestation of lexical elements

Next, let us discuss the essays with the highest and lowest LC. The essay with the highest LC according to the statistical analysis is undoubtedly the essay from participant ACALEX-1127. The essay's scores are higher than average in twenty out of the 25 indices, 16 of which its scores are the highest when compared to other Finnish students. Its scores are average in four indices and lower than average in one index. In terms of lexical diversity (LD), the essay achieves an average score of .50 when compared to other Finnish students. However, in terms of lexical sophistication, the essay achieves the highest score in each index. The essay's LS1 and LS2 are both .39, which means that 39 per cent of all words in the essay are sophisticated lexical words. LS1 and LS2 were not found to be statistically significant in the comparison between Finnish and Hungarian university-level English language students but verb sophistication was. The essay's VS1 is .51, which means that 51 per cent of all verb types are sophisticated verb types. Examining lexical variation, one notices that the essay also has the highest number of different words (NDW) with 295 words and the highest TTR with .50. The essay is extraordinary in LC to the point that it affects its readability. One might even wonder whether the writer wrote the entire essay specifically aiming for high LC. Example 5 displays a relatively long sentence with impressive use of vocabulary.

> (5) "A likelihood of never practicing an academic vocation in an organized society is conceivable and may be superseded by a scenario of fighting and struggling for my life every single day or being confined in a dark room".

The example sentence effectively demonstrates how the writer does little repetition when it comes to lexical words. The only repeated words in the sentence are "a", "an", "and" "in", and "of", which are all grammatical words. One can argue that grammatical words have to be repeated in longer texts because of the nature of them whereas lexical words can be usually used more freely. Since the writer avoids repetition of words, one may notice how they instead in one case might write something with two different words but with similar meanings. For example, in the sentence "fighting" and "struggling" convey relatively the same meaning in the context of the sentence and the overall essay. Thus, the meaning is repeated but not the actual words. The size of the vocabulary of the writer is impressive. However, the large number of unique words in one sentence makes the essay more difficult to

read and might distract the reader from the topic of the essay: the professional future of the student. Example 6 further exemplifies high LC by ACALEX-1127.

> (6) "Most of the people consider studies and erudition at university level to be the pinnacle stage of the education system endorsed with the most valiant grace and the power of influencing besetting society".

The writer not only uses a large number of different words, but they also tend to use sophisticated lexical words. The sentence in example 6 features six words that do not belong to the 2,000 most frequent words in English: "erudition", "pinnacle", "endorse", "valiant", "grace", and "besetting" (OED n.d.).

The essay with the lowest LC according to the statistical analysis is an essay by a Hungarian student ACALEX-1197. The essay's scores are lower than average in 15 indices, six of which in it has the lowest scores, and its scores are average in ten indices when compared to other Hungarian students. In terms of lexical diversity and lexical sophistication, the essay's scores are average when compared to other Hungarian students. The essay has the lowest NDW with 165 words, but importantly the essay is also the shortest in the data with 477 words. However, the essay's scores are also the lowest in the three transformations of NDW. This means that the writer repeats words often. The TTR of the essay is .34. Example 7 illustrates how certain words are repeated throughout the sentences.

> (7) "She was my favourite teacher because students always paid attention to her and her lessons. She always made her lessons interesting and exciting. I prepared conscientiously for her lessons. She got on well with her students and she was very helpful".

In Example 7, seven words are repeated throughout the sentences. Three of the repeated words are grammatical words: "she", "her", and "and". This type of repetition can be compared to the type of repetition the high LC essay has as well: repetition of words that are difficult to replace with other words because of their important grammatical functions. Four of the repeated words are lexical words: "was", "students", always, and "lessons". The repeated lexical and grammatical words make up around half of the whole text sample. Additionally, one can notice that most of the vocabulary is normal in terms of frequency. For example, the verbs used are all frequent words in English (OED n.d.). Surprisingly, the example includes one sophisticated word with a lower frequency in English:

"conscientiously" (OED n.d.). The following example illustrates beautifully how LC can be consistent but still multiplex.

(8) "In addition, if we want a good job we will have to be able to speak foreign languages and it is good if we can prove it with a certificate".

Example 8 further shows how the essay features a lot of repetition. In the example sentence the adjective "good" is used twice. The use of "good" leaves the intended meaning vague and simple and ultimately leaves the reader to wonder. For example, in the first case with "a good job", what is "a good job" according to the writer? Is it, for example, high-paying, fun, or humanitarian? Perhaps it is all three or something else entirely.

# 5 Discussion

In this section, the results of the analysis will be discussed further linking the findings to earlier research. However, first, the results will be discussed concerning the initial hypothesis made in Section 3 of the thesis. Furthermore, possible reasons for the results of the analysis will be hypothesised and any particular findings that ought to have further discussion will be highlighted. Finally, the generalisability and the limitations of the present study will be discussed.

This thesis aimed to compare syntactic and lexical complexity in essays written by Finnish and Hungarian university-level English language students and to highlight how syntactic structures and lexical elements manifest in the observed high-complexity and low-complexity essays. The initial hypothesis was that the two groups would differ, but no further speculations were made about what or how significant those differences would be as there was not sufficient prior research on the topic. Examining the results, it is apparent that there is a difference between Finnish and Hungarian university-level English language students; Finnish students, on average, achieve higher SC and LC compared to Hungarian students. Thus, the results support the initial hypothesis made.

According to the results, the difference between Finnish and Hungarian students is relatively clear but not necessarily with every index used in the analysis. In SC, Finnish students, on average, achieve higher scores with a statistical significance in half, five of the ten, indices: MLC, MLS, MLT, CN/C, and CN/T. These indices belong to the categories of the length of production and the relationship between syntactic structures and production units. No statistically significant difference was found in the rest of the categories: the sentence complexity ratio, the amount of subordination, and the amount of coordination.

Importantly, Hungarian university-level English language students achieve a higher mean score only in one index used in the whole analysis of the thesis, C/S, of which the result is not statistically significant according to the Mann-Whitney U test. However, perhaps Finnish and Hungarian university-level English language students have similarities in SC as well because half or five of the ten indices show no statistically significant difference. This means that the two groups show similarities in subordination, coordination, and in ratio of clauses and sentences. Thus, the differences lie mostly in the length of production, not necessarily in the choice of structures. To conclude, as Finnish students, on average, achieve higher scores in

half of the SC indices and similar scores to Hungarian students in the rest, it is safe to say that Finnish students, on average, display higher complexity in this regard.

In LC, the difference between Finnish and Hungarian students is much clearer than in SC. Finnish students, on average, achieve higher scores with a statistical significance in 18 of the 25 indices, which is 72 per cent of the indices used in the subanalysis. Hungarian students, on average, do not score higher in any of the LC indices. Thus, the difference between the two groups is considerably clearer in LC than in SC. Crossley and McNamara (2012, 120) state that LC may be one of the strongest predictors in detecting the L1 of a writer, which makes sense when observing the results of the present study as the difference appears to be stronger than in SC.

Even if the majority of the indices indicate a clear difference between the two groups, some indices indicate similarity. Two lexical sophistication indices, LS1 and LS2, show no statistically significant difference between the groups. Thus, the groups have a similar ratio of lexical words when compared to the whole word counts. Five of the indices with no statistically significant difference relate to lexical variation. One transformation of NDW, NDW-ER50, has a similar result with no statistically significant difference. However, as the other indices relating to the number of different words yield different results it is difficult to make conclusions about whether this similarity is important. The rest of the indices with no statistically significant difference are VV1, VV2, LV, and NV. The first two indices concern verb variation and the other two concern lexical words and nouns, respectively. It is important to note that the two other verb variation indices, SVV1 and CVV1, do show statistical differences and thus, it is not easy to assess the significance of the similarity in verb variation. However, concerning LV and NV, we can say that Finnish and Hungarian students, on average, use similar ratios of lexical types and nouns to the total number of lexical words. Furthermore, overall, there are more differences than similarities between Finnish and Hungarian university-level English language students when it comes to LC and it is safe to state that Finnish students display higher complexity in this regard.

The result that Finnish university-level English language students achieve higher scores than Hungarian students in both SC and LC is not surprising in light of limited earlier observations (e.g., European Commission, 2012 and EF, 2022) which suggest a clear difference between Finnish and Hungarian speakers of English. Crossley and McNamara (2012, 121) suggest that Finnish learners of English might score higher in lexical sophistication than in SC. Examining

the results of the present study, this may be possible as the results of the LC analysis show the difference between Finnish and Hungarian students significantly clearer. Additionally, the results of the present study align with earlier research suggesting that L1 may affect L2 complexity (e.g., Bernardini and Granfeldt, 2019, Lu and Ai, 2015, Phuoc and Barrot, 2022, and Ehret and Szmrecsanyi, 2019).

It is similarly not surprising that as the results indicate a difference between the two groups, the results of the two subanalyses correlate in that one group has higher scores in both SC and LC as it is likely that the two complexities develop alike. However, in the present study, it was also found that SC and LC do not always correlate perfectly. For example, the essay with the highest score in LC has an average score in SC. Thus, different types of complexity might manifest and develop differently. This speaks to the interplay of different types of complexities and how the present study provides a more holistic analysis of complexity by including two different types of complexities in the same analysis.

The extent of the difference between the groups is somewhat surprising as the groups compared are on the same academic level and have similarities overall. Furthermore, why the difference between the two groups is so significant, is compelling to speculate. As complexity is part of CAF, a framework used for L2 proficiency and L2 development, it is rational to discuss whether overall proficiency in English could explain the results. However, it is imperative to remember that even learners of similar proficiency levels may have differences in the development of the dimensions of CAF because of their differing L1 backgrounds (Barrot and Gabinete 2021, 227).

While proofreading the essays of the groups, it became apparent that there is a difference in overall proficiency with Finnish university-level English language students showing signs of higher proficiency. Finnish students appear to have higher accuracy as their essays required fewer edits than Hungarian students to prepare for the automatic analyses. There were other aspects which also suggested that Finnish students have a higher proficiency in writing in English than Hungarian students. For example, in the essays analysed, Finnish students also more often use the British or American conventions of quoting whereas Hungarian students more often use the Hungarian conventions of quoting. The overall proficiency of the two groups in the present study was not assessed, which should be acknowledged as the primary limitation of the thesis. However, as all the participants were university-level English students, one could assume the proficiency to be at least somewhat comparable.

The results in the comparison between groups on the same academic level raise important speculations on why the differences in overall proficiency are so significant. What are the differences between Finnish and Hungarian students that explain the difference? At first glance, the countries, Finland and Hungary, have many similarities: both are situated in Europe, are part of the European Union, and have a similar approach to foreign language instruction. This thesis did not research the differences between foreign language instruction, but there may be some aspects in the difference of instruction which may explain the results of the thesis. Furthermore, the reason for the differences could lie in socio-cultural factors, availability of resources or history.

Examining the specific SC measures, which show higher levels of linguistic complexity development, reveals that the results of the present study correspond to Bulté and Housen (2014, 53), who suggest that SC development manifests in the length of production, for example, in the length of clause, sentence, and T-unit. Kyle and Crossley (2018, 334) also observe that longer T-units are more common in higher proficiency L2 texts. Bulté and Housen (2014, 53) also note that SC development does not manifest in subordination. This is aligned with the results of the present study but differs from Ellis and Barkhuizen (2005, 139–140), who suggest a strong correlation between subordination and SC. Additionally, Bulté and Housen (2014, 54–55) suggest that there may be correlations between complexity measures and perceived writing quality. This does support the results of the present study, as the perceived writing quality observed during proofreading and the results of the analysis correlate. The results differ from Bulté and Housen (2014, 53) in that coordination was not found to explain differences in SC in the present study. Concerning the specific LC measures, the results of the present study are partially in correspondence with Kim, Crossley, and Kyle (2018, 135), who suggest that writing quality can be predicted by examining lexical sophistication. In the present study, in the category of lexical sophistication, only verb sophistication was found to be different in a statistically significant way.

It is intriguing to examine the results considering the Trade-off Hypothesis (Skehan 2009, 511), which suggests that learners may prioritise different CAF elements over the other elements. Whilst this is not the focus of the present study, it raises discussion, especially concerning the high LC essay examined concerning research question 2. Examining the essay, it becomes apparent that the writer prioritises LC. A question of whether this affected the other CAF dimensions, accuracy and fluency, remains unanswered, although based on the

number of edits which had to be made during proofreading, the essay seems to display high accuracy meaning that there may not have been an effect on it.

However, one may observe that the same essay's scores are average in terms of SC. The idea that prioritisation of LC affects the essay's SC negatively is very intriguing. Could there be a similar effect described by Skehan (ibid.) as in the Trade-off Hypothesis with prioritisation of specific types of complexities? Additionally, one should consider whether prioritisation of LC to the same extent is typical for the author also in their L1. It appears that the writer used dictionaries and/or thesauri, which affects the LC of their essay. If this is the case, this particular case may be aligned with Rahayu, Utomo, and Setyowati (2021, 259), who suggest that the use of such tools, which are more common in L2 writing than L1 writing, may result in the writer's LC being higher in their L2 compared to their L1. Furthermore, this raises a question about different factors influencing complexity. The essays analysed in the present study may display higher or lower complexity because the essays were written with the writers understanding that they were university course assignments which were to be archived for research purposes.

This thesis has highlighted written learner language and the multiplex nature of L2 complexity. Furthermore, this thesis has explored the cross-linguistic influence on linguistic complexity. Thus, the results of the study have helped to advance the understanding of the phenomena and have contributed to the field of SLA. The present study may be useful for people working with language learners. For example, in light of the results of the present study, professionals in the teaching field should understand how L1 background may affect the written language of a student and as students are not always a homogeneous group, instruction and language teaching curricula should be adjusted accordingly in learning environments with students of different backgrounds.

This study focuses on a relatively small sample of Finnish and Hungarian university-level English language students and is not meant to describe Finnish and Hungarian speakers of English as a whole but to highlight the possible differences and similarities between the two studied groups. Future studies may benefit from having a larger sample size or examining different types of texts, for example, spoken language. Additionally, the present study used mostly large-grained measures to investigate linguistic complexity and, thus, focusing on fine-grained measures may yield different results. The lack of assessment of the overall proficiency of the groups is a central limitation of the study, which any future study should

rectify, as this produced an additional variable in the study. However, the observed differences and similarities are evidence of the multiplexity of the interplay of written learner language, L1 background, and linguistic complexity in L2. Therefore, the results of the study are generalisable in that these phenomena can affect each other in various ways.

# 6  Conclusion

In this section, the thesis will be concluded by summarising the theoretical background, the methodology used, and the results of the analysis. Finally, after concluding the present study, suggestions for future research will be discussed.

Linguistic complexity in written learner language is a highly multiplex topic which has been researched with numerous viewpoints and methodologies. Of the two types of complexities researched in the present study, SC is much more studied and more difficult to define. SC involves examining the syntactic structures which the learner uses in their writing and assessing complexity by using measures such as length of production and amount of subordination or coordination. LC, in contrast, involves examining the lexical elements in learner language and assessing complexity by using measures such as lexical sophistication and variation. Earlier research involving a holistic and descriptive viewpoint on linguistic complexity in written learner language is scarce, but there is evidence of a multiplex relationship between the two phenomena.

The primary data in this thesis was the ACALEX corpus from which sixty essays written by Finnish and Hungarian university-level English language students were analysed. This pairing for the comparison was especially remarkable as Finnish and Hungarian, which are the participants' L1s, are related but only distantly. Furthermore, Finland and Hungary have some similarities, for example, with their approaches to foreign language instruction and EFL in general. The analysis consisted of using automatic analysers L2SCA and LCA by Xiaofei Lu to calculate syntactic and lexical complexity in the essays, further statistical analyses made with Microsoft Excel and IBM SPSS Statistics 27, and a more detailed qualitative analysis of four extreme case essays.

In the analysis, regarding research questions 1A and 1B,  it was found that Finnish university-level English language students, on average, display higher SC and LC compared to Hungarian university-level English language students. The difference in SC is explained by differences in the length of production (MLC, MLS, and MLT) and the relationship between syntactic structures and production units (CN/C and CN/T). Thus, Finnish students, on average, produce longer clauses, sentences, and T-units and produce more complex nominals per clause or T-unit than Hungarian students. On the contrary, similarities were found in the amount of subordination (DC/C and DC/T), the amount of coordination (CP/C and CP/T) and

in sentence complexity ratio (C/S). Thus, Finnish students, on average, score higher in half of the indices and the groups score comparatively in the other half.

The differences between the groups in LC are clearer. The difference between the groups is explained by the differences in lexical density (LD), verb sophistication (VS1, VS2, and CVS1), number of different words (NDW, NDWZ-50, and NDW-ES50), type-token ratios (TTR, MSTTR, CTTR, RTTR, LOGTTR, and UBER), verb variation (SVV1 and CVV1), adjective variation (ADJV), adverb variation (ADVV), and modifier variation (MODV). Thus, Finnish students, on average, tend to use more varied vocabulary and more sophisticated vocabulary than Hungarian students although some similarities were also found.

Regarding research question 2, it was found that syntactic structures and lexical elements manifest in various ways in observed high-complexity and low-complexity essays. The syntactic structures tend to be long in the high SC essay. Additionally, the essay features many complex nominals and intelligent use of grammar. On the contrary, the syntactic structures in the low SC essay tend to be short. Furthermore, there is a lack of subordination, which affects the relationship between the clauses. The lexical elements in the high LC essay are varied and sophisticated. Furthermore, repetition of lexical words in the essay is rare. The low LC essay features a lot of repetition of words and tends to have a more simplistic vocabulary.

The present study presents many opportunities for future research. As cross-linguistic influence on L2 linguistic complexity has not yet been thoroughly researched, many more comparisons between different L1 speakers should be made. This includes comparing L1 speakers of languages both linguistically related and unrelated to each other. It may also be interesting to compare the same writer in their L1 and L2. Additionally, future studies should highlight how the three dimensions of CAF interact with each other and compare Finnish and Hungarian learners of English of the same overall assessed proficiency level. Another future study opportunity is to research differences and similarities between the two groups in spoken language as it might bring a much more multiplex point of view to L2 complexity. One may also find it interesting to not focus solely on linguistic complexity but highlight the differences in overall proficiency between Finnish and Hungarian language learners and investigate the reasons behind the observed differences.

# References

Abdi Tabari, Mahmoud, Xiaofei Lu, and Yizhou Wang. 2023. "The Effects of Task Complexity on Lexical Complexity in L2 Writing: An Exploratory Study." *System* 114: 103021. Accessed 17 April 2023. Elsevier.

Abondolo, Daniel, and Riitta-Liisa Valijärvi. 2023. "Introduction to the Uralic Languages, with Special Reference to Finnish and Hungarian." In *The Uralic Languages*, edited by Daniel Abondolo, and Riitta-Liisa Valijärvi, 1–80. Abingdon: Routledge. Accessed 9 May 2023. T&F eBooks.

Ai, Haiyang, and Xiaofei Lu. 2010. "A web-based system for automatic measurement of lexical complexity." Paper presented at the 27th Annual Symposium of the Computer-Assisted Language Consortium (CALICO-10). Amherst, MA. June 8–12.

———. 2013. "A Corpus-Based Comparison of Syntactic Complexity in NNS and NS University Students' Writing." In *Automatic Treatment and Analysis of Learner Corpus Data*, edited by Ana Díaz-Negrillo, Nicolas Ballier, and Paul Thompson, 59: 249–264. John Benjamins: Amsterdam. Accessed 15 September 2023. ProQuest.

Alexopoulou, Theodora, Marije Michel, Akira Murakami, and Detmar Meurers. 2017. "Task Effects on Linguistic Complexity and Accuracy: A Large-Scale Learner Corpus Analysis Employing Natural Language Processing Techniques: Task Effects in a Large-Scale Learner Corpus." *Language Learning* 67, no. S1: 180–208. Accessed 17 September 2023. EBSCO.

Allaw, Elissa. 2021. "A Learner Corpus Analysis: Effects of Task Complexity, Task Type, and L1 & L2 Similarity on Propositional and Linguistic Complexity." *International Review of Applied Linguistics in Language Teaching, IRAL* 59, no. 4: 569–604. Accessed 3 May 2023. EBSCO.

Barrot, Jessie, and Mari Karen Gabinete. 2021. "Complexity, Accuracy, and Fluency in the Argumentative Writing of ESL and EFL Learners." *International Review of Applied Linguistics in Language Teaching, IRAL* 59, no. 2: 209–232. Accessed 20 March 2023. EBSCO.

Bernardini, Petra, and Jonas Granfeldt. 2019. "On Cross-linguistic Variation and Measures of Linguistic Complexity in Learner Texts: Italian, French and English." *International Journal of Applied Linguistics* 29, no. 2: 211–232. Accessed 20 March 2023. EBSCO.

Biber, Douglas, Susan Conrad, and Geoffrey N. Leech. [2002] 2019. *Longman Student Grammar of Spoken and Written English*. New York: Longman.

Bulté, Bram, and Alex Housen. 2012. "Defining and operationalising L2 complexity." In *Dimensions of L2 Performance and Proficiency: Investigating Complexity, Accuracy and Fluency in SLA*, edited by Alex Housen, Folkert Kuiken, and Ineke Vedder, 21–46. Amsterdam: John Benjamins. Accessed 17 April 2023. ProQuest.

———. 2014. "Conceptualizing and Measuring Short-Term Changes in L2 Writing Complexity." *Journal of Second Language Writing* 26, no. Dec: 42–65. Accessed 17 April 2023. Elsevier.

Casanave, Christine Pearson. 1994. "Language Development in Students' Journals." *Journal of Second Language Writing* 3, no. 3: 179–201. Accessed 22 October 2023. Elsevier.

Cooper, Thomas C. 1976. "Measuring Written Syntactic Patterns of Second Language Learners of German." *The Journal of Educational Research* 69, no. 5: 176–183. Accessed 22 October 2023. EBSCO.

Crossley, Scott A, and Danielle S. McNamara. 2012. "Detecting the First Language of Second Language Writers Using Automated Indices of Cohesion, Lexical Sophistication, Syntactic Complexity and Conceptual Knowledge." In *Approaching Language Transfer through Text Classification*, edited by Scott Jarvis and Scott A. Crossley, 106–126. Bristol: Multilingual Matters. Accessed 3 May 2023. ProQuest.

Cunningham, Kevin T., and Katarina L. Haley. 2020. "Measuring Lexical Diversity for Discourse Analysis in Aphasia: Moving-Average Type-Token Ratio and Word Information Measure." *Journal of Speech, Language, and Hearing Research* 63, no. 3: 710–721. Accessed 17 September 2023. EBSCO.

Dahl, Östen. 2004. *The Growth and Maintenance of Linguistic Complexity*. Vol. 71. Amsterdam: John Benjamins. Accessed 15 September 2023. ProQuest.

De Clercq, Bastien. 2015. "The Development of Lexical Complexity in Second Language Acquisition: A Cross-Linguistic Study of L2 French and English." *Eurosla Yearbook* 15: 69–94. Accessed 17 April 2023. EBSCO.

De Clercq, Bastien, and Alex Housen. 2017. "A Cross-Linguistic Perspective on Syntactic Complexity in L2 Development: Syntactic Elaboration and Diversity." *The Modern Language Journal* 101, no. 2: 315–334. Accessed 20 March 2023. EBSCO.

Díez-Bedmar, María Belén, and Pascual Pérez-Paredes. 2020. "Noun Phrase Complexity in Young Spanish EFL Learners' Writing: Complementing Syntactic Complexity Indices with Corpus-Driven Analyses." *International Journal of Corpus Linguistics* 25, no. 1: 4–35. Accessed 20 March 2023. EBSCO.

EF Education First. 2022. "EF English Proficiency Index: A Ranking of 111 Countries and Regions by English Skills." *EF Education First*. Accessed 18 May 2023. https://www.ef.com/assetscdn/WIBIwq6RdJvcD9bc8RMd/cefcom-epi-site/reports/2022/ef-epi-2022-english.pdf.

Ehret, Katharina, and Benedikt Szmrecsanyi. 2019. "Compressing Learner Language: An Information-Theoretic Measure of Complexity in SLA Production Data." *Second Language Research* 35, no. 1: 23–46. Accessed 20 March 2023. SAGE journals.

Ellis, Rod, and Gary Barkhuizen. 2005. *Analysing Learner Language*. Oxford: Oxford University Press.

European Commission. 2012. "Europeans and their Languages." *Special Eurobarometer* 386. Accessed 30 March 2023. https://europa.eu/eurobarometer/surveys/detail/1049.

EU = European Union. n.d. "Hungary." *European Union*. Accessed 18 September 2023. https://european-union.europa.eu/principles-countries-history/country-profiles/hungary_en.

Fergadiotis, Gerasimos. 2011. "Modeling Lexical Diversity Across Language Sampling and Estimation Techniques." Doctoral Dissertation, Arizona State University. Accessed 15 April 2024. https://core.ac.uk/download/pdf/79563501.pdf.

Igazságügyi Minisztérium. 2020. "A Nemzeti alaptanterv kiadásáról, bevezetéséről és alkalmazásáról szóló 110/2012." [On the Publication, Introduction, and Application of the National Core Curriculum110/2012]. *Magyar Közlöny*. Accessed 18 September 2023. https://magyarkozlony.hu/dokumentumok/3288b6548a740b9c8daf918a399a0bed1985db0f/letoltes.

Harkio, Noora, and Päivi Pietilä. 2016. "The Role of Vocabulary Breadth and Depth in Reading Comprehension: A Quantitative Study of Finnish Efl Learners." *Journal of Language Teaching and Research* 7, no. 6: 1079–1088. Accessed 17 September 2023. ProQuest.

Hunt, Kellogg W. 1965. "Grammatical Structures Written at Three Grade Levels." *NCTE Research Report No. 3.* Champaign: National Council of Teachers of English. Accessed 1 February 2024. ERIC.

Hwang, Hyun-Bin, and Charlene Polio. 2023. "Text Length Effects on the Reliability of Syntactic Complexity Indices." *Research Methods in Applied Linguistics* 2, no. 3: 100085. Accessed 22 October 2023. Elsevier.

Jiang, Jingyang, Peng Bi, and Haitao Liu. 2019. "Syntactic Complexity Development in the Writings of EFL Learners: Insights from a Dependency Syntactically-Annotated Corpus." *Journal of Second Language Writing* 46: 100666. Accessed 23 October 2023. Elsevier.

Kachru, Braj B. 1982. *The Other Tongue: English Across Cultures*. Urbana, IL: University of Illinois Press.

Kim, Minkyung, and Scott A. Crossley. 2018. "Modeling Second Language Writing Quality: A Structural Equation Investigation of Lexical, Syntactic, and Cohesive Features in Source-Based and Independent Writing." *Assessing Writing* 37: 39–56. Accessed 3 May 2023. Elsevier.

Kim, Minkyung, Scott A. Crossley, and Kristopher Kyle. 2018. "Lexical Sophistication as a Multidimensional Phenomenon: Relations to Second Language Lexical Proficiency, Development, and Writing Quality." *The Modern Language Journal* 102, no. 1: 120–141. Accessed 17 April 2023. EBSCO.

Kisselev, Olesya, Rossina Soyan, Dmitrii Pastushenkov, and Jason Merrill. 2022. "Measuring Writing Development and Proficiency Gains Using Indices of Lexical and Syntactic Complexity: Evidence from Longitudinal Russian Learner Corpus Data." *The Modern Language Journal* 106, no. 4: 798–817. Accessed 3 May 2023. Wiley Online Library.

KOTUS = Kotimaisten kielten keskus. n.d. "Kielet" [Languages]. *Kotimaisten kielten keskus*. Accessed 9 May 2023. https://www.kotus.fi/kielitieto/kielet.

KSH = Központi Statisztikai Hivatal. n.d. "Idegen nyelvet tanulók az általános iskolában" [Foreign Language Learners in Primary School]. *Központi Statisztikai Hivatal*. Accessed 18 September 2023. https://www.ksh.hu/stadat_files/okt/hu/okt0009.html.

Kuiken, Folkert, and Ineke Vedder. 2011. "Task Complexity and Linguistic Performance in L2 Writing and Speaking: The Effect of Mode." In *Second Language Task Complexity*: *Researching the Cognition Hypothesis of Language Learning and Performance*, edited by Peter Robinson, 91–104. Amsterdam: John Benjamins. Accessed 18 May 2023. EBSCO.

———. 2012. "Syntactic Complexity, Lexical Variation and Accuracy as a Function of Task Complexity and Proficiency Level in L2 Writing and Speaking." In *Dimensions of L2 Performance and Proficiency*, edited by Alex Housen, Folkert Kuiken, and Ineke Vedder, 143–169. Amsterdam: John Benjamins.

———. 2019. "Syntactic Complexity Across Proficiency and Languages: L2 and L1 Writing in Dutch, Italian and Spanish." *International Journal of Applied Linguistics* 29, no. 2: 192–210. Accessed 20 March 2023. EBSCO.

Kyle, Kristopher. 2016. "Measuring Syntactic Development in L2 Writing: Fine Grained Indices of Syntactic Complexity and Usage-Based Indices of Syntactic Sophistication." Dissertation, Georgia State University. Accessed 2 October 2023. https://scholarworks.gsu.edu/cgi/viewcontent.cgi?article=1035&context=alesl_diss.

Kyle, Kristopher, and Scott A. Crossley. 2018. "Measuring Syntactic Complexity in L2 Writing Using Fine-Grained Clausal and Phrasal Indices." *The Modern Language Journal* 102, no. 2: 333–349. Accessed 15 September 2023. EBSCO.

Lan, Ge, Qiusi Zhang, Kyle Lucas, Yachao Sun, and Jie Gao. 2022. "A Corpus-Based Investigation on Noun Phrase Complexity in L1 and L2 English Writing." *English for Specific Purposes* 67: 4–17. Accessed 3 May 2023. Elsevier.

Larson-Hall, Jenifer. 2016. *A Guide to Doing Statistics in Second Language Research Using SPSS and R*. 2nd ed. London: Routledge. Accessed 21 February 2024. EBSCO.

Larsson, Tove, and Henrik Kaatari. 2020. "Syntactic Complexity Across Registers: Investigating (in)formality in Second-Language Writing." *Journal of English for Academic Purposes* 45: 100850. Accessed 23 October 2023. Elsevier.

Lee, Jongbong. 2021. "Using Corpus Analysis to Extend Experimental Research: Genre Effects in L2 Writing." *System* 100: 102563. Accessed 20 March 2023. Elsevier.

Lei, Lei, Ju Wen, and Xiaohu Yang. 2023. "A Large-Scale Longitudinal Study of Syntactic Complexity Development in EFL Writing: A Mixed-Effects Model Approach." *Journal of Second Language Writing* 59: 100962. Accessed 23 October 2023. Elsevier.

Liu, Liming, and Lan Li. 2016. "Noun Phrase Complexity in EFL Academic Writing : A Corpus-Based Study of Postgraduate Academic Writing." *Journal of Asia TEFL* 13, no. 1: 48–65. Accessed 20 March 2023. ProQuest.

Lu, Xiaofei. 2010. "Automatic analysis of syntactic complexity in second language writing." *International Journal of Corpus Linguistics* 15 no. 4: 474–496. Accessed 11 April 2023. EBSCO.

———. 2011. "A corpus-based evaluation of syntactic complexity measures as indices of college-level ESL writers' language development." *TESOL Quarterly* 45, no. 1: 36–62. Accessed 15 September 2023. EBSCO.

———. 2012. "The Relationship of Lexical Richness to the Quality of ESL Learners' Oral Narratives." *The Modern Language Journal* 96, no. 2: 190–208. Accessed 11 April 2023. EBSCO.

Lu, Xiaofei, and Haiyang Ai. 2015. "Syntactic Complexity in College-Level English Writing: Differences Among Writers with Diverse L1 Backgrounds." *Journal of Second Language Writing* 29, no. Sep: 16–27. Accessed 15 September 2023. Elsevier.

Mehnert, Uta. 1998. "The Effect of Different Lengths of Time for Planning on Second Language Performance." *Studies in Second Language Acquisition* 20, no. 1: 83–108. Accessed 14 September 2023. CambridgeCore.

Mäkynen, Iida. 2023. "L2 Development During a Short English Course : Studying Finnish University Students' Written Syntactic Complexity." Master's Thesis, University of Turku. Accessed 17 October 2023.

Nasseri, Maryam, and Paul Thompson. 2021. "Lexical Density and Diversity in Dissertation Abstracts: Revisiting English L1 Vs. L2 Text Differences." *Assessing Writing* 47: 100511. Accessed 23 October 2023. Elsevier.

Nation, I. S. P. 2001. *Learning Vocabulary in Another Language*. Cambridge: Cambridge University Press. Accessed 23 October 2023. CambridgeCore.

Norris, John M., and Lourdes Ortega. 2009. "Towards an Organic Approach to Investigating CAF in Instructed SLA: The Case of Complexity." *Applied Linguistics* 30, no. 4: 555–578. Accessed 20 March 2023. EBSCO.

OED = *Oxford English Dictionary*. 2nd edition. Oxford: Oxford University Press. Available online by subscription at http://www.oed.com.

OPH = Finnish National Agency for Education. 2014. "Perusopetuksen opetussuunnitelman perusteet 2014" [The National Core Curriculum of Basic Education 2014]. *Finnish National Agency for Education.* Accessed 18 September 2023. https://www.oph.fi/sites/default/files/documents/perusopetuksen_opetussuunnitelman_ perusteet_2014.pdf.

———. 2019. "Perusopetuksen opetussuunnitelman perusteiden 2014 muutokset ja täydennykset koskien A1 kielen opetusta vuosiluokilla 1–2" [The Changes and Additions to the National Core Curriculum 2014 of Basic Education Regarding the Teaching of the A1 Language in Grades 1–2]. *Finnish National Agency for Education.* Accessed 18 September 2023. https://www.oph.fi/sites/default/files/documents/perusopetuksen_vuosiluokkien_1-2_a1-kielen_opetussuunnitelman_perusteet.pdf.

Ortega, Lourdes. 2003. "Syntactic Complexity Measures and Their Relationship to L2 Proficiency: A Research Synthesis of College-level L2 Writing." *Applied Linguistics* 24, no. 4: 492–518. Accessed 20 March 2023. EBSCO.

Pallotti, Gabriele. 2009. "CAF: Defining, Refining and Differentiating Constructs." *Applied Linguistics* 30, no. 4: 590–601. Accessed 15 September 2023. EBSCO.

Peltomäki, Nelli. 2018. "Vocabulary Analysis of Finnish University Students' English Essays: a Lexical Sophistication and Content Analysis Study." Master's Thesis, University of Turku. Accessed 17 October 2023.

Phuoc, Vo Dinh, and Jessie S. Barrot. 2022. "Complexity, Accuracy, and Fluency in L2 Writing Across Proficiency Levels: A Matter of L1 Background?" *Assessing Writing* 54: 100673. Accessed 20 March 2023. Elsevier.

Qin, Wenjuan, and Paola Uccelli. 2020. "Beyond Linguistic Complexity: Assessing Register Flexibility in EFL Writing Across Contexts." *Assessing Writing* 45: 100465. Accessed 18 September 2023. Elsevier.

Rahayu, Famala Eka Sanhadi, Aries Utomo, and Ririn Setyowati. 2021. "Syntactic and Lexical Complexity of Undergraduate Students' Essays: A Comparison Study Between L1 and L2 Writings." *Indonesian Journal of English Language Teaching and Applied Linguistics* 5, no. 2: 251–263. Accessed 18 September 2023. ERIC.

Read, John. 2000. *Assessing Vocabulary*. Cambridge: Cambridge University Press.

Ricci, Vito. 2005. "Fitting Distributions with R." Release 0.4-21. Accessed 21 February 2024. http://cran.r-project.org/doc/con-trib/Ricci-distributions-en.pdf.

Robinson, Peter. 2001. "Task Complexity, Cognitive Resources, and Syllabus Design: A Triadic Framework for Examining Task Influences on SLA." In *Cognition and Second Language Instruction*, edited by Peter Robinson, 287–318. Cambridge: Cambridge University Press. Accessed 18 May 2023. CambridgeCore.

Sarte, Kayla Marie, and Ksenia Gnevsheva. 2022. "Noun Phrasal Complexity in ESL Written Essays Under a Constructed-Response Task: Examining Proficiency and Topic Effects." *Assessing Writing* 51: 100595. Accessed 20 March 2023. Elsevier.

Seidinejad, Leila, and Zohreh Nafissi. 2018. "Developing Lexical Complexity in EFL Students' Essays via Creative Thinking Techniques." *Pertanika Journal of Social Science & Humanities* 26, no. 3: 1697–1712. Accessed 23 October 2023. Pertanika.

Skehan, Peter. 1996. "A Framework for the Implementation of Task-Based Instruction." *Applied Linguistics* 17, no. 1: 38–62. Accessed 11 May 2023. EBSCO.

———. 2009. "Modelling Second Language Performance: Integrating Complexity, Accuracy, Fluency, and Lexis." *Applied Linguistics* 30, no. 4: 510–532. Accessed 15 September 2023. EBSCO.

Skehan, Peter, and Pauline Foster. 2001. "Cognition and Tasks." In *Cognition and Second Language Instruction*, edited by Peter Robinson, 183–205. Cambridge: Cambridge University Press.

SUKOL = Suomen kieltenopettajien liitto ry. n.d. "Tilastotietoa kielivalinnoista" [Statistics about Language Choices]. *SUKOL*. Accessed 18 September 2023. https://www.sukol.fi/liitto/tilastot/tilastotietoa_kielivalinnoista

Thi, Nang Kham, De Van Vo, and Marianne Nikolov. 2023. "Investigating Syntactic Complexity and Language-Related Error Patterns in EFL Students' Writing: Corpus-Based and Epistemic Network Analyses." *Language Learning in Higher Education* 13, no. 1: 127–151. Accessed 14 September 2023. De Gruyter.

Tsai, Yea-Ru. 2021. "Exploring the Effects of Corpus-Based Business English Writing Instruction on EFL Learners' Writing Proficiency and Perception." *Journal of Computing in Higher Education* 33, no. 2: 475–498. Accessed 23 October 2023. ProQuest.

Wolfe-Quintero, Kate, Shunji Inagaki, and Hae-Young Kim. 1998. *Second Language Development in Writing : Measures of Fluency, Accuracy and Complexity*. Honolulu: Second Language Teaching & Curriculum Center.

Yoon, Hyung-Jo. 2017. "Linguistic Complexity in L2 Writing Revisited: Issues of Topic, Proficiency, and Construct Multidimensionality." *System* 66: 130–141. Accessed 23 October 2023. Elsevier.

Yoon, Hyung-jo, and Charlene Polio. 2017. "The Linguistic Development of Students of English as a Second Language in Two Written Genres." *TESOL Quarterly* 51, no. 2: 275–301. Accessed 17 September 2023. EBSCO.

Zhang, Chao, and Shumin Kang. 2022. "A Comparative Study on Lexical and Syntactic Features of ESL Versus EFL Learners' Writing." *Frontiers in Psychology* 13: 1002090. Accessed 3 May 2023. DOAJ.

Zhang, Ronggen. 2022. "Corpus-Based Study on Syntactic Complexity of Texts by L1 and L2 Learners." *ITM Web of Conferences* 45: 01081. Accessed 3 May 2023. DOAJ.

## Finnish Summary

Tämä pro gradu -tutkielma vertailee *lingvististä kompleksisuutta* suomalaisten ja unkarilaisten englannin kielen yliopisto-opiskelijoiden esseissä. Tutkielman tarkoituksena on kartoittaa *syntaktista* ja *leksikaalista kompleksisuutta oppijakielessä* ja siihen vaikuttavia tekijöitä kuten kielitaustoja ja sosiokulttuurisia tekijöitä sekä tarkastella syntaktisten rakenteiden ja leksikaalisten elementtien ilmenemistä. Pro gradu käsittelee ilmiötä tutkimalla korpusaineistoa määrällisesti käyttäen hyödyksi automaattisia kielen analysaattoreita sekä tarkastelemalla muutamaa esimerkkitekstiä laadullisesti. Oppijakielen kompleksisuutta on tutkittu runsaasti eri näkökulmista ja erilaisilla metodeilla. Kuitenkin olemassa oleva tutkimus keskittyy tavallisesti yhteen kompleksisuustyyppiin eikä näin pysty tarkastelemaan ilmiötä kokonaisvaltaisesti. Tämä tutkielma pyrkii vastaamaan tutkimusaukkoon sisällyttämällä vertailuun kaksi eri kompleksisuuden tyyppiä. Lisäksi tutkielma keskittyy mielenkiintoiseen vertailuun, jossa verrattavat ryhmät edustavat ensikieliltään lingvistisesti kaukaisia sukukieliä. Tutkielma pyrkii tuottamaan uutta tietoa kirjallisesta oppijakielestä ja sen moninaisesta luonteesta, mikä voi olla hyödyksi varsinkin kielenoppijoiden parissa työskenteleville. Pro gradun tutkimuskysymykset ovat seuraavanlaiset:

1A. Miten suomalaisten ja unkarilaisten yliopisto-opiskelijoiden englanninkieliset tekstit eroavat syntaktisessa kompleksisuudessa?

1B. Miten suomalaisten ja unkarilaisten yliopisto-opiskelijoiden englanninkieliset tekstit eroavat leksikaalisessa kompleksisuudessa?

2. Miten syntaktiset rakenteet ja leksikaaliset elementit ilmenevät esseissä, joiden on havaittu osoittavan korkeaa tai heikkoa kompleksisuutta?

Lingvistinen kompleksisuus *toisen kielen oppimisessa* on monitahoinen käsite, joka on jokseenkin hankala määritellä. Tämä tutkielma käsittelee *absoluuttista kompleksisuutta* ja määrittelee sen liittyvän oppijan taitoon tuottaa monimutkaista kieltä. Kompleksisuuden tarkastelu oppijakielessä juontaa juurensa Skehanin (1996) teoriaan, jossa oppijan kielitaito jaetaan kolmeen ulottuvuuteen: *sujuvuuteen*, *virheettömyyteen* ja kompleksisuuteen (CAF). Malli on saanut kritiikkiä, mutta on suosittu lähestymistapa oppijakielen tutkimisessa. Tieteellinen tutkimus kuitenkin osoittaa, että edistyneemmät L2-oppijat tuottavat sekä sujuvampaa, virheettömämpää että kompleksisempaa kieltä kuin vähemmän edistyneet L2-oppijat.

Oppijakielen kompleksisuus voidaan jakaa eri tyyppeihin, joista tämä tutkielma keskittyy kahteen: syntaktiseen ja leksikaaliseen kompleksisuuteen. Syntaktinen kompleksisuus on hyvin monitahoinen ilmiö, josta toisinaan puhutaan myös "kieliopillisena kompleksisuutena". Toisinaan nämä kaksi termiä käsitetään myös kahtena eri ilmiönä. Tämä tutkielma käyttää termiä syntaktinen kompleksisuus. Syntaktinen kompleksisuus esiintyy kielenoppijan taitona käyttää laajaa valikoimaa sekä tavanomaisia että hienostuneita syntaktisia rakenteita. Syntaktista kompleksisuutta mitataan usein tarkastelemalla lauseiden välisiä alistus- ja rinnastussuhteita, rakenteiden pituutta ja *T-yksiköitä* (T-unit). Tieteellinen tutkimus osoittaa, että edistyneemmät kielenoppijat käyttämät enemmän alistusrakenteita, pidempiä rakenteita ja pidempiä T-yksiköitä kuin vähemmän edistyneet kielenoppijat. Leksikaalinen kompleksisuus esiintyy kielenoppijan taitona käyttää laajaa valikoimaa sekä tavanomaisia että hienostuneita sanoja. Leksikaalista kompleksisuutta mitataan usein tarkastelemalla kielenoppijan sanaston laajuutta, vaihtelevuutta ja hienostuneisuutta. Tieteellinen tutkimus osoittaa, että edistyneemmät kielenoppijat käyttävät laajempaa, vaihtelevampaa ja hienostuneempaa sanastoa kuin vähemmän edistyneet kielenoppijat.

Kielten välistä vaikutusta lingvistiseen kompleksisuuteen on tutkittu aiemmin, mutta vain rajallisesti. Aikaisempi tutkimus osoittaa, että kielenoppijan kielitaustalla voi olla vaikutusta lingvistiseen kompleksisuuteen toisessa kielessä ja että eri kieliyhdistelmillä voi olla siihen erilaisia vaikutuksia. Myös eroja *englanti toisena kielenä* -oppijoiden (ESL) ja *englanti vieraana kielenä* -oppijoiden (EFL) välillä on kartoitettu. Lisäksi tulee muistaa, että myös sosiokulttuuriset tekijät ja opetukseen liittyvät seikat selittävät samankaltaisuuksia ja eroja eri oppijoiden välillä.

Tämä tutkielma keskittyy englannin kielen oppijoihin, jotka ovat taustaltaan kahdenlaisia: ensikielenään suomen kieltä puhuvia ja ensikielenään unkarin kieltä puhuvia. Suomen kieli ja unkarin kieli kuuluvat molemmat samaan kieliryhmään eli suomalais-ugrilaisiin kieliin ja jakavat jonkin verran samankaltaisuuksia eri kielen osa-alueita tarkastellessa. Kielet ovat lingvistisesti katsoen myös kuitenkin etäisiä. Suomen kieli on itämerensuomalainen kieli ja unkarin kieli on ugrilainen kieli, mikä tarkoittaa, että vaikka kielet kuuluvatkin samaan kieliryhmään ne ovat tarpeeksi kaukana toisistaan, jotta niiden puhujat eivät keskenään ymmärrä toisiaan. Englannin kieltä opetetaan sekä Suomessa että Unkarissa vieraana kielenä ja on molemmissa valtioissa kouluissa eniten opiskelluin vieras kieli. Joidenkin tutkimusten mukaan suomalaiset ovat edistyneempiä englannin kielen käyttäjiä verrattuna unkarilaisiin englannin kielen käyttäjiin.

Tutkielman hypoteesin mukaan suomalaisten ja unkarilaisten englannin kielen yliopisto-opiskelijoiden välillä on eroja lingvistisessä kompleksisuudessa, mutta tarkempia hypoteeseja on vaikea asettaa olemassa olevan tutkimuksen pohjalta. Tutkielman primäärisenä lähteenä käytetään ACALEX-korpusta, joka sisältää sekä kirjallisia että suullisia oppijakielinäytteitä aikaväliltä 2016–2018 kolmesta eurooppalaisesta yliopistosta: Turun yliopistosta, Åbo akademista ja Szegedin yliopistosta. Kielinäytteen antaneiden opiskelijoiden henkilöllisyys on piilotettu korpuksessa käyttämällä osallistujista koodeja. Tässä tutkielmassa hyödynnetään korpuksesta kolmeakymmentä suomalaisten ja kolmeakymmentä unkarilaisten englannin kielen yliopisto-opiskelijoiden kirjoittamia esseitä, joissa opiskelijat pohtivat tulevaisuuden urasuunnitelmiaan.

Esseitä analysoitiin käyttämällä Xiaofei Lun oppijakielen tarkastelua varten kehittämiä automaattisia analysaattoreita, sillä näin pystyttiin prosessoimaan tehokkaasti ja tarkasti suuri määrä oppijakieltä. Muun ohella Lun kehittämät analysaattorit ovat aikaisemman tutkimuksen mukaan lupaavia reliabiliteetin suhteen, käyttäjäystävällisiä ja soveltuvat erinomaisesti tutkielman näkökulman kanssa yhteen. Tässä tutkielmassa käytettiin analysaattoreiden ilmaisia selainversioita. Syntaktisen kompleksisuuden analysointiin käytettiin *The L2 Syntactic Complexity Analyzer* -ohjelmaa (L2SCA), joka analysoi kirjallista oppijakieltä käyttämällä 14 mittaria. L2SCA-ohjelmaa on käytetty laajasti aikaisemmissa tutkimuksissa ja sen on todettu antavan luotettavia tuloksia oppijakielen syntaktisesta kompleksisuudesta. Tässä tutkielmassa hyödynnettiin kymmentä Lun (2011) suosittelemaa mittaria, jotka mittaavat rakenteiden pituutta (MLC, MLS ja MLT), lause–virkesuhdetta (C/S), lauseiden alistuksen määrää (DC/C ja DC/T), lauseiden rinnastuksen määrää (CP/C ja CP/T) ja syntaktisten rakenteiden ja tuotettujen yksiköiden suhdetta (CN/C ja CN/T). Leksikaalisen kompleksisuuden analysointiin käytettiin *The Lexical Complexity Analyzer* -ohjelmaa (LCA), joka käyttää kirjallisen oppijakielen analysointiin 25 mittaria, jotka voidaan jaotella kolmeen ryhmään: leksikaalinen tiheys, leksikaalinen hienostuneisuus ja leksikaalinen vaihtelu. Vaikka LCA-ohjelmaa on käytetty laajasti tutkimuksissa, sen luotettavuutta ei olla arvioitu laajasti.

Ennen esseiden analysoimista analysaattoreilla, ne oikoluettiin ja niihin tehtiin mahdollisia muokkauksia kuten kirjoitusvirheiden korjauksia, jotta automaattisen analyysin tarkkuus voitiin varmistaa. Esseet myös kategorisoitiin LCA:n kahden tukeman kielistandardin mukaan joko brittienglanniksi tai amerikanenglanniksi. Tämän jälkeen esseet prosessoitiin analysaattoreissa ja analysaattoreiden antamia tuloksia käsiteltiin vielä Microsoft Excelillä sekä IBM SPSS Statistics 27 -ohjelmalla. Tulosten tilastollista merkitsevyyttä mitattiin Mann-

Whitney U -testillä. Määrällisen analyysin jälkeen esseet, jotka oli havaittu syntaktisesti ja leksikaalisesti kompleksisimmiksi analysoitiin yhä laadullisesti.

Tutkielman analyysin tulokset osoittavat, että suomalaiset englannin kielen yliopisto-opiskelijat tuottavat keskimäärin lingvistisesti kompleksimpaa kirjallista kieltä verrattuna unkarilaisiin englannin kielen yliopisto-opiskelijoihin. Syntaktisen kompleksisuuden osalta mielenkiintoista on muun muassa, että unkarilaiset opiskelijat tuottavat keskimäärin lähes kaikkia L2SCA:n huomioon ottamia syntaktisia rakenteita enemmän kuin suomalaiset opiskelijat, vaikka verrattujen ryhmien keskimääräiset sanamäärät ovat lähes identtiset. Tämä tarkoittaa, että unkarilaisten opiskelijoiden tuottamat rakenteet ovat siis keskimäärin lyhyempiä verrattuna suomalaisten opiskelijoiden tuottamiin rakenteisiin. Suomalaiset opiskelijat osoittavat tilastollisesti merkitsevästi korkeampaa syntaktista kompleksisuutta puolessa käytetyistä L2SCA:n mittareista. Selkein tilastollisesti merkitsevä ero syntaktisessa kompleksisuudessa suomalaisten ja unkarilaisten opiskelijoiden välillä on suomalaisten opiskelijoiden alttius kirjoittaa pidempiä rakenteita (MLC, MLS ja MLT) kuin unkarilaiset opiskelijat. Verrattujen ryhmien väliltä löytyy myös samankaltaisuuksia syntaktisessa kompleksisuudessa. Tulokset lause–virkesuhteelle (C/S) ovat samankaltaiset. Myös lauseiden alistuksen määrä (DC/C ja DC/T), lauseiden rinnastuksen määrä (CP/C ja CP/T) olivat ryhmillä keskimäärin samankaltaiset. Kuitenkin keskimäärin suomalaiset opiskelijat osoittavat unkarilaisia opiskelijoita korkeampaa kompleksisuutta syntaktisten rakenteiden suhteen.

Leksikaalisen kompleksisuuden analyysin tulokset ovat yhä selkeämmät kuin syntaktisen kompleksisuuden analyysin tulokset. Suomalaiset opiskelijat osoittavat tilastollisesti merkitsevästi keskimäärin korkeampaa leksikaalista kompleksisuutta 18:ssa 25:stä LCA:n mittarista verrattuna unkarilaisiin opiskelijoihin. Suomalaisten opiskelijoiden kirjoittamat esseet ovat keskimäärin unkarilaisten opiskelijoiden kirjoittamia esseitä leksikaalisesti tiheämpiä ja esseissä käytettävät verbit ovat unkarilaisten opiskelijoiden käyttämiä verbejä hienostuneempia. Eroa suurentaa yksi poikkeuksellinen suomalaisen opiskelijan kirjoittama essee, jossa esiintyy erityisen hienostuneita sanoja. Ryhmien välillä ei kuitenkaan löydetty eroa kahdessa leksikaalisen hienostuneisuuden mittarissa (LS1 ja LS2) eli voidaan todeta, että ryhmillä on myös samankaltaisuuksia leksikaalisessa hienostuneisuudessa.

Leksikaalisen vaihtelun suhteen suomalaiset opiskelijat saavat tilastollisesti merkitsevästi keskimäärin korkeampia tuloksia kuin unkarilaiset opiskelijat lähes kaikista eri sanojen

määrää (NDW) ja tyyppi–esiintymäsuhdetta (TTR) arvioivista mittareista. Tämä tarkoittaa, että suomalaisten opiskelijoiden käyttämä sanasto on keskimäärin vaihtelevampaa kuin unkarilaisten käyttämä sanasto. Suomalaisten opiskelijoiden verbien vaihtelu on tilastollisesti merkitsevästi keskimäärin unkarilaisiin opiskelijoihin verrattuna suurempaa osan mittareista mukaan. Ryhmien välillä on samankaltaisuuksia sisältösanojen ja substantiivien vaihtelussa. Eroja löytyy myös adjektiivien, adverbien ja määritteiden vaihtelussa, joissa suomalaiset opiskelijat saavat tilastollisesti merkitsevästi keskimäärin korkeamman tuloksen kuin unkarilaiset opiskelijat. On siis mielekästä todeta, että suomalaiset opiskelijat osoittavat keskimäärin korkeampaa leksikaalista kompleksisuutta kuin unkarilaiset opiskelijat.

Syntaktisen ja leksikaalisen kompleksisuuden arvioinnin lisäksi tutkimuksessa tarkasteltiin ääriesimerkkejä eli esseitä, jotka osoittavat erityisen korkeaa tai heikkoa syntaktista tai leksikaalista kompleksisuutta. Esseistä tarkasteltiin, miten syntaktiset rakenteet ja leksikaaliset elementit ilmenevät niissä. ACALEX-1102:n, kirjoittama essee osoittaa keskivertoa korkeampaa syntaktista kompleksisuutta puolen käytettyjen mittareista mukaan, kun esseetä verrataan muihin suomalaisiin opiskelijoihin. Loppujen mittareiden mukaan essee suoriutuu keskimääräisesti. Esseen syntaktiset rakenteet ovat erityisen pitkiä ja siinä käytetään myös paljon komplekseja nominaalimuotoja. Esseen syntaktinen kompleksisuus ei siis johdu ainoastaan pitkistä rakenteista vaan myös niiden taitavasta ja hienostuneesta käytöstä. Alhaisimman tuloksen syntaktisessa kompleksisuudessa saanut ACALEX-1211:n essee saa keskivertoa heikomman tuloksen kaikista L2SCA:n mittareista ja sisältää paljon lyhyitä rakenteita ja vähän lauseiden välisiä alistussuhteita. Vähäinen kompleksisuus syntaktisissa rakenteissa vaikuttaa negatiivisesti myös esseen kykyyn välittää merkitystä.

Leksikaalisen kompleksisuuden analyysin mukaan ACALEX-1127:n essee on leksikaalisesti kompleksein ja sen voi huomata nopeasti esseetä tarkastelemalla. Essee on suurimman osan LCA:n mittareista mukaan keskivertoa kompleksimpi leksikaalisesti. Esseessä käytetään sanastoa vaihtelevasti ja sanasto on suurelta määrin hienostunutta. Korkean leksikaalisen kompleksisuuden vuoksi esseetä on jopa hankala lukea ja merkitystä on paikoin hankala löytää. On syytä epäillä, että esseen kirjoittaja on erityisesti pyrkinyt korkeaan leksikaaliseen kompleksisuuteen. Esseetä on mielekästä verrata ACALEX-1197:n esseeseen, joka osoittaa heikkoa leksikaalista kompleksisuutta. Essee osoittaa keskivertoa heikompaa leksikaalista kompleksisuutta yli puolessa käytetyistä mittareista verrattaessa muihin unkarilaisiin opiskelijoihin. Essee toistaa leksikaalisia sanoja usein ja käyttää yleisimmin tavanomaista sanastoa.

Tutkimuksen tulokset osoittavat, miten kirjallinen oppijakieli, kielitausta ja lingvistinen kompleksisuus toisen kielen oppimisessa vaikuttavat toisiinsa ja miten monimutkaisia ilmiöt ovat. Tulokset voivat olla hyödyllisiä erityisesti opetusalalla työskenteleville, joiden tulee ymmärtää ilmiöiden vuorovaikutuksen laajuus. Myös uusia opetussuunnitelmia kehittäessä tulisi huomioida vaihtelu oppijoiden taustoissa. Jatkotutkimus voi esimerkiksi perehtyä lisää erilaisten kielitaustojen vaikutukseen, tarkastella kompleksisuuden eroja oppijan ensikielessä ja toisessa kielessä, tutkia kompleksisuuden, virheettömyyden ja sujuvuuden välistä suhdetta, verrata kirjallista ja puhuttua oppijakieltä tai keskittyä eroihin ja selittäviin tekijöihin suomalaisten ja unkarilaisten oppijoiden yleisissä taitotasoissa.