

**INVESTIGATING THE  
CONTENT VALIDITY OF THE  
IELTS LISTENING TEST  
THROUGH THE USE OF  
LEXICAL BUNDLES**

**by Phuoc Ha Thien Nguyen**

British Council's Master's Dissertation Awards 2023  
Special Commendation

**Nottingham Trent University  
MA in Teaching English to Speakers of Other Languages**

**INVESTIGATING THE CONTENT VALIDITY OF THE IELTS LISTENING TEST  
THROUGH THE USE OF LEXICAL BUNDLES**

by

**PHUOC HA THIEN NGUYEN  
N1055851**

Dissertation submitted to Nottingham Trent University School of Arts and Humanities, in partial fulfilment of the requirements of the degree of Master of Arts (Teaching English to Speakers of Other Languages)

**28<sup>th</sup> September, 2022**

**14,957 words**

## Table of Contents

<i>Copyright Statement</i> .....	<i>iii</i>
<i>Acknowledgement</i> .....	<i>iv</i>
<i>Abstract</i> .....	<i>v</i>
<b>1. Introduction</b> .....	<b>1</b>
1.1. Contextual background and research gap .....	1
1.2. Research aims and research questions .....	3
1.3. Overview of the following chapters .....	4
<b>2. Literature review</b> .....	<b>5</b>
2.1. Formulaic language.....	5
2.2. Lexical bundles .....	7
2.2.1. Definition and general characteristics .....	7
2.2.2. Structural and functional categorisation .....	8
2.3. Lexical bundles in academic contexts .....	11
2.4. Corpus analysis of lexical bundles .....	14
2.5. The IELTS test.....	16
2.5.1. Overview of the IELTS test.....	16
2.5.2. Validity evaluation of the IELTS test.....	17
2.6. Summary of the Pilot study.....	18
<b>3. Methodology</b> .....	<b>18</b>
3.1. Corpus compilation.....	19
3.2. Lexical bundle extraction .....	23
3.2.1. Computational extraction .....	23
3.2.2. Refinement .....	26
3.3. Theoretical frameworks for data analysis.....	27
3.3.1. Frequency analysis.....	27
3.3.2. Structural analysis.....	29
3.3.3. Functional analysis.....	31
<b>4. Results and discussions</b> .....	<b>34</b>
4.1. Frequency distribution of bundles across corpora .....	34
4.2. Structural patterns of bundles across corpora .....	40
4.2.1. Between IELTS S3 and Interactive ASE .....	40
4.2.2. Between IELTS S4 and Monologic ASE .....	47
4.3. Functional patterns of bundles across corpora .....	51
4.3.1. Between IELTS S3 and Interactive ASE .....	51
4.3.2. Between IELTS S4 and Monologic ASE .....	60

<b>5. Conclusion .....</b>	<b>67</b>
5.1. Summary .....	67
5.2. Pedagogical implications.....	68
5.3. Limitations.....	70
5.4. Future research and development .....	70
<b>References.....</b>	<b>72</b>
<b>Appendices.....</b>	<b>80</b>
Appendix 1. Sources for the IELTS Section 3 sub-corpus .....	80
Appendix 2. Sources for the IELTS Section 4 sub-corpus .....	81
Appendix 3. Sources for the Interactive ASE sub-corpus .....	82
Appendix 4. Sources for the Monologic ASE sub-corpus .....	84
Appendix 5. Structural classification of lexical bundles across corpora.....	85
Appendix 6. Functional classification of lexical bundles across corpora.....	89

## Copyright Statement

(1) Copyright in the text of this dissertation rests with the Author. Copies (by any process) either in full, or of extracts, may be made only in accordance with instructions given by the Author. This page must form part of any copies made. Further copies (by any process) of copies made in accordance with such instructions may not be made without the written permission of the Author.

(2) The ownership of any intellectual property rights which may be described in this dissertation is vested in Nottingham Trent University, subject to prior agreement to the contrary, and may not be made available for use by third parties without the written permission of the university, which will prescribe the terms and conditions of any such agreement.

## Acknowledgement

The completion of this dissertation is largely due to the help and support of many people who have encouraged and guided me over the past year. Special thanks must be given to my parents, who first started me on this journey of discovery and has supported me unconditionally through this process of learning. I am also particularly grateful to Amy, my dear tutor, who has been generous with her time in reading and re-reading my long drafts, and stimulating in the many fruitful discussion of its contents.

I would also like to extend my gratitude to my husband and other family members for their patience and understanding over the past year. Thanks also to all the classmates, especially Vy, Thao, and An for their support, guidance and friendship throughout the MA course.

Finally, love and thanks to my family and friends, especially my dear son, without whom this dissertation would not have been finished.

## Abstract

Listening comprehension is an essential skill in higher education context (Flowerdew, 1995), yet ESL students are often found to have difficulty understanding academic lectures (Anderson-Mejias, 1986) due to their failure to identify the speech's macro-structure and general meaning (Olsen and Huckin, 1990). This problem can be linked to students' restricted awareness of *lexical bundles*, defined by Biber et al. (1999) as "recurring sequences of three or more words" that signal forthcoming information (Neely and Cortes, 2009). Despite its importance to academic listening, however, few attempts have been made to examine the characteristics of lexical bundle use in academic speech, especially in EAP listening tests. This study thus addressed this gap by investigating the frequency, structural and functional patterns of lexical bundles in the International English Language Testing System (IELTS) listening test, which is required for entry to most English-medium universities. A corpus-driven analysis was conducted on the IELTS corpus, which consists of transcripts of official IELTS practice tests, and the reference corpus which was composed of selected transcripts of authentic academic speech events from the Michigan Corpus of Academic Spoken English (MICASE). Lexical bundles were computationally extracted and classified into specific structural and functional categories before being compared across corpora in terms of frequency distribution, structural and discursal patterns. The data analysis revealed that the IELTS listening test underuses lexical bundles and uses them differently from authentic academic speech. It is more direct in expressing obligation and evaluation, but less conversational and listener-centred. It also contains fewer hedging and idea expansion phrases than real-life discourse. These differences raise doubts to the test as a replica of academic spoken discourse, thus weakening its content validity.

# 1. Introduction

## 1.1. Contextual background and research gap

Listening comprehension is an essential communication skill in higher education context (Flowerdew, 1995), yet it is widely agreed to pose considerable challenges to ESL learners (e.g. James, 1977; Flowerdew 1992). Researchers such as Olsen and Huckin (1990) and Mendelsohn (2002) believe that although ESL students can comprehend individual words of a lecture, they struggle to follow its overall structure and extract the key points. This problem is argued to stem from listeners' reliance on bottom-up rather than top-down listening processing. While the former is characterised by an attention to every detail of the acoustic input (Buck, 2001; Morley, 2001), the later requires listeners to make associations with various types of prior knowledge to make sense of the general contents they are hearing.

The lack of top-down processing in academic listening can be linked to students' restricted awareness of *lexical bundles*. This is a term coined by Biber et al. (1999, p.990) to address frequently "recurring sequences of three or more words" in a text that are "the building blocks of extended strings of language". With this characteristic, lexical bundles can operate as a frame to signal forthcoming information, thus are argued to promote top-down listening (Neely and Cortes, 2009). This argument is strengthened by Chaudron and Richards (1986), who discovered that macro-markers similar to Biber et al. (2004)'s list of lexical bundles such as "*what I'm going to talk about today*" assist high-level information processing, which is crucial to lecture comprehension. This viewpoint is also shared by Biber & Barbieri (2007) and Chen & Chen (2020), who acknowledge that failure to understand the various functions of lexical bundles can impact



successful speech comprehension. Therefore, it could be concluded that lexical phrases play a key role in academic listening competence.

Despite its importance to academic performance, few studies have been done on lexical bundles in academic speech and EAP testing. An abundance of research in this field has investigated the characteristics of lexical bundles in various types of academic written discourse: academic prose (Biber et al., 1999), textbooks (Biber et al., 2004), L2 versus L1 writing (Chen & Baker, 2010), to name a few. Yet, few research explores the patterns in academic speech (Biber et al., 2004). Even these small number of studies are rather limited in their focus or problematic in their methodology. Nesi & Basturkmen (2006) and Neely & Cortes (2009) investigated the functional characteristics of a few lexical bundles used in university lectures. Biber et al. (2004), Biber & Barbieri (2007) and Simpson-Vlach and Ellis (2010) examined more aspects of various lexical bundles in different academic spoken registers. Nevertheless, the selection of corpus materials and lexical bundle extraction criteria in these studies raise some concerns over representativeness and reliability. Similarly, little has been done on the phraseological aspect of the language used in test papers. To the researcher's knowledge, although several studies have considered the use of lexical bundles in test takers' written and spoken response (Cooper, 2013; Read and Nation, 2006), existing research has not dealt with the test content to examine whether it captures the nature of language use in the real world. This gap generates a motivation for the current research to explore EAP assessment materials from the perspective of lexical bundles.

Being one of the world's most popular English proficiency tests, the IELTS test is thus chosen for this study. It has four modules, corresponding to four communication skills: listening, speaking, reading, and writing. Compared to the other three modules, the IELTS listening test remains

relatively under-researched (Cooper, 2013). A major line of research in this area evaluates the test's different types of validity, including predictive validity, construct irrelevance, and cognitive validity. However, none have dealt with the test's content validity, especially from the lexical dimension. Content validity reflects how fully the assessment materials represent the nature of knowledge and skills they are intended to measure (Green, 2014). From the IELTS listening test's point of view, its content validity involves the extent to which the test materials correspond to academic spoken language use, among other criteria. One key feature of the language used in this register is, as concluded above, lexical bundles. Hence, investigating lexical bundle use in the IELTS listening module would provide some evidence for the evaluation of this globally recognised test, which can then help test designers to make well-informed decisions.

## **1.2. Research aims and research questions**

Addressing the gap mentioned previously, this study aims to investigate the use of four-word lexical bundles in Sections 3 and 4 of the IELTS listening test as a mean to evaluate how close the test's language use simulates corresponding authentic situations. The scope is narrowed down to Sections 3 and 4 because unlike Sections 1 and 2, they focus exclusively on academic English, thus are relevant and comparable to the target language use domain of IELTS' student test takers (i.e. academic environment). The study investigated the following research questions:

1. To what extent do the frequency distributions of lexical bundles in the IELTS listening test differ from those of authentic academic spoken English?
2. To what extent do the structural patterns of lexical bundles in the IELTS listening test differ from those of authentic academic spoken English?

3. To what extent do the functional patterns of lexical bundles in the IELTS listening test differ from those of authentic academic spoken English?

A corpus-driven research approach was employed to allow the automatic extraction of bundles from large sets of data. First, two main corpora, the IELTS corpus and the Academic Spoken English (ASE) corpus were compiled. Lexical bundles were then computationally extracted and refined before being classified into appropriate categories. A detailed analysis of three aspects of lexical bundle use (i.e. frequency, structural and functional characteristics) between two corpora was conducted to reveal information about the two discourses.

### **1.3. Overview of the following chapters**

The following chapter, Literature Review, reviews the literature regarding formulaic language, lexical bundles, and the research approaches used in this area. It also provides a brief introduction of the IELTS test and its validation. Chapter three, Methodology, opens with a description of the data collection process. It then discusses the approaches to bundle extraction and categorisation of the collected data. The third section draws this chapter to a close by exploring the theoretical frameworks to analysing frequency, structural, and functional characteristics of lexical bundles. Chapter four presents the main findings and discussions of the results. Finally, Chapter five summarises the research's key conclusions and implications, followed by a consideration of limitations and suggestions for further research.

## 2. Literature review

This chapter opens with an examination of formulaic language, which offers a background to the concept of lexical bundles and its importance in language acquisition in the next section. From there, the chapter looks at the use of lexical bundles in academic environments discovered by previous research. This leads to the next area, approaches to investigating lexical bundles, with a focus on corpus analysis. The fifth section reflects on studies into the IELTS test to justify the need for conducting research in this field, and the final section finishes with a summary of the pilot study.

### 2.1. Formulaic language

With the advent of large corpora in the academic world (e.g. the Michigan Corpus of Academic Spoken English (MICASE)), more studies have focused on formulaic language to explore the linguistic patterns in academic speech. Although there are different terms for this phenomenon, formulaic sequence can be considered the most comprehensive term (Schmitt and Carter, 2004). It is defined as a “a sequence, continuous or discontinuous, of words or other elements, which is, or appears to be, prefabricated” (Wray, 2002, p.9).

The two most important features of formulaic language emerge from the above definition: it is associated with sequences of lexis, and processed in the mind as wholes (Schmitt and Carter, 2004). This understanding aligns with functional and usage-based theories of language, which assert that language is represented by constructions, or conventionalised word sequences, in a speaker’s mind that are stored as separate processing units and used frequently enough to be accessed together (Bybee and Clay, 2009). Such a perspective is also shared by Schmitt (2000, p.105), who point out “language is not constructed word by word, but key word by key word”. In

other words, it is not strictly compositional but consists of prefabricated phrases (Biber et al., 2004), and dealt with as single “big words” (Ellis, 1996, p.111), or “single choices” (Sinclair, 1991, p.110). These expressions are so pervasive that they constitute up to a half of our discourse depending on the measurement (Conklin and Schmitt, 2012).

Formulaic language is widely accepted to be highly useful in communication because they serve various referential, discursal, and social purposes (Schmitt and Schmitt, 2020). Besides framing, conveying, and elaborating ideas and meanings (ibid.), many sequences, especially fixed expressions, are used as discourse organising tools (Schmitt and Carter, 2004). In spoken discourse, they also imply speech acts such as showing politeness, making requests, and facilitate conversation (Nattinger and DeCarrico, 1992). Another notable role of these strings of words is they enhance fluency. As preconstructed phrases are stored as chunks in language users’ mind, they are often retrieved with little effort from memory, thus reducing the cognitive demands (Kuiper, 1996, and Ellis, 1996) and improving processing speed both in reading (Underwood et al., 2004) and speaking (Yorio, 1980; Dechert, 1983; Schmitt et al., 2004). Overall, contrary to Pinker (1994)’s assertion that prefabricated chunks is merely a peripheral concern that brings little contribution to the understanding of language processing, these functions confirm the importance of formulaic language in both communication and language acquisition, making it a topic worthy of examination.

Formulaic language exists in a range of forms with distinct characteristics such as collocations, phrasal verbs, phrasal expressions, idioms, and lexical bundles. Lexical bundles are extended word sequences that reoccur regularly in discourse, serving diverse important stance, referential, and discursal functions depending on the circumstances. This category is so predominant in

both conversation and written texts that it deserves more rigorous study as a class (Biber et al., 1999; Vilkaite, 2016).

## **2.2. Lexical bundles**

### **2.2.1. Definition and general characteristics**

There is little consensus on how to name and identify lexical bundles. Several terminologies have been used to address these uninterrupted word sequences: recurrent word-combinations (Altenberg, 1998), phrasicon (DeCock et al., 1998), lexical clusters (Hyland, 2008; Schmitt et al., 2004), n-grams (Banerjee & Pedersen 2003), and lexical bundles (Cortes, 2002, 2004; Biber et al., 1999). However, until recently, the term ‘lexical bundles’ appears to receive most preference as it is used in many studies in this field (e.g. Chen and Baker, 201, 2016; Neely and Cortes, 2009; Chen and Chen, 2020). Biber et al. (1999, p.990) gives the first definition of this concept: “Lexical bundles are recurrent expressions, regardless of their idiomaticity, and regardless of their structural status”. Cortes (2004, p.400) further develops this definition, highlighting that these “extended sequences” must “statistically co-occur in a register”.

Three characteristics of lexical bundles arise from the previous descriptions. First, they occur frequently in a particular register. In many studies (e.g. Biber et al., 2004; Biber and Barbieri, 2007; Chen and Chen, 2020), a sequence must occur at least 40 times per 1 million words to qualify as a lexical bundle. This characteristic distinguishes lexical bundles from other types of formulaic language. In Vilkaite (2016)’s study, while idioms account for only around 2.5% of a text, lexical bundles are the largest subset of formulaic language, covering from 14% up to 50% of a text depending on the register. Besides frequency, most researchers also set a minimum range of texts in a corpus where a sequence appears for it to be classified as a bundle. The range

thresholds vary from three to five texts to guard against idiosyncratic use of bundles in a single text (e.g. Biber et al., 1999; Cortes, 2004; Hyland, 2008). However, there are wide disparities in the frequency and range cut-off points applied in previous research. For example, in Biber et al. (1999)'s study of English conversation and academic prose, and Simpson-Vlach and Ellis (2010)'s investigation of academic discourse, lexical bundles must have a minimum frequency of 10 times/million words; meanwhile, Vo (2019) sets this threshold at 200 times/million words in her comparison between L1 and L2 writing. This difference may result from the dissimilarity in corpus sizes, research aims, and research scope of these studies.

The second characteristic of lexical bundles is idiomaticity, the extent to which a phrase is transparent in meaning. The majority of lexical bundles can be easily understood by retrieving the meaning of the component words (e.g. "I'd like to", "take a look at") (Cortes, 2004). The third characteristic is structural incompleteness. Most studies observe that the majority of lexical bundles are syntactic fragments that bridge two structural units, such as in the case of "don't want to", "the end of the" (Biber et al., 1999; Hyland, 2008; Byrd and Coxhead, 2010). Biber et al. (1999) found only 5% of the recurrent sequences in academic prose are syntactically complete. It is this third attribute of lexical bundles that reduces their salience, hence raising challenges for teaching and learning.

### **2.2.2. Structural and functional categorisation**

Lexical bundles also exhibit a multitude of structural patterns. Two structural frameworks for lexical bundles have been suggested by Biber et al. (1999) and Chen and Baker (2010), as shown in Table 1 below.

Biber et al. (1999)		Chen and Baker (2010)	
Conversation	Academic prose		
1. personal pronoun + lexical verb phrase	1. anticipatory it + verb phrase/adjective phrase	1. anticipatory it + verb phrase/adjective phrase	VP-based
2. pronoun/noun phrase + be +	2. passive verb + prepositional phrase fragment	2. passive verb + prepositional phrase fragment	
3. verb phrase with active verb	3. copula be + noun phrase/adjective phrase	3. copula be + noun phrase/adjective phrase	
4. yes-no question fragments	4. pronoun/noun phrase + be (+ . . .)	4. verb phrase with active verb	
5. wh-question fragments	5. (verb phrase +) that-clause fragment	5. (verb phrase +) that-clause fragment	
6. lexical bundles with wh-clauses	6. (verb/adjective +) to-clause fragment	6 (verb/adjective +) to-clause fragment	
7. lexical bundles with to-clauses	7. adverbial clause fragment	7. other expressions	
8. verb + that-clause fragments	8. noun phrase with of-phrase fragment	8. preposition + noun phrase fragment	PP-based
9. adverbial clause fragments	9. noun phrase with other post-modifier fragment	9. noun phrase with post-modifier fragment	NP-based
10. noun phrase expressions	10. prepositional phrase with of-phrase fragment		
11. prepositional phrase expressions	11. other prepositional phrase fragment		
12. quantifier expressions	12. other expressions.		
13. other expressions			
14. meaningless sound bundles			

*Table 1. Comparison of structural categories between Biber et al. (1999) and Chen and Baker (2010)*

Biber et al. (1999) provides the first and most comprehensive structural categorization of lexical bundles. In this study, lexical bundles are classified into 14 types of structures for conversation register, and 12 structures for academic prose, with some structures shared between them. Later, Chen and Baker (2010) simplify and reorganize Biber et al. (1999)'s taxonomy into a single system, with three main categories: verb phrase-based (VP-based) (e.g. "it's going to be"), noun phrase-based (NP-based) (e.g. "one of the things"), and prepositional phrase-based bundles (PP-based) (e.g. "at the end of"). The sub-categories are the same as in the original system although some are removed. This will be explained in more detail in the Chapter 3. Previous research shows that structural patterns of lexical bundles vary across registers and writers of different language or professional backgrounds (e.g. Biber et al., 1999, 2004; Chen and Baker, 2010; Vo, 2019). For example, spoken discourse uses far more lexical bundles than written discourse, and is characterised by the dominance of VP-based phrases over the other structures, while written registers heavily rely on noun and prepositional phrases (Biber et al., 1999, 2004).



Lexical bundles can also serve a range of discursual functions, as shown in Table 2.

	<b>Biber et al. (2004)</b>	<b>Simpson-Vlach &amp; Ellis (2010)</b>
<b>Stance bundles</b>	1. Epistemic stance 2. Attitudinal/modality stance 2.1. Desire 2.2. Obligation/directive 2.3. Intention/prediction 2.4. Ability	1. Epistemic stance 2. Obligation/directive 3. Intention/volition & prediction 4. Ability & possibility 5. Evaluation 6. Hedges
<b>Referential bundles</b>	1. Specification of attributes 1.1. Intangible framing attributes 1.2. Tangible framing attributes 1.3. Quantity 2. Identification/focus 3. Imprecision 4. Time/place/text reference 4.1. Time reference 4.2. Place reference 4.3. Text deixis 4.4. Multifunctional reference	1. Specification of attributes 1.1. Intangible framing attributes 1.2. Tangible framing attributes 1.3. Quantity 2. Identification/focus 3. Vagueness markers 4. Deictics and locatives 5. Contrast and comparison
<b>Discourse organising bundles</b>	1. Topic introduction 2. Topic elaboration/clarification	1. Topic introduction and focus 2. Topic elaboration 2.1. Non-causal 2.2. Cause and effect 3. Metadiscourse and textual reference 4. Discourse markers
<b>Special conversational functions</b>	1. Politeness 2. Simple inquiry 3. Reporting	

*Table 2. Comparison of functional categories between Biber et al. (2004) and Simpson-Vlach and Ellis (2010)*

The most widely adopted functional taxonomy so far is probably the one by Biber et al. (2004) in a study of lexical bundle use in academic language. In this model, there are four main functional categories: stance, referential, discourse organizing, and special conversation expressions. Stance expressions “express attitudes or assessments of certainty” of the speaker (Biber et al., 2004, p. 384). Referential bundles are used to identify or specify attributes of entities or the

textual context, while discourse organisers signal and link previous and upcoming discourse (ibid.). Special conversation functions include showing politeness, inquiring, and reporting (ibid.). This categorisation system was later modified by Simpson-Vlach and Ellis (2010). They merged some sub-functions together (e.g. Stance-Desire merged into Stance-Intention/Directive), eliminated special conversation functions, and added five extra sub-categories, although the three main functions (i.e. Stance, Referential, and Discourse organising function) still remain, as presented in Table 2. They argue these modifications are necessary to account for the much larger number of recurrent phrases extracted from their study, and identify useful bundles for pedagogical purposes. Hyland (2008) also developed a totally different functional classification composed of three functional groups. Research-oriented bundles are those used to structure activities and procedures, such as “at the same time”, “the purpose of the”. Text-oriented expressions facilitate text organisation and signal relationships between sections. Participant-oriented phrases assist writers to express their attitudes. This taxonomy is established to serve the investigation of lexical bundles common in research articles; hence, it may not match the objectives of research in other discourse communities such as academic conversations and lectures as in this study.

### **2.3. Lexical bundles in academic contexts**

Much work in the field of lexical bundles in academic discourse has centred on the use of these phrases in written forms, leaving spoken registers relatively under-explored. There are two major lines of lexical bundle research in academic writing. The first line contrasts and compares compositions by writers of different linguistic backgrounds, language proficiency, and professional superiority to draw useful pedagogical implications. For instance, DeCock (2000)

found non-native undergraduate students use more lexical bundles in total than their native counterparts, and each group favours disparate sets of bundles. Read & Nation (2006) and Vidakovic & Barker (2010) reported a higher number of formulaic word strings in high-performing test takers than in lower-level students. The second line of research analyses the patterns of lexical bundles used in other academic written registers such as in published research articles (Hyland, 2008), textbooks (Biber et al., 2004), course syllabi and university brochures (Biber and Barbieri, 2007), EAP textbooks (Oshima and Hogue, 2004; Williams, 2005; Wood, 2005). In general, these studies suggest there are clear variations in the use of lexical bundles across different types of writing, which indicates the need for more attention to the teaching and learning of this linguistic aspect in academic environment.

Most studies of academic spoken discourse focus on university lectures. DeCarrico and Nattinger (1988, p.94) investigates the role of lexical phrases as “macro-reorganisers” in authentic lectures on various disciplines. With the goal to providing valuable insights into teaching this language aspect for better lecture comprehension among students, they categorised lexical phrases into eight functional groups (e.g. Topic markers, Topic shifters, Summarisers, Exemplifiers, Relators). However, these phrases were identified mainly based on the researchers’ intuition about their salience, thus a number of structurally incomplete phrases were not included (Nesi and Basturkmen, 2006). Biber et al. (2004) compares lexical bundle patterns between classroom teaching, normal conversation, academic prose, and textbooks and draws some significant conclusions. Small differences in registers create a continuum from the more conversational or oral registers to the more informational or literate ones. Classroom teaching uses a much larger number of lexical bundle types than conversation. It also uses more stance and discourse

organising bundles than conversation, and also more referential bundles than academic prose. This indicates classroom teaching's reliance on both oral and literate bundles. It is argued that the complex communicative nature of lectures, which involves both impromptu, time-constrained production and informative elements, is the reason behind this attribute. Nevertheless, the researchers do not distinguish between different types of classroom teaching. Most of the speech events in their classroom teaching sub-corpus drawn from T2K SWAL are interactive like regular lessons rather than pre-planned monologic lectures (Nesi and Basturkmen, 2006). This may compromise the representativeness and homogeneity of the corpus, hence the final findings.

Biber and Barbieri (2007) is one of the few studies that explore other spoken registers besides teaching in academic contexts. They include the language used in classroom teaching, class management, office hours, study groups, and service encounters, alongside a range of written registers. They discovered that each register relies on a dissimilar collection of lexical bundles corresponding to its respective communicative needs. Classroom management and service encounters employ more lexical bundles than any other spoken registers including classroom teaching. This study also questions the observation in previous research (Biber et al., 1999, 2004) that lexical bundles are more frequent in spoken than in written discourse, stating that the frequency of this lexical device depends on speakers' communicative purposes in specific contexts. Similar to the previous study by Biber et al. (2004), the materials for the classroom teaching sub-corpus in this study include both lecture-style and interactive lessons, which again raises concern over corpus representativeness.

Given the lack of a comprehensive list of formulaic sequences for pedagogical purposes, Simpson-Vlach and Ellis (2010) examined a 4.2-million-word corpus of academic spoken and written discourse to compile the Academic Formula Lists (AFL). These lists consist of expressions of three or more words, divided into three sub-lists: the Core List, the Written AFL, and the Spoken AFL. The researchers found that most bundles in academic language are referential, which is consistent with Chen and Chen (2020)'s observation in their study of academic lectures. This, however, contrasts Biber et al. (2004)'s and Biber & Barbieri (2007)'s findings that stance expressions are the largest functional group. This may be attributed to the different functional taxonomies used in these studies, leading to different categorisation. Another possible explanation lies in the lexical bundles themselves, as while the AFL includes all three-, four -, and five-word bundles, those in Biber et al. (2004) and Biber & Barbieri (2007) are four-word bundles only. Their slightly different approaches to corpus analysis of lexical bundles are also a potential reason behind their contrasting findings, which will be explained in the coming section. Besides, Simpson-Vlach and Ellis (2010) argue there exist a number of core lexical bundles shared across multiple academic disciplines. This contradicts Hyland (2008)'s view, which says that it is impossible to compile a list of academic lexical expressions common to various disciplines.

#### **2.4. Corpus analysis of lexical bundles**

Corpus linguistics is a research approach that allows empirical investigation of multiple language aspects using statistical analysis of a large collection of texts called 'corpus' (Biber and Reppen, 2015). There are two major corpus approaches to investigating lexical bundles: the phraseological and the distributional approaches (Granger and Paquot, 2008). In the phraseological approach, recurrent expressions are extracted based on a preconstructed list of

phrases previously identified based on the researchers' intuition or other means. As mentioned earlier, this approach relies on the perceptual salience of word combinations, thus ignoring syntactically fragmented phrases that tend to slip humans' observation (Nesi and Basturkmen, 2006). Because lexical phrases are identified manually, this approach is only feasible for small-scaled research (e.g. DeCarrico and Nattinger, 1988; Nattinger & DeCarrico, 1992) (Wood, 2005). By contrast, the distributional approach determines lexical bundles automatically on the basis on their statistical frequency using concordancer software, making it suitable for research involving a large corpus (ibid.). Indeed, most influential studies in this field involve enormous corpus sizes (e.g. 8.89 million words (Chen & Chen, 2020); 4.2 million words (Simpson-Vlach & Ellis, 2010); 2.6 million words (Biber & Barbieri, 2007)).

However, the distributional approach is not without limitations. Firstly, to extract lexical bundles, selection criteria such as minimum frequency of occurrence of a bundle and range of texts where it occurs must be applied to the software. These thresholds are arbitrary as they largely depend on the research objectives, scope and the nature of the inspected corpus. Secondly, the reliability of a corpus analysis is greatly influenced by the quality of the corpus, or the extent to which it represents the target language use domain (Hunston, 2002). The reason is corpus data cannot present the language features available in the target domain but not present in the corpus (ibid.). The third drawback relates to how concordancers work in distributional studies. Simpson-Vlach & Ellis (2010) argue that since word strings are extracted based on statistical frequency alone, too many strings of undifferentiated value can be produced. This is especially true for shorter strings as the number of three-word bundles (25%) is found to far exceed that of four-word bundles (3%) (Biber et al., 1999), and there would also be far more overlaps, such as "I don't"

and “don’t know”. Finally, identifying lexical bundles based on frequency links to the fixedness of lexical bundles, which causes the undesirable elimination of some combinations which would otherwise be counted as bundle. For instance, although “I am going to” and its contracted form “I’m going to” are in effect the same structure, they are treated by concordancers as two distinct combinations, affecting their frequency counts, hence their possibility to satisfy the selection criteria of lexical bundles in the study. In such cases, manual investigation of each possible case using the “Key word in context” function is needed to account for the ignored instances.

## **2.5. The IELTS test**

### **2.5.1. Overview of the IELTS test**

The International English Testing System (IELTS), co-organised by Cambridge English, British Council, and IPD IELTS Australia, is probably the most popular assessment of learners’ English language proficiency as a prerequisite to English-medium universities. The test assesses all four communication skills of language learners and grades their competency on a banded scale from 1 (non-user) to 9 (expert user). It offers two separate testing choices, IELTS Academic for those who aim to join higher education, thus the interest of this paper, and General Training for non-academic purposes such as immigration. IELTS listening test consists of four sections, with Sections 1 and 2 representing general English and Sections 3 and 4 featuring language use in typical academic contexts. Hence, due to the limited scope, this study focuses on the latter sections alone to compare the test’s academic language use to that of actual academic contexts which often causes troubles for international students. These sections belong to different spoken registers. Section 3 is a conversation between two or three speakers about an academic issue.

Meanwhile, Section 4 is a monologic lecture by one speaker only. This difference was taken into consideration when designing the reference corpus for comparison.

### **2.5.2. Validity evaluation of the IELTS test**

Validity is the extent to which a test can measure what it is supposed to measure (McCall, 1922).

There are several types of validity, so conclusions about this aspect require examination of all of these dimensions (Weir, 2005). A commonly researched dimension is the relationship between test results and actual performance of test takers in the target domain, also known as predictive validity (Weir, 2005). While most studies report positive correlations between IELTS listening test scores and GPA of students from various backgrounds (eg. Woodrow, 2006; Yen and Kuzma, 2009; Schoepp and Garinger, 2016; Dang and Dang, 2021), some record the opposite results (e.g. Denham and Oner, 1992). However, this approach is questionable because each faculty places different values on GPA calculation, thus this score may not reliably represent academic performance of students across disciplines (Sawaki and Nissan, 2009). Construct relevance is another aspect of the IELTS test's validity explored in existing literature. It reflects the extent to which the knowledge and skills required in the assessment is relevant to those it aims to measure (Messick, 1989). Aryadoust (2012) argue the listening construct in this test was under-represented. This is because most test items primarily evaluate the ability to understand details and explicit information, whereas effective listening comprehension demands other essential subskills such as making inferencing and drawing conclusions. Regarding cognitive validity, Field (2009) maintains that the simultaneous listen-read-write format of the test's Section 4 and test-taking conditions place unreasonable pressure on students' cognitive processing. Another significant dimension of test validity is content validity, which seems largely under-researched.



Content validity is how fully the assessment “represents the full range of knowledge, skills or abilities it is intended to cover” (Green, 2014, p. 78). For a language listening test, linguistic features of the language used in its transcripts, including lexical bundles, can be one indicator of its content validity. Nevertheless, it seems no study has taken this approach to validating the IELTS listening test.

## **2.6. Summary of the Pilot study**

Prior to this study, a pilot study was conducted to determine the best corpus compilation approach for the main study. It investigated what should be the constituents of the IELTS corpus and ASE corpus. Regarding the former corpus, this study compared between Official and Non-official IELTS practice tests, while for the latter corpus, partial and full versions of the transcripts in MICASE were compared to see the extent these two pairs of corpora differ in lexical bundle use. The results suggest that there exist significant disparities between two pairs of corpora, so to maintain corpus representativeness and reliability, the IELTS corpus will consist of Official practice tests alone, and the ASE corpus will be compiled from transcripts of the entire speech events selected in MICASE.

## **3. Methodology**

This chapter presents the methodology adopted in this research. First, it clarifies the corpus compilation process, then continues with an explanation of how the lexical bundles were extracted, first using computational software, followed by manual refinement and categorisation of extracted results. The chapter closes with an overview of the theoretical frameworks employed for the data analysis.

Given the earlier gap in existing literature, this study aims to investigate the content validity of the IELTS listening test through the lens of lexical bundle use in its Sections 3 and 4. A comparative analysis between a corpus of transcripts from IELTS listening practice tests and a corpus of spoken English in academic contexts was conducted to address the research questions regarding the extent to which frequency distribution, structural and functional patterns of lexical bundles in IELTS differ from those in authentic academic spoken discourse.

### **3.1. Corpus compilation**

This study adopted a corpus-driven approach to researching lexical bundles. It is an inductive procedure involving the identification of word sequences from corpus analysis (Biber and Rippen, 2015) followed by qualitative categorisation and interpretation.

To allow for comparison between IELTS language and authentic academic language, two main corpora were compiled. The IELTS corpus consists of two sub-corpora (IELTS S3 and IELTS S4) corresponding to two sections (Sections 3 and 4) of the listening test targeted in this study. The ASE corpus is the reference corpus representing academic language use in real-life conditions. It is also comprised of two corresponding sub-corpora for comparative purposes: the Interactive ASE containing transcripts of highly interactive speech events and the Monologic ASE incorporating transcripts of highly monologic events. Two pairs of sub-corpora are required because as explained above, the two listening sections of the IELTS test differ fundamentally in registers. Section 3, with the participation of multiple speakers, demonstrates more conversational styles of spoken discourse, whereas Section 4 as monologic lectures is more informational and literate in nature. As these registers have been proven to exhibit dissimilar characteristics (Biber et al., 1999; 2004; Biber and Barbieri, 2007), this separation of corpora is

essential to ensure homogeneity of corpus compilation, a key requirement for corpus analysis (Sinclair, 2005).

The following principles were followed during the corpus design procedure. The corpora must guarantee representativeness, “including the full range of variability” in the target population as much as possible (Biber, 1993, p.243). It should also contain a balanced distribution of subsections illustrating different types in the target domain (Hunston, 2002), preserve complete documents of samples, and be homogeneous (i.e. having no major differences between the components) (Sinclair, 2005). Table 3 summarises the components of the two main corpora built based on the above principles.

	IELTS corpus		ASE corpus	
	Section 3	Section 4	Interactive	Monologic
Word count	67,790	61,134	653,061	134,893
No. of texts	85	85	48	13

*Table 3. Components of the IELTS corpus and the ASE corpus*

The IELTS corpus was compiled from the transcripts of the IELTS practice tests officially published by the test’s co-owners (i.e. Cambridge University Press, British Council, IDP)<sup>1</sup> to be as representative of the actual tests as possible. 85 texts of 600-900 words for each IELTS sub-corpus were collected through this selection process, which is also the total number of texts accessible to the researcher. This amounts to 67,790 words for IETLS S3 and 61,134 words for IELTS S4 (Table 3), which are relatively similar sizes.

---

<sup>1</sup> See Appendices 1 and 2 for the sources for the IELTS corpus.

The ASE corpus was built of transcripts selected from MICASE, a 1.7-million-word corpus of academic spoken English used in various registers such as lectures, seminars, service encounters, etc. at universities around the US (Simpson and Swales, 2001). MICASE was chosen because it comprises both interactive conversations and monologic speech in academic environments. This is a feature of interest in this study but is absent from the British Academic Spoken English (BASE) corpus, another candidate for the reference corpus. The transcripts were filtered from MICASE based on the following criteria. Figure 1 illustrates the screenshot of text selection criteria for MISCASE.

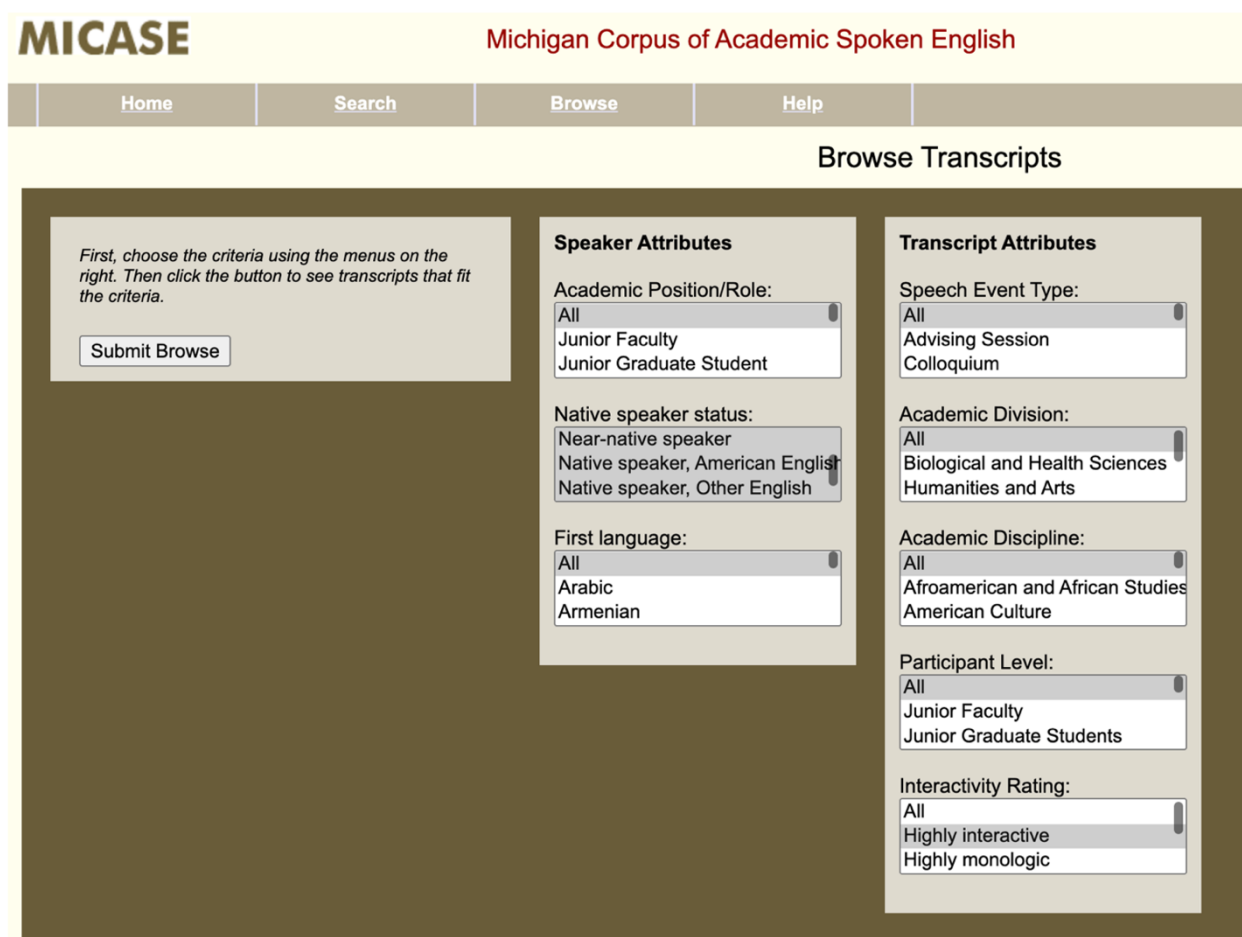


Figure 1. Transcript selection criteria for the Interactive ASE sub-corpus

As shown, the transcripts constituting Interactive ASE must have the interactivity rating of ‘highly interactive’ and be participated in by native or near-native speakers to best assimilate the standard in IELTS. It is worth clarifying that among six interactivity ratings available on MICASE’s interface, the options ‘highly interactive’ and ‘mostly interactive’ are both relevant to Interactive ASE. However, only ‘highly interactive’ transcripts were shortlisted because they contain conversations with consistently high levels of interactions throughout the events, which resembles the interactivity of IELTS’ Section 3, while ‘mostly interactive’ transcripts have occasional monologic turns. On the contrary, Monologic ASE was formed from all ‘highly monologic’ transcripts (13 texts) to match the purely monologic lecture style in IELTS’ Section 4. In fact, the BASE corpus provides a much larger number of texts that meet the above requirements for Monologic ASE (over 60 texts), but it was not chosen as it represents British English, different from the American variety used in Interactive ASE. Therefore, to maintain homogeneity, the sub-corpora of ASE contain only transcripts taken from MICASE. In summary, Interactive ASE consists of 40 texts of highly interactive speech events of different formats (e.g. advising sessions, discussions, tutorials, lab sessions), equivalent to 654,061 words. Monologic ASE is composed of 13 texts of highly monologic lectures, totalling 134,893 words (Table 3).

This corpus compilation approach has some limitations. There are significant differences in size and text lengths between the IELTS corpus and the ASE corpus, which can raise some concerns over comparability (Bestgen, 2019). Another concern is the different variety of English used between the IELTS test, which employs British English, and the reference corpus, which features American English. This dissimilarity may also somewhat affect comparability across corpora because two English varieties can syntactically and functionally differ from each other (Biber,

1987). However, given the resources available for a researcher without premium access to specialized materials and the limited scope of an MA dissertation, this approach is likely to be the optimal. Moreover, text length variations between two corpora are hardly avoidable because as an assessment, the IELTS test has to compromise between validity and practicality (Weir, 2005). In other words, each transcript can only present a 9-to-10-minute recording to maintain a practical duration for a language test, even though this length is a long way from an authentic academic lecture or discussion. To partly reduce potential contamination of the research findings caused by dissimilar English varieties, a data cleaning process was conducted at the end of this stage. Words spelled or spoken in American standards were converted to British English conventions (e.g. “-ize” changed into “-ise”; “gonna”, “wanna”, “kinda”, “coz”, “gotta” transformed into their full forms: “going to”, “want to”, “kind of”, “cause”, “got to”).

## **3.2. Lexical bundle extraction**

### **3.2.1. Computational extraction**

The n-grams function in the concordance software Antconc 4.0.6 (Anthony, 2022) was run to extract lexical bundles. Four-word sequences were targeted as they have a clearer range of structures and functions than longer bundles and are not as recurrent and numerous as shorter chunks (Biber et al., 1999). It is worth mentioning that contractions such as “’s” and “’t” were treated as one word by Antconc, so the phrase “I don’t know”, for instance, was also identified as a four-word bundle in the current study.

To be considered as a lexical bundle in this study, a sequence must occur at the rate of at least 40 times per million words and dispersed across 5 texts within a sub-corpus. This frequency cut-off point was chosen for several reasons. First, compared to the various frequency thresholds

adopted in previous research (e.g. 10 times/million words in Biber et al. (1999); 20 times/million words in Cortes (2004) and Hyland (2008); 40 times/million words in Biber et al. (2004), Biber and Barbieri (2007); 200 times/million words (Vo, 2019), this figure seems relatively conservative, thus minimizing the number of bundles selected by chance. In fact, this threshold is probably the most commonly used in pertinent literature. Second, after several trials with different frequency cut-off values, 40 was finally chosen because it allowed a number of bundles to be extracted which seems both reasonable to serve the research aims and manageable within the research scope (i.e. from 19 to 274 bundles were identified in each sub-corpus). A restriction on dispersion of phrases was applied to eliminate idiosyncratic use of any expressions in a particular text (e.g. Biber et al., 1999; 2004). The minimum dispersion rates in previous studies vary from 3 to 5 texts (Chen and Baker, 2010). Judging from the number of texts in the corpora and the degree of conservatism adopted in this study, the standard dispersion rate was set at 5.

To allow for comparison across corpora of different sizes, dynamic thresholds were adopted as discussed in Biber and Barbieri (2007) and elaborated in Chen and Baker (2010). Accordingly, the standard frequency threshold of 40 times/million words was normalized according to the corpora's sizes using the formula: *Normalised frequency threshold = (corpus' word count \* 40) / 1,000,000*. The resulting normalized figures were fractional numbers, so they were rounded up to the nearest absolute value to facilitate frequency counting. Table 4 summarises the frequency and range cut-off points for each sub-corpus.

	IELTS corpus		ASE corpus	
	Section 3	Section 4	Interactive	Monologic
Word count	67,790	61,134	653,061	134,893
No. of texts	85	85	48	13
Normalised frequency cut-off points	2.7	2.4	26.1	5.4
Rounded frequency cut-off points	3	3	27	6
Normalised range cut-off points	5	5	2.8	0.8
Rounded range cut-off points	5	5	3	2
Lexical bundle types (before refinement)	80	19	274	109

*Table 4. Normalised and rounded frequency and range cut-off points for each sub-corpus.*

As illustrated in Table 4, the normalized frequency threshold of 2.7 times for IELTS S3 was achieved by multiplying its word count of 67,790 by 40 divided by 1,000,000. This means a sequence must occur at least 2.7 times, or 3 times in effect as frequency counts must be an integer value, to be extracted as a bundle. Using this approach, the rounded frequency thresholds for IELTS S4, Interactive, and Monologic ASE are 3, 27, and 6 times respectively. The normalisation principle was also applied to dispersion rates as the numbers of texts vary significantly across corpora from 13 to 85 texts, as shown in Table 4. A base range threshold of 5 texts was set to IELTS' sub-corpora, as explained above, which was then normalised as follows: *Normalised range threshold = (corpus' number of texts \* 5) / 85*. The fractional figures were rounded up to 3 texts for Interactive ASE and 2 texts for Monologic ASE. In fact, the normalised range cut-off point for Monologic ASE was 0.8, which should have been rounded up to 1. However, a dispersion threshold of 1 text is a meaningless restriction, thus to allow some limitation on idiosyncrasies, the threshold was raised to the next integer value (i.e. 2 texts) for this sub-corpus. Nevertheless, although the normalisation approach to calculating dynamic thresholds seems intuitively logical to deal with corpus size variation, and thus is popular in many studies (e.g. Chen and Baker, 2010; Chen and Chen, 2020), some quantitative linguistics studies question its fairness in accounting



for the relationship between corpus size and distribution of word sequences (Bestgen, 2019). This will be discussed further in Section 3.4.1.

### **3.2.2. Refinement**

With the above specifications, 80, 19, 274 and 109 lexical bundles were identified from IELTS S3, IELTS S4, Interactive, and Monologic ASE sub-corpora respectively (Table 4). The bundles then underwent a manual refinement process to filter out only desirable items. The extracted lists were removed of sequences interrupted by punctuation marks (e.g. “I think. It’s”), transcribing markings (e.g. “S: Speaker information restricted”, self-repetitions (e.g. “it’s it’s”), and context-dependent phrases (e.g. “the D-N-A”, “times ten to the”) because they either do not meet the characteristics requirements of lexical bundles or are not relevant to the target language use. For simplification, contracted and complete forms of a sequence (e.g. “’re going to be” and “are going to be” were treated a single bundle type represented by the more frequent form, the total frequency of which is the sum of its two variations’ frequency.

An important decision during refinement concerns overlapping phrases. As clarified by Chen and Baker (2010), there are two types of overlaps. Complete overlaps are cases where two bundles are parts of a longer phrase and co-occur in all instances. For instance, “have a look at” and “a look at the”, each occurring twice, are parts of a longer chunk “have a look at the”, whose frequency is also 2. In this example, the shorter bundles were combined into one long phrase and counted as one bundle type with a total frequency of 2 so as not to inflate the number of different bundles and frequency counts. The second type is complete subsumption, where two phrases are overlapped in some but not all circumstances. For example, “we are going to”, whose frequency is 50, and “are going to see”, whose frequency is 10 times, occur together twice in a

five-word sequence “we are going to see”. In the remaining instances, each phrase accompanies other combinations. In this case, two partial overlapping sequences were regarded as different four-word bundles, but 2 overlapping instances were deducted from their aggregate frequency (i.e. Total frequency of two bundles = 50 + 10 - 2 = 58). Studies accounting for overlapping bundles (e.g. Chen and Baker, 2010; Chen and Chen, 2020) just distinguish between types of partially overlaps but do not elaborate on how to count the occurrences of such cases, thus this study proposes a more thorough procedure based on the belief that these complications could significantly impact the results if not properly addressed. After refinement, total lexical bundles in IELTS S3, IELTS S4, Interactive, and Monologic ASE sub-corpora are 77, 19, 140, and 132 respectively.

### **3.3. Theoretical frameworks for data analysis**

The refined lists of lexical bundles of the corpora were analysed in three aspects corresponding to three research questions: frequency distribution, structural, and functional characteristics.

#### **3.3.1. Frequency analysis**

Statistical frequency of lexical bundles was investigated from two perspectives: the number of individual bundles (types) and their total occurrences within each corpus (tokens). Both are necessary because while the former measures the range, or diversity, of lexical bundle use, the latter measures its density. A corpus can use a wider range of bundles, each occurring less frequently, while others can rely on highly frequent use of a smaller number of bundles (e.g. Biber et al., 2004). The type-token ratio, which is normally employed to measure lexical richness (Nation, 2013) based on the rate of repetitions of a word type within a set of lexis, is used in some studies as an indicator of lexical bundle density (e.g. Cooke, 2017; Hyland, 2018). However, this

ratio is “notoriously sensitive to text length so smaller corpora are likely to be more densely packed with repeated types” (Hyland, 2018). Hence, it was not used in this study where corpora vary greatly in size.

There are some divergences in the approach to comparing frequency distribution of multi-word sequences across corpora of different sizes that need to be discussed. In Vo (2019)’s study about the use of lexical features in non-native academic writing, she uses the number of bundle tokens normalised to a rate of 50,000 words instead of raw numbers to compare across corpora with different text lengths. Table 5 presents the raw and normalised numbers of bundle tokens in her study.

Corpus	Lexical Bundles: Tokens		Lexical Bundles: Types	Total Words
	Raw Frequency	Normed per 50,000		
EPT 101B	348	263	30	57,990
EPT 101C/D	851	242	31	156,790
EPT Pass	744	195	32	171,671
Total	1943	700	32	386,451

*Table 5. Frequency distribution of lexical bundles across corpora in Vo (2019)*

However, this method seems flawed. The raw numbers of bundle tokens in her study were achieved after a bundle extraction process had been run across the sub-corpora using the same standard frequency and range cut-off points. As these criteria had already been normalised to the respective corpus sizes, the raw number of tokens should not be normalised again as corpus size differences had been recognised during the extraction process. Following the above argument, bundle tokens in this study were compared using raw rather than normalised values, which is also the approach used in the majority of studies (e.g. Biber et al., 2004; Chen and Chen,

2020). Secondly, a potential limitation of this research is linked to the assertion that comparing distributions of lexical bundles across corpora of different sizes using normalised thresholds is likely to yield unreliable results (e.g. Bestgen, 2018, 2019; Cortes, 2008, 2015; Gray, 2016). Cortes (2002) found that the number of bundles is proportional to corpus size; the smaller the size, the more bundles there are. Bestgen (2019) explains that this happens due to the Zipf's law of word distribution in language: "a few words occur with very high frequency while many words occur but rarely" (Zipf, 1935, p.40-41). This disproportionality is clearer for smaller corpora (ibid.), causing more bundles to be extracted from these corpora when the same normalised frequency cut-off point is applied (Bestgen, 2019). Bestgen (2019) thus suggests that comparison should take place with corpora of similar sizes, and setting a high threshold can reduce the size effect. The first solution is not applicable to the current study due to the limited resources available to the researcher. However, compared with previous studies, the thresholds applied here (minimum frequency of 40 times/million words and minimum range of 5 texts) can be considered conservative, which can strengthen the reliability of the findings.

### **3.3.2. Structural analysis**

Most studies on grammatical structures of lexical bundles adopt the taxonomy by Biber et al. (1999), albeit with some adaptations (e.g. Chen and Baker, 2010; Vo, 2019). The current study also followed this convention as it is currently the most extensive in literature. Table 6 shows the structural categorisation of lexical bundles in the present study, compared with that of previous studies.

Biber et al. (1999)		Chen and Baker (2010)	Present study	Examples for present study
1. personal pronoun + lexical verb phrase	VP-based	1. anticipatory it + verb phrase/adjective phrase	1. (connector+) personal pronoun + lexical verb phrase	<i>I don't think</i>
2. pronoun/noun phrase + be +		2. passive verb + prepositional phrase	2. (connector+) pronoun/noun phrase	<i>that's a good</i>
3. verb phrase with active verb		3. copula be + noun phrase/ adjective	3. verb phrase with active verb	<i>have a look at</i>
4. yes-no question fragments		4. verb phrase with active verb	4. yes-no question fragments	<i>do you think it</i>
5. wh-question fragments		5. (verb phrase +) that-clause	5. wh-question fragments	<i>what did you think</i>
6. lexical bundles with wh-clauses		6 (verb/adjective +) to-clause	6. lexical bundles with wh-clauses	<i>know what I mean</i>
7. lexical bundles with to-clauses		7. other expressions	7. lexical bundles with to-clauses	<i>to be able to</i>
8. verb + that-clause fragments			8. verb + that-clause fragments	<i>don't think that</i>
9. adverbial clause fragments			9. adverbial clause fragments	<i>if you want to</i>
10. noun phrase expressions			10. copula be + noun phrase/adjective phrase	<i>be a good idea</i>
11. prepositional phrase expressions	NP-based	8. noun phrase with post-modifier fragment	11. noun phrase with of-phrase fragment	<i>the end of the</i>
12. quantifier expressions			12. Noun phrase with other post-modifier fragment	<i>the best way to</i>
13. other expressions			13. other noun phrases	<i>or something like that</i>
14. meaningless sound bundles	PP-based	9. preposition + noun phrase fragment	14. prepositional phrase with of-phrase fragment	<i>at the end of</i>
			15. other prepositional phrases	<i>in the eighteen nineties</i>
			16. others	<i>but I don't</i>

*Table 6. Structural categorisation of lexical bundles in the present study.*

As shown in Table 6, following Chen and Baker (2010)'s system, the structures were grouped into three major categories (i.e. VP-based, NP-based, and PP-based phrases) to facilitate the identification of general patterns. 12 structures were chosen from 14 grammatical structures in conversation in Biber et al. (1999), with the exclusion of "quantifier expressions" and "meaningless sound bundles" as no bundles in this study were identified to have these structures. One structure common in academic prose in Biber et al. (1999)'s taxonomy, "copula be + noun/adjective phrase", was added to this study to account for phrases such as "is a good idea". This addition is necessary as academic speech is believed to exhibit characteristics of both regular conversations and written registers (Biber et al., 2004). Unlike Chen and Baker (2010), NP-based bundles in the present study were distinguished between phrases with "-of" fragments, phrases with other post modifiers, and other phrases. The same principle also applied to PP-based expressions.

### 3.3.3. Functional analysis

This study primarily based its functional analysis on Biber et al. (2004)'s taxonomy, with the addition of some functions identified by Simpson-Vlach and Ellis (2010), to form a revised functional categorisation of lexical bundles more suited to the current investigation, which is shown in Table 7 below.

	<b>Biber et al. (2004)</b>	<b>Simpson-Vlach &amp; Ellis (2010)</b>	<b>Present study</b>	<b>Examples for present study</b>
<b>Stance bundles</b>	1. Epistemic stance	1. Epistemic stance	1. Epistemic	<i>I don't know</i>
	2. Attitudinal/modality		2. Attitudinal/Modality	
	2.1. Desire	2. Intention/volition & prediction	2.1. Intention/prediction	<i>we're going to</i>
	2.3. Intention/prediction		2.2. Obligation/directive	<i>need to think about</i>
	2.2. Obligation/directive	3. Obligation/directive	2.2. Obligation/directive	<i>and you can see</i>
	2.4. Ability	4. Ability & possibility	2.3. Ability	<i>it's a good</i>
	5. Evaluation	2.4. Evaluation		
		6. Hedges		
<b>Referential bundles</b>	1. Specification of attributes	1. Specification of attributes	1. Specification of attributes	
	1.1. Intangible framing	1.1. Intangible framing	1.1 Framing attributes	<i>in terms of the</i>
	1.2. Tangible framing	1.2. Tangible framing		
	1.3. Quantity	1.3. Quantity	1.2. Quantity	<i>a little bit more</i>
	2. Identification/focus	2. Identification/focus	2. Identification/Focus	<i>that's what you</i>
	3. Imprecision	3. Vagueness markers	3. Imprecision	<i>it's kind of</i>
	4. Time/place/text reference	4. Deictics and locatives	4. Time/Place/Text reference	<i>in the middle of</i>
	4.1. Time reference			
	4.2. Place reference			
	4.3. Text deixis			
4.4. Multifunctional				
	5. Contrast and comparison			
<b>Discourse organising bundles</b>	1. Topic introduction	1. Topic introduction and focus	1. Topic introduction	<i>going to talk about</i>
	2. Topic elaboration	2. Topic elaboration	2. Topic elaboration	<i>I mean it's</i>
		2.1. Non-causal		
		2.2. Cause and effect		
	3. Metadiscourse and textual reference			
	4. Discourse markers			
<b>Special conversational functions</b>	1. Politeness		1. Politeness	
	2. Simple inquiry		2. Simple inquiry	
	3. Reporting		3. Reporting	<i>I told you that</i>
<b>Others</b>				<i>I just don't</i>

Table 7. Functional categorisation of lexical bundles in the present study.

As demonstrated in Table 7, lexical bundles were divided into five main groups: stance expressions, referential expressions, discourse organising expressions, expressions with special conversation functions, and others, which include phrases that do not belong to any other groups. The functions served by these categories have been explained in Section 2.2.2, so this part of the paper will focus on how revisions were made to the current classification system.

The major sub-categories of stance expressions, epistemic and attitudinal/modality phrases, are preserved as in the original taxonomy. Three sub-functions under attitudinal/modality bundles, namely obligation/directive, intention/prediction, and ability bundles were also retained. The sub-function of desire bundles in Biber et al. (2004) was removed from the revised categorisation as it can be considered strongly similar to intention/prediction, making it hard to distinguish between these two functions if both remain (Simpson-Vlach and Ellis, 2010). Instead, a new sub-function identified by Simpson-Vlach and Ellis (2010), evaluation bundles, was added to the attitudinal/modality class to recognise phrases such as “it’s a good”. However, hedging expressions (e.g. “a kind of”), another new function created in their study, was not included because they considerably overlap with imprecision bundles in Biber et al. (2004)’s taxonomy.

Referential bundles and their sub-categories are mostly the same as in the original system. Minor changes were the collapse of intangible and tangible framing bundles into one single sub-function of framing attributes bundles, and time/place/text referential bundles being kept as a single sub-function. The purpose was to simplify the functional analysis, as the differentiation of these sub-categories seems insignificant to the comparison between IELTS and authentic language.

The remaining functions remain similar to the original taxonomy. Accordingly, discourse organising expressions consist of topic introducing bundles and topic elaborating bundles.

Simpson-Vlach and Ellis (2010) argue that discourse organisers should also be distinguished between *metadiscourse and textual reference* and *discourse markers*. Nevertheless, this study holds the view that such arrangements may be unnecessary, as the bundles classified under the former sub-function in their research (e.g. “I’m talking about”) can also be considered as topic introducers, and there seems a lack of justifications for most discourse organisers in their study (e.g. “and if you”) to be treated as a distinct functional group. Special conversational functions include politeness, simple inquiry and reporting bundles, all of which are absent from Simpson-Vlach and Ellis (2010)’s taxonomy.

In some cases, lexical bundles can perform multiple functions in different contexts. For instance, “what’s going on” is a referential identification bundle in “I have an idea of what’s going on”, but can be considered as a simple inquiry bundle in “alright what’s going on”. In such cases, the more frequent function was chosen to represent the general function of the bundle to make the classification process more manageable.



## 4. Results and discussions

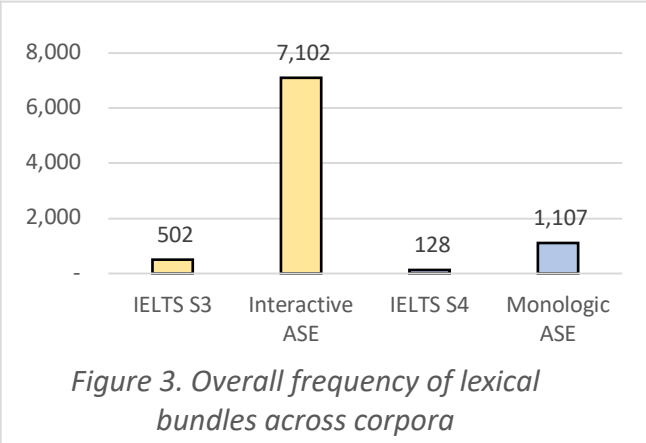
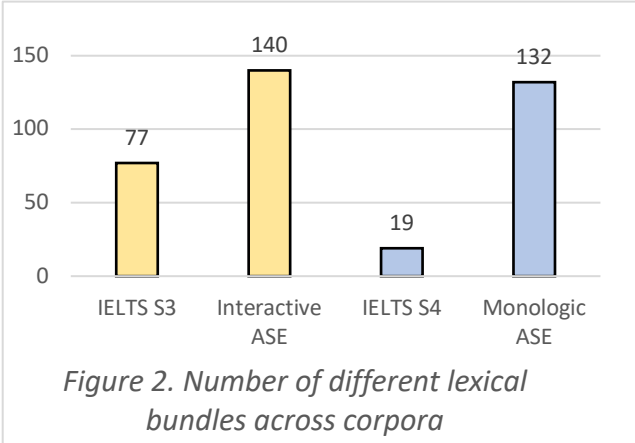
This chapter presents the main findings from the corpus analysis of lexical bundles between the IELTS corpus and the ASE corpus, together with some discussions on the results. Frequency distribution, structural, and functional patterns of bundles across corpora will be analysed sequentially in the following sections.

### 4.1. Frequency distribution of bundles across corpora

The distribution of lexical bundles in four sub-corpora is summarised in Table 8, Figure 2, and Figure 3 below.

Categories	Interactive speech		Monologic speech	
	IELTS S3	Interactive ASE	IELTS Section 4	ASE Monologic
Bundle tokens	502	7102	128	1107
Bundle types	77	140	19	132
Corpus size (words)	67,790	653,061	61,134	134,893

Table 8. Frequency distribution of lexical bundles across corpora



The data shows that regarding interactive speech, with the same standard bundle frequency and dispersion thresholds applied (as presented in Section 3.3.1), only 77 different lexical bundles were identified from IELTS Sections 3, while 140 bundles were found in Interactive ASE, almost

double the former figure. The difference is even more significant in the total occurrences of bundles, with the figure for Interactive ASE (7102 tokens) far exceeding that of IELTS S3 (502 tokens). A similar pattern is also true for monologic speech. While only 19 recurrent phrases were extracted from the 61,134-word IELTS S4 sub-corpus, nearly 7 times as many as that number of bundles were retrieved from authentic lectures. Also, there are only 128 occurrences of bundles in IELTS S4 whereas the figure for Monologic ASE is 1107 occurrences. These results indicate that IELTS academic listening sections use a smaller set of lexical bundles and use them much less frequently than real-life academic speech. It seems EAP material designers do not pay enough attention to authentic lexical bundle usage, as this finding also shares some similarities to Wood (2005)'s conclusion of the under-representation of lexical bundles in EAP textbooks. A possible explanation is that with the purpose of gauging candidates' language competency, IELTS test designers tend to avoid repetitions by paraphrasing ideas, thus reducing occurrences of fixed sequences. This, however, can unfortunately compromise the test's authenticity.

Compared to IELTS S3, IELTS S4 uses remarkably fewer lexical bundles in both type and token terms, despite coming from the same test. This is worth noticing particularly when the corresponding differences are much smaller between authentic interactive and monologic discourse. As in Table 6, Monologic ASE contains 8 bundle types fewer than Interactive ASE, while IELTS S4 have 58 bundle types fewer than its interactive counterpart. This difference may result from the fact that Section 3 still resembles natural conversations by involving interactions among speakers, whereas Section 4 sounds more like a scripted presentation with minimal informal, conversational interactions, unlike a real lecture. As lexical bundles are often more diversely and frequently used in speaking than in writing (Biber et al., 2004), such under-representation of

lexical bundles in Section 4 indicates the tendency towards written discourse of the language used here. As a test needs to reflect the nature of “the situations and text genres that candidates are likely to encounter” to be considered valid (Field, 2009, p.20), this raises doubts to the validity of Section 4 in assessing academic listening comprehension.

The 20 most frequent lexical bundles in IELTS S3 and Interactive ASE shown in Table 9 below reveal evidence which supports earlier findings.

		IELTS S3		Interactive ASE	
No.	Types	Freq.	Range	Types	Freq. Range
1	<b>I don't think</b>	25	19	<b>I don't know</b>	916 46
2	I'd like to	21	17	<b>I don't think</b>	309 45
3	what do you think	18	17	you're going to	227 40
4	<b>we're going to</b>	17	12	I'm going to	215 44
5	I'm not sure	15	15	it's going to	203 39
6	need to think about	14	10	<b>I think it's</b>	164 35
7	<b>I don't know</b>	11	10	<b>is going to be</b>	164 39
8	that's a good	11	10	you know what I	164 30
9	we have to do	11	7	<b>we're going to</b>	163 41
10	<b>don't want to</b>	10	8	I think that's	139 34
11	<b>have a look at</b>	10	8	don't know if	138 35
12	<b>I think it's</b>	10	9	that's what I	125 37
13	I've got a	10	9	you don't have	119 35
14	<b>is going to be</b>	10	9	I don't have	113 33
15	you'll need to	10	9	know what I mean	106 23
16	I think we should	9	7	I mean it's	100 29
17	is a good idea	9	8	<b>don't have to</b>	95 32
18	<b>but I don't</b>	8	6	don't know what	95 34
19	do we need to	8	7	<b>but I don't</b>	94 37
20	<b>don't have to</b>	8	8	<b>don't want to</b>	94 34

*Table 9. 20 most frequent lexical bundles in IELTS S3 and Interactive ASE*

Nearly a half of the 20 most frequent bundles are shared between IELTS S3 and Interactive ASE. This number is close to a half of the total number of bundles mutually present in both sub-corpora. These observations suggest that Section 3 of the test does capture the nature of

academic conversations to some extent in terms of the specific bundle types. Interestingly, however, the bundles in Interactive ASE are recurrent in a far wider range of texts than those in IELTS S3, despite the much smaller number of texts in the sub-corpus. The most frequent bundle in Interactive ASE is present in 46 among 48 transcripts within the corpus, while that of IELTS S3 is seen in only 19 out of 85 texts. This implies that there may be low consistency in IELTS' lexical bundle use. Looking at the dispersion of bundles within a corpus can reveal meaningful information regarding their distribution, yet previous studies seem to not give much attention to this parameter. Nevertheless, this conclusion should be treated with caution because each text in ASE is much longer than that in IELTS. The average text length in Interactive ASE is 13,600 words, while that of IELTS S3 is only 800 words. The longer an average text, the more chances a bundle can occur.

Disparities can also be seen in the discoursal functions of the bundles. In Table 7, most bundles existing only in Interactive ASE are associated with: (1) intention/prediction (e.g. "pronoun + be going to"), (2) obligation (e.g. "pronoun + don't have (+ to)"), (3) topic elaboration (e.g. "you know what I"; and (4) epistemic stance (e.g. "don't know + modifier"). This pattern is consistent with what Biber and Barbieri (2007) found in spoken university registers, with stance epistemic, stance obligation, and stance intention bundles being the most popular, followed by discourse organising bundles including topic elaborators. Meanwhile, the most frequent bundles in IELTS S3 only perform stance functions, namely: (1) intention (e.g. "I'd like to" – this is the second most frequent bundle), (2) obligation (e.g. "I think we should", "you'll need to"), and (3) evaluation (e.g. "that's a good"). These differences can be associated with the contextual characteristics of the IELTS Section 3 which often revolve around a discussion about an academic issue,

necessitating the need to give explicit intention, advice, and judgemental expressions. This will be explained in detail in Section 4.3.1. For now, some differences can be identified here. First, “I’d like to” is a more polite phrase to express intentions. Second, the obligation bundles in IELTS S3 are mainly in positive forms to suggest someone to do something, opposite to the negative structures “don’t have to” in Interactive ASE (this will be explained more clearly in Section 4.3.1). Finally, evaluation bundles are absent from the most used bundles in Interactive ASE. These findings indicate that the IELTS test favours some bundles that are not as common in authentic conversations.

The 20 most frequent bundles in IELTS S4 and Monologic ASE can be seen in Table 10 below.

IELTS S4				Monologic ASE		
No.	Types	Freq.	Range	Types	Freq.	Range
1	<b>I'm going to</b>	32	24	<b>we're going to</b>	89	13
2	I'd like to	14	13	<b>I'm going to</b>	53	10
3	am going to talk	12	11	you're going to	44	11
4	<b>going to talk about</b>	10	9	I don't know	24	6
5	let's look at	10	9	<b>going to talk about</b>	23	9
6	<b>we're going to</b>	10	10	are going to talk	23	7
7	<b>at the same time</b>	9	8	the end of the	23	8
8	one of the most	9	8	is going to be	21	9
9	today I'm going	9	9	to be able to	21	7
10	on the other hand	8	8	to make sure that	21	7
11	parts of the world	7	5	at the end of	20	7
12	a wide range of	6	5	it's going to	20	7
13	all over the world	6	5	<b>at the same time</b>	18	6
14	today we're going	6	6	are going to be	18	8
15	have been looking at	6	6	one of the things	16	5
16	going to look at	5	5	a little bit about	15	6
17	the bottom of the	5	5	has to do with	15	8
18	we've been looking	5	5	is going to be	15	8
19	will be able to	5	5	a little bit of	14	6
20				of the things that	14	5

Table 10. 20 most frequent lexical bundles in IELTS S4 and Monologic ASE

In this list, only 4 out of 19 bundles in IELTS S4 (those in bold) are shared in authentic lectures. This number is much smaller than the corresponding figure for Section 3 analysed above, suggesting a relatively weaker correspondence of Section 4 to real-life conditions. Similar to the trend in Section 3, lexical bundle use in IELTS S4 clearly shows less consistency than Monologic ASE. The most frequent bundle in this sub-corpus, “I’m going to”, is dispersed across around 30% of the texts within the corpus, while its counterpart in Monologic ASE, “we’re going to”, is present in every text. However, as mentioned earlier, it is important to bear in mind the effect of text length when interpreting these data. Looking at the specific bundles in more detail, intention/prediction stance bundles, particularly those with “be going to” (e.g. “I’m going to”, “we’re going to”), topic introducers (e.g. “let’s look at”), and referential bundles (e.g. at the same time) are mutually dominant in both sub-corpora. Nevertheless, there are some interesting dissimilarities in the exact wording and usage. Real-life lecturers tend to introduce a topic using the collective pronoun “we” as in “(we) are going to talk”, while the singular first-person pronoun “I” is more common in this combination in IELTS lectures. This is probably associated with the difference in the target audience of each speaker. As real lecturers were speaking directly to actual students in the theatre, using the inclusive “we” can help engage students in the lesson, whereas IELTS lectures are more like a presentation, causing the language to be more teacher-centred. Also, genuine lecturers frequently add “a little bit about” to topic introducing phrases containing “talk” or “tell”, as is also observed in Neely and Cortes (2009), while this bundle is absent from the entire IELTS S4. It seems that IELTS do not fully capture the informal, conversational elements of real spoken discourse. In addition, scrutiny of the concordance lines revealed that almost all instances of “pronoun + be going to” in IELTS S4 are used to signal a new

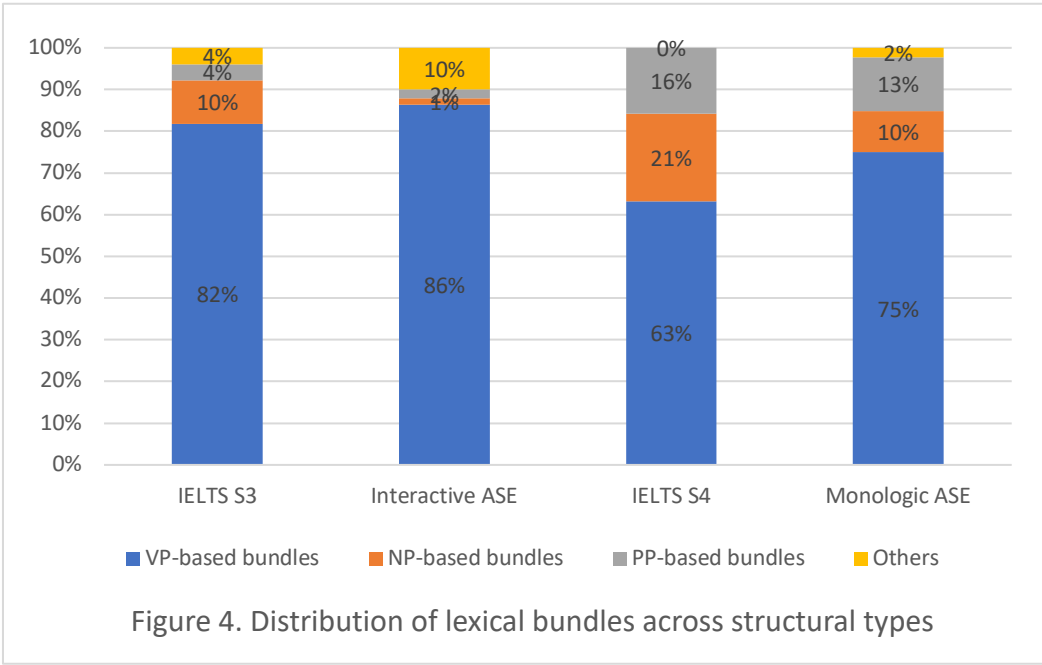
topic, while those in Monologic ASE express various kinds of intentions as well as predictions. The above findings suggest a potential lack of authentic language features and usage in the IELTS Section 4.

To sum up, the frequency analysis shows that IELTS’ academic listening sections use a narrower range of lexical bundles, with lower frequency and consistency. They are also lacking in the natural use of common phrases compared to authentic academic spoken discourse.

**4.2. Structural patterns of bundles across corpora**

**4.2.1. Between IELTS S3 and Interactive ASE**

This section analyses lexical bundles used in interactive speech corpora, namely IELTS S3 and Interactive ASE, in terms of grammatical structures. The structural distribution of lexical bundles across the sub-corpora in percentage term is illustrated in Figure 4.



From a broad perspective, it is clear that both interactive sub-corpora, IELTS S3 and Interactive ASE, employ comparably extensive proportions of VP-based bundles (e.g. “I don’t think), 82% and

86% respectively, while NP-based (e.g. “the end of the”) and PP-based bundles (e.g. “in the middle of the”) accounts for very small percentages. This distribution pattern is typical of academic spoken discourse (Biber et al., 2004) and also of general spoken discourse (Biber et al., 1999). Indeed, verb phrase bundles constitute 87% of the bundles in the conversation register in Biber et al. (1999), which is just slightly higher than the earlier figures. Despite not being directly comparable, as the data shows, it supports the point that IELTS’ Listening section 3 replicates authentic spoken discourse in academic situations in terms of the dependence on VP-based structures.

The detailed distribution of lexical bundles across all structural categories is shown in Table 11 below.



Structures	IELTS S3		Interactive ASE	
	Bundle types	Percentage	Bundle types	Percentage
<b>1. Verb phrase expressions</b>	<b>63</b>	<b>82%</b>	<b>121</b>	<b>86%</b>
1.1. (connector+) personal pronoun + lexical verb phrase <i>e.g. "I don't think"</i>	27	35%	49	35%
1.2. (connector+) pronoun/noun phrase + be <i>e.g. "that's a good"</i>	13	17%	21	15%
1.3. verb phrase with active verb <i>e.g. "have a look at"</i>	5	6%	23	16%
1.4. Yes-no question fragments <i>e.g. "do you think it"</i>	3	4%	4	3%
1.5. Wh-question fragments <i>e.g. "what did you think"</i>	4	5%	4	3%
1.6. Lexical bundles with wh-clause fragments <i>e.g. "know what I mean"</i>	0	0%	9	6%
1.7. Lexical bundles with to-clause fragments <i>e.g. "to be able to"</i>	8	10%	3	2%
1.8. (Verb) + that-clause fragments <i>e.g. "don't think that"</i>	0	0%	5	4%
1.9. Adverbial clause fragments <i>e.g. "if you want to"</i>	0	0%	3	2%
1.10. Copula be + noun phrase/adjective phrase <i>e.g. "be a good idea"</i>	3	4%	0	0%
<b>2. Noun phrase expressions</b>	<b>8</b>	<b>10%</b>	<b>2</b>	<b>1%</b>
2.1. Noun phrase with of-phrase fragment <i>e.g. "the end of the"</i>	5	6%	1	1%
2.2. Noun phrase with other post-modifier fragment <i>e.g. "the best way to"</i>	3	4%	0	0%
2.3. Other noun phrases <i>e.g. "or something like that"</i>	0	0%	1	1%
<b>3. Prepositional phrase expressions</b>	<b>3</b>	<b>4%</b>	<b>3</b>	<b>2%</b>
3.1. Preposition phrase with of-phrase fragment <i>e.g. "at the end of"</i>	2	3%	1	1%
3.2. Other prepositional phrases <i>e.g. "in the eighteen nineties"</i>	1	1%	2	1%
<b>4. Others</b>	<b>3</b>	<b>4%</b>	<b>14</b>	<b>10%</b>
<i>e.g. "but I don't"</i>				
TOTAL	77	100%	140	100%

*Table 11. Structural distribution of lexical bundles in IELTS S3 and Interactive ASE*

Examination of the distribution of bundles across the VP-based structure's subcategories, as shown in Table 11, reveals further findings. The distributions of "personal pronoun + lexical verb phrase" and "pronoun/noun phrase + be" structures in IELTS S3 are very similar to those in

authentic discourse: both being the largest structural groups, with virtually identical percentages. The proportions of two structures in IELTS S3 are 35% and 17%, while the corresponding figures in Interactive ASE are 35% and 15%. However, nearly 40% of the “personal pronoun + lexical verb phrase” bundles in IELTS S3 are also used by real-life speakers (see Appendix 5 for a full list of bundles in this sub-category). By contrast, the majority of IELTS’ “pronoun/noun phrase + be” bundles are in the form of “pronoun + be + evaluative adjective” (e.g. “that’s a good”, “that’s right it”, “it’s hard to) whereas those in Interactive ASE follow the structure of “pronoun + be + kind of/like” (e.g. “that’s kind of”, “it’s just like”), which is used as a hedging device or filler. All bundles with this structure in IELTS S3 and Interactive ASE are shown in Table 12 below, with those used for hedging and fillers highlighted in bold.

<b>IELTS S3</b>	<b>Interactive ASE</b>
<i>that's a good</i>	<i>that's what I</i>
<i>that's what I</i>	<i>that's what it</i>
<i>that's right it</i>	<i>that's what you</i>
<i>it's a good</i>	<i>that's a good</i>
<i>it's good to</i>	<i>that's why I</i>
<i>it's hard to</i>	<b><i>that's kind of</i></b>
<i>well it's a</i>	<i>okay so that's</i>
<i>there's a lot, then</i>	<i>yeah that's what</i>
<i>there's the</i>	<b><i>it's kind of</i></b>
<i>I'm not sure</i>	<i>it's not a</i>
<i>but I'm not</i>	<b><i>it's like a</i></b>
<i>I'm sure you</i>	<b><i>it's just like</i></b>
	<b><i>it's not like</i></b>
	<i>it's the same</i>
	<i>it's in the</i>
	<i>but it's not</i>
	<i>so it's not</i>
	<i>no it's not</i>
	<b><i>so it's like</i></b>
	<i>I'm not sure</i>

*Table 12. Lexical bundles with the structure “pronoun/noun phrase + be” in interactive speech*

The previous finding suggests IELTS' conversations may underuse common hedging language, a characteristic of authentic conversations. The scripted nature of IELTS speech may reduce the tendency to resort to these phrases as a means to minimize certainty and occupy the awkward pauses as often seen in instantaneous speech. The following examples illustrate how hedging phrases are used in spontaneous conversations in Interactive ASE.

1. S2: *you may want to take, um what they call sort of an intensive Latin when you're a junior or senior, uh it's kind of, it's kind of the equivalent of three semesters in a year.*  
(ADV700JU023)

2. S1: *pro- a pro- pr- a prophage is just um, like something that infects something else.*

S3: *so it's just like a generic term.* (OFC175JU145)<sup>2</sup>

Overall, these observations suggest that although there is strong congruence between IELTS' listening section 3 and actual academic conversations in the dominance of "pronoun/noun phrase + lexical verb phrase/be" structures, divergences exist in the use of specific expressions such as the overuse of "pronoun + be + evaluative adjective" bundles and the underuse of "pronoun + be + kind of/like" bundles.

Disparities can also be seen in lexical bundles with "to-clause fragments" and "wh-clause fragments". 10% (8 bundles) of the bundles in IELTS S3 falling into the former category as opposed to only 2% (3 bundles) in academic conversations. Even in Biber et al. (1999)'s general conversation, the figure is merely 5%. This may indicate IELTS' overuse of to-clause bundles in comparison with typical spoken discourse. Regarding wh-clause fragments, while there is no bundle with this structure in IELTS S3, 6% of the bundles in authentic academic interactions are

---

<sup>2</sup> The codes in brackets represent the transcripts' file name. See Appendix 3 for the full list of transcripts.

of this structural type, such as “what + I’m/you’re saying”, “what I mean like”. Biber et al. (1999) also found a quite similar percentage of lexical bundles with *wh*-clause fragments, 4%, in general conversations. It can be argued that authentic conversations use more of these phrases to reaffirm the content delivered and ensure mutual understanding. To test candidates’ ability to distinguish between false and correct information, this communicative need is also essential in the IELTS listening test, but the speakers seem to either prefer expressions different from those in genuine contexts, as in the example below, or not demonstrate it at all.

1. Robert: *Well, I guess that means a bit more work for people. I mean, they have to separate the organic and inorganic waste themselves before they take it out to the compost bin,...* (Official guide to IELTS, Test 3)

2. *but i think, it's the fact that like, maybe like when this when this um cell is like reproducing or when like this viral, stuff takes over, you know what I mean like maybe the cell doesn't do normally what it does.* (SGR175MU126)<sup>3</sup>

The shorter phrase “I mean” is often favoured in IELTS’ transcripts, while authentic conversations see a longer, more informal chunk “what I mean like”, with “like” now used frequently as a filler and hedging device by regular English speakers (Wolfson, 2022). Although opinions differ as to whether or not the overuse of “like” should be discouraged, it is still a distinctive feature of conversational discourse. Thus, a lack of such informal linguistic features in a language proficiency test can weaken the test’s authenticity as a simulation of the target language use domain.

---

<sup>3</sup> The codes in brackets represent the transcripts’ file name. See Appendices 1 and 3 for the full list of transcripts.

It is also evident from [Figure 4](#) that IELTS S3 use noticeably more NP-based bundles than authentic discourse. 10% (8 bundle types) of the bundles in IELTS S3 are noun phrases, compared to only 1% (2 bundle types) in the reference corpus. The frequent use of NP/PP-based bundles is, however, associated with written discourse (e.g. textbooks, academic prose), not conversational discourse (Biber et al., 1999; 2004). This characteristic of IELTS S3 may result from the fact that it is a scripted talk, so it may be difficult for the test designer to look from the perspective of a real speaker to incorporate subtle elements of authentic speech.

The NP-based bundles in two corpora are presented in Table 13 below.

	<b>IELTS S3</b>	<b>Interactive ASE</b>
1. Noun phrase with of-phrase fragment	<i>that a lot of</i> <i>a lot of work</i> <i>the end of the</i> <i>a bit of a</i> <i>that sort of thing</i>	<i>the end of the</i>
2. Noun phrase with other post-modifier fragment	<i>a good idea to</i> <i>a look at the</i> <i>the best way to</i>	
3. Other noun phrases		<i>or something like that</i>

*Table 13. NP-based lexical bundles in interactive sub-corpora*

It is interesting to note that nearly two-third of these phrases are not recorded in Interactive ASE as well as previous studies in lexical bundles in academic spoken discourse. For example, IELTS S3 uses “that a lot of” and “a lot of work” (Table 13), whereas common expressions with “a lot” in other lexical bundle lists are “quite a lot of” (Biber et al., 1999), “a lot of people”, “a lot of the”, “there’s a lot of” (Chen and Chen, 2020). This also challenges the authenticity of IELTS listening section 3.

#### 4.2.2. Between IELTS S4 and Monologic ASE

As shown in [Figure 4](#) above, VP-based phrases are still the largest structural group in IELTS S4, but its proportion in this sub-corpus (63%) is relatively lower than the corresponding figure in Monologic ASE (75%). Similar to the pattern observed in IELTS S3, the greatest difference in structural distribution between the IELTS test and real-life language can be seen in the use of NP-based phrases. While 21% of the bundles in IELTS Section 4 are NP-based, the figure for Monologic ASE is only a half as much, indicating an overuse of NP-based bundles. As mentioned earlier, Biber et al. (2004) maintain that the emphasis on noun phrase and prepositional phrase bundles is the key feature of written discourse, as 63% of the bundles in academic prose and over 70% of the bundles in textbooks have these structures. Hence, this general structural distribution of lexical bundles in IELTS S4 substantiates the point made in Section 4.2.1. about the presence of typical features of written discourse in IELTS' listening tests. This may be partly because Section 4 is designed to be an informative presentation, hence its reliance on NP-based bundles to make identifications, references and under-presence of interactive elements as in a real lecture with the presence of a true target audience. Also, the disparities between IELTS' and authentic language is more glaring in Section 4 than in Section 3 due to the wider gaps between the corresponding proportions (Figure 4).

Table 14 shows the breakdown of lexical bundles in IELTS S4 and Monologic ASE in terms of grammatical structures.

Structures	IELTS S4		Monologic ASE	
	Bundle types	Percentage	Bundle types	Percentage
<b>1. Verb phrase expressions</b>	<b>12</b>	<b>63%</b>	<b>99</b>	<b>75%</b>
1.1. (connector+) personal pronoun + lexical verb phrase <i>e.g. "I don't think"</i>	6	32%	41	31%
1.2. (connector+) pronoun/noun phrase + be <i>e.g. "that's a good"</i>	0	0%	14	11%
1.3. verb phrase with active verb <i>e.g. "have a look at"</i>	5	26%	19	14%
1.4. Yes-no question fragments <i>e.g. "do you think it"</i>	0	0%	0	0%
1.5. Wh-question fragments <i>e.g. "what did you think"</i>	0	0%	0	0%
1.6. Lexical bundles with wh-clause fragments <i>e.g. "know what I mean"</i>	0	0%	5	4%
1.7. Lexical bundles with to-clause fragments <i>e.g. "to be able to"</i>	1	5%	11	8%
1.8. (Verb) + that-clause fragments <i>e.g. "don't think that"</i>	0	0%	2	2%
1.9. Adverbial clause fragments <i>e.g. "if you want to"</i>	0	0%	5	4%
1.10. Copula be + noun phrase/adjective phrase <i>e.g. "be a good idea"</i>	0	0%	2	2%
<b>2. Noun phrase expressions</b>	<b>4</b>	<b>21%</b>	<b>13</b>	<b>10%</b>
2.1. Noun phrase with of-phrase fragment <i>e.g. "the end of the"</i>	4	21%	6	5%
2.2. Noun phrase with other post-modifier fragment <i>e.g. "the best way to"</i>	0	0%	6	5%
2.3. Other noun phrases <i>e.g. "or something like that"</i>	0	0%	1	1%
<b>3. Prepositional phrase expressions</b>	<b>3</b>	<b>16%</b>	<b>17</b>	<b>13%</b>
3.1. Preposition phrase with of-phrase fragment <i>e.g. "at the end of"</i>	0	0%	10	8%
3.2. Other prepositional phrases <i>e.g. "in the eighteen nineties"</i>	3	16%	7	5%
<b>4. Others</b>	<b>0</b>	<b>0%</b>	<b>3</b>	<b>2%</b>
<i>e.g. "but I don't"</i>				
TOTAL	19	100%	132	100%

Table 14. Structural distribution of lexical bundles in IELTS S4 and Monologic ASE

Within the VP-based category, “personal pronoun + lexical verb phrase” bundles represent the largest group in both sub-corpora, accounting for 63% and 75% of all bundles in IELTS S4 and

Monologic ASE respectively. Specific lexical bundles in this sub-category are shown in Table 15 below.

	<b>IELTS S4</b>	<b>Monologic ASE</b>
personal pronoun + lexical verb phrase	<i>I'm going to today I'm going we're going to today we're going we've been looking I'd like to</i>	<i>we're going to, I'm going to, you're going to, it's going to, that's going to, they're going to and we're going, and you're going, so we're going, that we're going it's not going, I'm not going to you don't have, it has to do, it has to be I don't know you want to make, so I want to, that I want to, I want to do, we don't want I mean it's, I mean that's, I think it's we'll talk about, I'll talk about, we're talking about you look at the you think about it, we need to think of course you know, you know it's I'd like to I have to say I told you that we've got a, you've got to you don't get we're looking at you're trying to</i>
pronoun/noun phrase + be	<i>N/A</i>	<i>that's one of, and that's what, this is one of, and this is the, so this is the, and this is a it's a very, and it's not, but it's not there's a lot he was able to</i>
noun phrases	<i>one of the most a wide range of the bottom of the parts of the world</i>	<i>one of the things, and one of the, um one of the a little bit of, a little bit about, a little bit more, little bit about the the end of the a whole bunch of the rest of the the best way to the ways in which, ways in which we</i>

**Table 15. Some lexical bundles in monologic sub-corpora**

As shown in Table 15, most of these bundles in IELTS S4 incorporate “be going to”, which are invariably used to introduce a new topic, such as “today I’m going”. Meanwhile, those in Monologic ASE include a wider variety of combinations (e.g. bundles with “be going to”, “want to”, “have to”). Even the bundles with “be going to” are combined with different verbs (e.g. “we’re going to” + talk/see/get/do) and used with more varied purposes such as topic introducing, stating intentions, making prediction. Interestingly, IELTS S4 employs no “pronoun/noun phrase + be” structure, while this category accounts for 11% of the bundles in



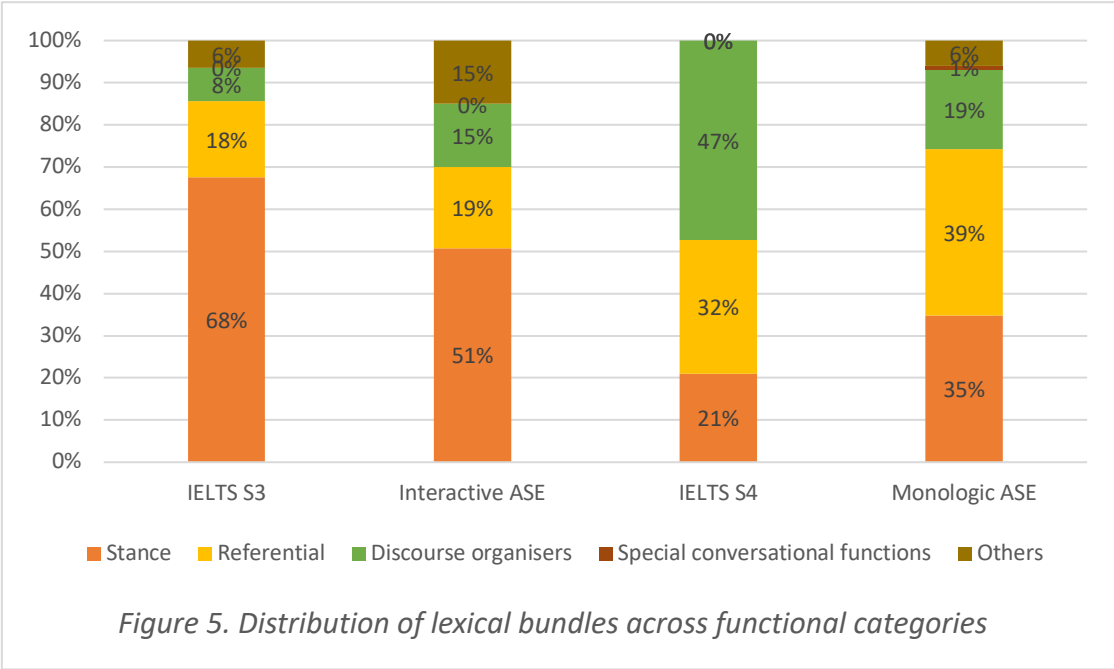
authentic lectures. Those used in Monologic ASE are consistently referential bundles such as “this is one of”, “it’s a very”, “and that’s what”, which direct references to previous discourse. This function is critical as it helps convey content, highlight key points, elaborate on preceding ideas, and guide listeners through lengthy and highly specialised lectures in real life (Chen and Chen, 2020). The analysis shows that IELTS’ listening section 4 seems deficient in a diverse use of common VP-based bundles. It also under-represents the complex nature of authentic lectures. In addition, a closer examination of NP-based bundles supports the previous assertion that IELTS’ listening section 4 overuses some features of written discourse. 3 in 4 NP-based bundles (i.e. “one of the most”, “parts of the world”, “a wide range of”) are identified by Biber et al. (1999) and Simpson-Vlach and Ellis (2010) as typical bundles in academic writing.

In general, the structural analysis reveals that although the speech in IELTS listening test resembles academic spoken language in the dominant use of VP-based bundle, there is noteworthy incongruence in the specific expressions used, especially in terms of formality and diversity, and the overuse of some written linguistic features. Section 4 shows relatively weaker conformity to authentic discourse than Section 3.

### 4.3. Functional patterns of bundles across corpora

#### 4.3.1. Between IELTS S3 and Interactive ASE

The following section analyses the functional characteristics of lexical bundles in IELTS S3 and Interactive ASE. The distribution of bundles across five main categories in four sub-corpora is represented in Figure 5 below.



Overall, similarities can be seen in the relative proportions of bundles across the major categories between two interactive sub-corpora: stance bundles are by far the most common, followed by referential bundles and discourse organising bundles. This result is not surprising as according to Biber et al. (2004), there is a strong correlation between functional and structural categories of lexical bundles. Most stance expressions are VP-based phrases (e.g. “I think that’s”), while most referential expressions are noun phrases (e.g. “that sort of thing”) or prepositional phrases (e.g. “at the end of”). As mentioned in Section 4.2.1, VP-based bundles are the largest structural group

in all sub-corpora, so it can be predicted that stance bundles also account for the majority of the bundles extracted.

The breakdown of individual functional categories is shown in Table 16, which reveals considerable differences between two corpora.

Functions	IELTS S3		Interactive ASE	
	Bundle types	Percentage	Bundle types	Percentage
<b>1. Stance expressions</b>	<b>52</b>	<b>68%</b>	<b>71</b>	<b>51%</b>
1.1. Epistemic <i>e.g. "I don't know"</i>	18	23%	30	21%
1.2. Attitudinal/Modality	34	44%	41	29%
1.2.1. Obligation/Directive <i>e.g. "need to think about"</i>	15	19%	9	6%
1.2.2. Intention/Prediction <i>e.g. "we're going to"</i>	9	12%	28	20%
1.2.3. Ability <i>e.g. "and you can see"</i>	0	0%	1	1%
1.2.4. Evaluation <i>e.g. "it's a good"</i>	10	13%	3	2%
<b>2. Referential expressions</b>	<b>14</b>	<b>18%</b>	<b>27</b>	<b>19%</b>
2.1. Identification/Focus <i>e.g. "that's what you"</i>	5	6%	16	11%
2.2. Imprecision <i>e.g. "it's kind of"</i>	1	1%	7	5%
2.3. Specification of attributes	4	5%	0	0%
2.3.1. Framing attributes <i>e.g. "in terms of the"</i>	0	0%	0	0%
2.3.2. Quantity <i>e.g. "a little bit more"</i>	4	5%	0	0%
2.4. Time/Place/Text reference <i>e.g. "in the middle of"</i>	4	5%	4	3%
<b>3. Discourse organisers</b>	<b>6</b>	<b>8%</b>	<b>21</b>	<b>15%</b>
3.1. Topic introduction/focus <i>e.g. "going to talk about"</i>	6	8%	4	3%
3.2. Topic elaboration <i>e.g. "I mean it's"</i>	0	0%	17	12%
<b>4. Special conversational functions</b>	<b>0</b>	<b>0%</b>	<b>0</b>	<b>0%</b>
4.1. Politeness <i>e.g. "thank you very much"</i>	0	0%	0	0%
4.2. Simple inquiry <i>e.g. "what are you doing"</i>	0	0%	0	0%
4.3. Reporting <i>e.g. "I said to him"</i>	0	0%	0	0%
<b>5. Others</b> <i>e.g. "I just don't"</i>	<b>5</b>	<b>6%</b>	<b>21</b>	<b>15%</b>
TOTAL	77	100%	140	100%

Table 16. Functional distribution of lexical bundles in IELTS S3 and Interactive ASE

As shown in Table 16, the largest gap (17%) is seen in stance expressions, with IELTS Section 3 having more of these bundles (68%) than authentic conversations (51%). This difference is shown to be caused by attitudinal bundles, which are used to deliver speakers' obligatory statements (e.g. "you'll need to"), intention (e.g. "we're going to"), ability (e.g. "to be able to"), and evaluation (e.g. "it's good to"). The largest attitudinal sub-group in IELTS is obligation, with 18% of total bundles, while the figure for this category in AEC is merely 6%, among the least popular functional category. The overuse of obligatory expressions may be attributed to the contextual characteristics of IELTS' conversations. IELTS' listening section 3 replicates discussions among students, or between tutors and students about an academic issue such as preparing for a group presentation, seeking advice for an assignment. They involve language functions and communicative needs such as arguing, expressing agreement or disagreement, giving suggestions to reach mutual consensus and clear actions points, hence the heavy use of obligation/directive bundles, as illustrated in the following example.

*JOHN: Yeah that's right. So we need things to measure the time and the area with, right ... what else do we need to think about? (Road to IELTS, Test 6)*

*MADDIE: Well, to compare the beaches properly we'll need to visit them all first, won't we? (Road to IELTS, Test 6)<sup>4</sup>*

In this extract from IELTS S3, two obligation bundles, "need to think about" and "we'll need to", appear within a conversation turn to facilitate the final consensus in the discussion between two students. Such situations are frequently presented in the test, which explains the abundance of obligatory expressions.

---

<sup>4</sup> The codes in brackets represent the transcripts' file name. See Appendix 1 for the full list of transcripts.

The specific obligation/directive bundles in IELTS S3 and Interactive ASE are shown in Table 17.

<b>IELTS S3</b>	<b>Interactive ASE</b>
<i>need to think about</i>	<i>don't have to</i>
<i>do we need to</i>	<i>doesn't have to</i>
<i>we'll need to</i>	<i>not have to be</i>
<i>you'll need to</i>	<i>it has to be</i>
<i>need to look at</i>	<i>going to have to</i>
<i>you don't need</i>	<i>you might want to</i>
<i>don't need to</i>	<i>you don't want</i>
<i>we have to do</i>	<i>you don't need</i>
<i>don't have to</i>	<i>don't need to</i>
<i>you'll have to</i>	
<i>I'll have to</i>	
<i>we'll have to</i>	
<i>we've got to</i>	
<i>I think we should</i>	

Table 17. Obligation/directive bundles in interactive sub-corpora

Not only do IELTS' speakers use more obligatory expressions, they also use them differently from real-life speakers. Most bundles exclusive to IELTS have the positive form (e.g. "you'll need to", "we'll have to", "I think you should"), while those in AEC are negative (e.g. don't have to, doesn't have to, not have to be), or use more polite forms such as "you might want to" (Table 17). This feature indicates that IELTS speakers tend to be more direct when instructing others. This may stem from the fact that some solutions or conclusions regarding the issue need to be reached in IELTS' conversations, creating more chances for strong, explicit suggesting expressions to be employed.

Another sub-category of attitudinal bundles witnessing considerable dissimilarities between two corpora is evaluation bundles, which are listed in Table 18.

<b>IELTS S3</b>	<b>Interactive ASE</b>
<i>that's a good</i>	<i>that's a good</i>
<i>it's a good</i>	<i>it's the same</i>
<i>is a good idea</i>	<i>it doesn't matter</i>
<i>be a good idea</i>	
<i>would be a good</i>	
<i>a good idea to</i>	
<i>that's right it</i>	
<i>yes that's right</i>	
<i>it's good to</i>	
<i>it's hard to</i>	

*Table 18. Evaluation bundles in interactive sub-corpora*

IELTS S3 uses substantially more evaluation stance expressions than Interactive ASE, 13% as compared to only 2% (Table 16). However, as shown in Table 18, only one bundle, “that’s a good”, is shared between two sub-corpora. Interestingly, most of the bundles in IELTS are associated with “a good (idea)”. Further examination of the corpora reveals that there are 22 hits of “a good idea” in IELTS S3, but only 20 hits in Interactive ASE. Judging by the much smaller size of the IELTS S3 sub-corpus, this data indicates its over-presence of evaluative expressions. IELTS speakers extensively use these phrases possibly to clarify their opinions, so that listeners can distinguish whether the speakers agree or disagree with each other in a discussion. In reality, more combinations other than “that’s a good (idea)” are used to express evaluation. As seen in the concordance lines containing “that’s a good” in Interactive ASE in Figure 6 below, some examples are “that’s a good one”, “that’s a good point”. By contrast, IELTS speakers do not use “it doesn’t matter”, an evaluative bundle found in Interactive ASE and also in academic spoken language (Simpson-Vlach and Ellis, 2010). These analyses strengthen the point made in Section 4.1 that IELTS’ listening section 3 exhibits some features uncommon to natural language use.

Total Hits: 56 Page Size 100 hits 1 to 56 of 56 hits				
	File	Left Context	Hit	Right Context
1	ADV700JU047.docx	like one–forty one–forty–one. S6: right. S2: and	that's a good	sign. S6: and then like for
2	LAB175SU026.docx	or a, meadowlark. SU–m: ooh Evening Grosbeak	that's a good	one SU–f: or a, American
3	LAB175SU026.docx	over the road. SU–f: oh nice SU–f: yeah	that's a good	one SU–f: do you want
4	LAB175SU032.docx	, [SU–f: really? ] with ten centimeter lines, [SU–f:	that's a good	idea. ] to to help make it
5	LAB175SU032.docx	nests in uh inside (Hook) Point. SU–f: mhm yep.	that's a good	place there're a few, snails
6	LAB175SU032.docx	ng about (xx) for the snorkeling? SU–f: well, (xx)	that's a good	motivating factor (for the) scuba diving. [
7	LAB500SU044.docx	re an amygdala on our (tract...) S8: alright S7: so	that's a good	slice. S6: i don't (if) (
8	LES175SU031.docx	ust mean we have to do more samples? S1: well,	that's a good	question. let's think of some
9	LES220SU140.docx	re no we're talking about like, this whole [S8:	that's a good	one ] article and it it affects
10	LES220SU140.docx	damentali– yeah ] suicide, political suicide. okay	that's a good	example yeah. S14: how about cults?
11	LES320SU085.docx	it was it was good. it was good SU–m:	that's a good	example of (xx) S5: yeah S1:
12	LES565SU137.docx	to what extent they have one yeah. yeah so no	that's a good	point. good. okay well let's
13	LES565SU137.docx	kinds of problems come up in her work. so yeah	that's a good	point Molly did you, S15: i
14	MTG400MX008.docx	i'll see you all in about three weeks. S5:	that's a good	thing S3: not me S2: not
15	MTG425JG004.docx	think we've decided yet but i think that's	that's a good	place to start. S4: yeah um,
16	MTG485SG142.docx	with no plates inside and bake it S8: right okay,	that's a good	thought. S1: so it's only
17	MTG999ST015.docx	you know i don't know S1: yeah. it that	that's a good	thing i wanted to ask this

Figure 6. Examples of combinations with “that’s a good” in Interactive ASE

On the contrary to the previous two attitudinal sub-functions, IELTS uses fewer intention/prediction bundles than authentic language. 20% of the bundles in Interactive ASE express speakers’ intentions or predictions, while the figure for IELTS is just 12% (Table 16). Table 19 below presents the bundles in this functional category found in both sub-corpora.

IELTS S3	Interactive ASE
<i>we're going to</i>	<i>you're going to , I'm going to , it's going to , we're going to , they're going to , that's going to</i>
<i>it's going to</i>	<i>to , I was going to</i>
<i>is going to be</i>	<i>is going to be , are going to be , am going to be , are going to have , are going to do</i>
<i>you're going to</i>	<i>is not going to , are not going to , you're not going , I'm not going , am not going to , not</i>
<i>don't want to</i>	<i>going to be</i>
<i>to make sure that</i>	<i>going to be like , going to be a , are you going to</i>
<i>to get hold of</i>	<i>don't want to , I don't want , do you want to , if you want to , you want to do , so you want to</i>
<i>I'd like to</i>	<i>I'm trying to</i>
<i>I'll do that</i>	

Table 19. Intention/prediction bundles in interactive sub-corpora



As seen in Table 19, bundles with the negative form “not going to” (highlighted in bold), which are relatively common in Interactive ASE (6 out of 28 intention/prediction bundles), do not appear in IELTS S3. This may explain the smaller number of intention/prediction bundles used in the test. Other studies, such as Biber et al. (1999) also found a number of “not going to” bundles in genuine conversations (e.g. “you’re/he/we not going to”, “not going to be”), which supports the point that IELTS’ lack of these bundles is problematic. Instead of “not going to”, IELTS addressers seem to prefer using phrases with “won’t” to express future predictions or intentions, as they occur at a rate of 55 times per 100,000 words in this sub-corpus as opposed to only 22 times in Interactive ASE. This difference might partly be influenced by the general inclination towards using “be (not) going to” in American English conversations (Biber et al., 1999, p.488), which is the variety used in Interactive ASE, compared to British English conversations, which the IELTS test follows. Yet, “will/won’t” is still used more frequently than “be (not) going to” in American English conversations, suggesting that the influence caused by cross-corpora dialect differences may not be significant. Another possible explanation is that IELTS speakers, who tend to use “will/won’t” which expresses higher certainty, may not be used to hedging their statements of intentions and predictions. Meanwhile, real-life speakers depend more on “be (not) going to” to avoid the face threatening force of certain intention/prediction expressions. In contrast to the pattern in Stance expressions, IELTS S3 uses fewer discourse organisers than real-life conversations, 8% compared to 15%. Not only do the bundles differ in quantity, divergences are evident in their usage. The discourse organising bundles in both corpora are presented in Table 20 below.

	<b>IELTS S3</b>	<b>Interactive ASE</b>
Topic introduction/focus	<i>let's think about to look at the have a look at let's have a a look at the let's look at</i>	<i>was going to say what do you think I have a question you're talking about</i>
Topic elaboration		<i>you know what I, you know it's, you know you're, what I'm saying, know what I'm, know what I mean, what I mean like, do you know what you see what I, see what I'm I mean it's, I mean that's, I mean I don't, I mean I think, I mean if you what do you mean does that make sense</i>

**Table 20. Discourse organising bundles in interactive sub-corpora**

Table 20 shows that all discourse organising bundles in IELTS S3 are used to introduce a new topic, opposite to those in authentic conversations, where the majority of discourse organisers are topic elaborators. It may be because real-life speakers often need to clarify and expand on their opinions to be able to fully express their thoughts in an unplanned conversation, so phrases such as “what I mean like”, “you know what I” can signal this intention and also act as hesitation markers to earn them some time for thinking. Indeed, “you know” and “I mean” are reported a typical feature of informal speech (Ostman, 1981). This is unlikely the case for IELTS speakers who have to follow a preplanned script. The planned nature of IELTS conversations also leads to the deficiency of imprecision referential bundles associated with “kind of” and “like” which are extremely common in authentic conversations.

Overall, evidence of functional analysis of lexical bundles reveals that IELTS’ listening section 3 deviates from authentic conversational discourse. On the one hand, it is more direct in expressing obligation and evaluation. Yet, it shows less caution and informality than authentic speech.

#### **4.3.2. Between IELTS S4 and Monologic ASE**

This section explores the functional pattern of lexical bundles in IELTS S4 and authentic lectures.

Table 21 below represents the breakdown of functional categories in these sub-corpora.

Functions	IELTS S4		Monologic ASE	
	Bundle types	Percentage	Bundle types	Percentage
<b>1. Stance expressions</b>	<b>4</b>	<b>21%</b>	<b>46</b>	<b>35%</b>
1.1. Epistemic <i>e.g. "I don't know"</i>	0	0%	3	2%
1.2. Attitudinal/Modality	4	21%	43	33%
1.2.1. Obligation/Directive <i>e.g. "need to think about"</i>	0	0%	5	4%
1.2.2. Intention/Prediction <i>e.g. "we're going to"</i>	3	16%	34	26%
1.2.3. Ability <i>e.g. "and you can see"</i>	1	5%	4	3%
1.2.4. Evaluation <i>e.g. "it's a good"</i>	0	0%	0	0%
<b>2. Referential expressions</b>	<b>6</b>	<b>32%</b>	<b>52</b>	<b>39%</b>
2.1. Identification/Focus <i>e.g. "that's what you"</i>	1	5%	27	20%
2.2. Imprecision <i>e.g. "it's kind of"</i>	0	0%	1	1%
2.3. Specification of attributes	1	5%	15	11%
2.3.1. Framing attributes <i>e.g. "in terms of the"</i>	0	0%	9	7%
2.3.2. Quantity <i>e.g. "a little bit more"</i>	1	5%	6	5%
2.4. Time/Place/Text reference <i>e.g. "in the middle of"</i>	4	21%	9	7%
<b>3. Discourse organisers</b>	<b>9</b>	<b>47%</b>	<b>25</b>	<b>19%</b>
3.1. Topic introduction/focus <i>e.g. "going to talk about"</i>	8	42%	13	10%
3.2. Topic elaboration <i>e.g. "I mean it's"</i>	1	5%	12	9%
<b>4. Special conversational functions</b>	<b>0</b>	<b>0%</b>	<b>1</b>	<b>1%</b>
4.1. Politeness <i>e.g. "thank you very much"</i>	0	0%	0	0%
4.2. Simple inquiry <i>e.g. "what are you doing"</i>	0	0%	0	0%
4.3. Reporting <i>e.g. "I said to him"</i>	0	0%	1	1%
<b>5. Others</b> <i>e.g. "I just don't"</i>	<b>0</b>	<b>0%</b>	<b>8</b>	<b>6%</b>
<b>TOTAL</b>	<b>19</b>	<b>100%</b>	<b>132</b>	<b>100%</b>

Table 21. Functional distribution of bundles in IELTS S4 and Monologic ASE

Compared to Section 3, there are more striking differences in the functional distribution of lexical bundles between IELTS Section 4 and authentic lectures represented by the Monologic ASE sub-corpus. Discourse organising bundles are the largest functional group in IELTS (47%), followed by referential and stance bundles, accounting for 32% and 21% of total bundles respectively (Table 21). In contrast, referential bundles (39%) and stance bundles (35%) are the largest groups in Monologic ASE, with discourse organiser being the most restricted functional group (19%). A similar distributional pattern of lexical bundles' functions is also reported by Chen and Chen (2020) in their study of academic lectures. 43%, 38%, and 19% are the respective proportions of referential, stance, and discourse organising bundles in their study, which are quite close to the figures for Monologic ASE. Since their corpus contains authentic lectures only, not other interactive teaching situations, it represents a comparable benchmark for the present study and supports the point that IELTS S4's lexical bundle use can be unnatural. In fact, IELTS S4's functional distribution of bundles does not align with the patterns found in any type of spoken or written discourse identified in previous research. This further evidence a deviation of IELTS' language use away from natural language use. However, as mentioned in Chapter 3 and Section 4.1, this result should be treated with consideration of the differences in corpus size between two sub-corpora. A closer look at the discourse organising bundles used in these corpora shown in Table 22 reveals further information.

	<b>IELTS S4</b>	<b>Monologic ASE</b>
Topic introduction/ focus	<i>am going to talk going to talk about going to look at today I'm going today we're going have been looking at we've been looking let's look at</i>	<i>going to talk about, we'll talk about, I'll talk about, we're talking about if you look at, you look at the, let's look at, we're looking at if you think about, you think about it a little bit about, talk a little bit, little bit about the</i>
Topic elaboration	<i>on the other hand</i>	<i>to do with the, have to do with, has to do with, it has to do I mean it's, I mean that's of course you know, you know it's on the other hand, so in other words let's say that, and let's say</i>

*Table 22. Discourse organising bundles in monologic sub-corpora*

Similar to Section 3, nearly all discourse organisers in IELTS are topic introducers (Table 22), which seems reasonable as topic introduction and transitions between sections in a speech are often made explicit to ensure test takers can follow the recordings. However, in real lectures, discourse organisers serve both sub-functions almost equally frequently, with several elaboration bundles shared with conversation register such as those incorporating “I mean”, “you know”. This may be due to the higher spontaneity of genuine lecturers when giving a non-scripted talk, which raises the need for self-explanation and elaboration. It may also be due to the higher level of interactivity associated with the presence of in-person target audience, which is not the case for IELTS lectures. The student-oriented nature of authentic lecture discourse is also shown in the use of transition signal bundles such as “if you look at” and “we’re looking at”, which are also identified by Biber et al. (2004) and Simpson-Vlach and Ellis (2010) as common expressions in academic speaking. As Biber et al. (2004) pointed out, using second person pronouns as in “if you look at” helps draw students’ attention to the topic and also invites their participation. Interestingly, the closest phrase found in IELTS S4 is “we’ve been looking at”, which summarises

previous discourse and at the same time notifies topic transition. It thus seems that IELTS lecturers are more explicit in signalling topic transitions. However, it lacks idea expansion signals (e.g. “I mean it’s”, “you know it’s”) and direct interactive phrases with listeners to engage them in the talk such as “if you look at” or “if you think about”.

Another functional category worth exploring is referential bundles. One of the most significant differences in referential bundles between IELTS and authentic lectures is seen in the identification/focus sub-category, which accounts for the largest proportion of bundles in Monologic ASE (20%), but is one of the smallest functional sub-groups in IELTS S4 (5%) (Table 16).

Table 23 lists all bundles having this sub-function in IELTS S4 and Monologic ASE.

<b>IELTS S4</b>	<b>Monologic ASE</b>
<i>one of the most</i>	<i>what we're going            what I want to            that I want to            one of the things            of the things that            and/um one of the            that's one of            this is one of            is one of the            it's a very            and it's not            but it's not            and this is the            so this is the            and this is a            that there's a            so there's a            if there's a            we've got a            and that's what            that's part of            what's going on            what you find is            what you get is            when it comes to            the best way to</i>

*Table 23. Identification/focus referential bundles in monologic sub-corpora*

As shown above, most of these bundles in authentic discourse are in the forms of “what + personal pronoun + verb” (e.g. “what I want to”) and “demonstrative or dummy subject + be” (e.g. “that’s one of”). The high proportion of this sub-function in real-life contexts may be because the content delivered is more complicated in nature, thus lecturers need to elaborate by giving examples and make use of more identification phrases to effectively draw listeners’ attention to the target content. This point can be illustrated in the examples below.

1. Lecturer: *okay does everybody see that? .... Okay so, what I want to explain is what these models have to do with the cannonball arrangement and sphere packings.* (COL385MU054)

2. Lecturer: *...and if you happen to have a couple of cancer cells sticking together, it will be almost impossible for them, to pass through the capillary. so this is the first place, where the cancer cells are going to have a hard time getting through the plumbing...* (LEL175SU106)<sup>5</sup>

In these examples, the lecturers use identification/focus bundles to identify their intention (Example 1) and focus students’ attention to the props required for the topic (Example 2). These phrases thus assist them to elaborate on the previous content more easily.

IELTS’ infrequency of referential bundles may also link to the effect of a real audience mentioned earlier, thus the need for frequent interactions through identification/focus referential phrases. Such evidence again undermines the IELTS listening test as a replica of genuine academic discourse in terms of audience-oriented interactions.

To conclude, the functional analysis revealed that compared to authentic academic conversations, IELTS interactive conversations are more direct in delivering obligations but less cautious and informal. Meanwhile, IELTS lectures tend to be less listener-oriented and

---

<sup>5</sup> The codes in brackets represent the transcripts’ file name. See Appendix 4 for the full list of transcripts.



interactive. Stronger divergences are found in Section 4 as opposed to Section 3, which corresponds to the previous finding in the structural analysis. These data again call into question the language used in the IELTS listening test as a replica of academic discourse, thus challenging the test's content validity.

## 5. Conclusion

This chapter summarises the main findings of this study, followed by some suggestions for pedagogical implications of these findings. It also discusses the study's potential limitations before sharing some thoughts on future research and development.

### 5.1. Summary

The purpose of this study was to investigate the content validity of the IELTS' academic listening sections through the use of lexical bundles. A corpus-driven analysis was conducted to identify the frequency distribution, structural and functional characteristics of lexical bundles in the IELTS test as compared to authentic academic spoken discourse. To do so, comparisons were made between the transcripts of IELTS Listening sections 3 and 4 and the transcripts of highly interactive and highly monologic speech events selected from the publicly available MICASE corpus. These cross-corpora comparisons revealed the following main findings.

Regarding the general frequency, IELTS uses a more restricted range of lexical bundles and use them less frequently than authentic spoken discourse. The bundles are also less dispersed across different texts in IELTS, indicating lower consistency in lexical bundle use in IELTS compared to natural language use. This can be due to the tendency to avoid repetition and diversify expressions to allow more robust assessment of test takers' language competency.

The structural analysis showed that IELTS shares some structural features of authentic speech characterised by the dominance of VP-based bundles, particularly the "pronoun + lexical verb phrase" structure. However, it overuses NP-based phrases, a feature of written discourse. Examination of the structural sub-categories points to a lack of hedging and idea elaboration but an over-presence of obligatory and evaluative expressions in IELTS.

The functional analysis revealed further divergences of IELTS from the target language use domain. Similar to authentic academic conversations, IELTS Section 3 uses stance bundles the most, followed by referential and discourse organising bundles. However, there are more stance bundles in IELTS S3 than Interactive ASE, particularly in the obligation/directive and evaluation sub-functions. Its obligation/directive expressions are also shown to be more direct and explicit. In contrast, fewer intention/prediction bundles are found in IELTS Section 3, which may be attributable to the underuse of “be going to” expressions as a hedging device typical in genuine conversations. Both listening sections also use fewer discourse organising bundles, with the absence of phrases to signal topic elaboration, indicating the possibly lower complexity of IELTS’ situational characteristics and the lower interactivity with listeners. Compared to Section 3, IELTS Section 4 shows more deviations from authentic speech due to the complete dissimilarity in functional distribution of bundles.

To conclude, IELTS listening test considerably differs from the target language use domain in lexical bundle use. These differences challenge the test’s purpose of simulating real academic speech, thus possibly weakening its content validity as an assessment of listening comprehension in academic environment.

## **5.2. Pedagogical implications**

A number of implications can be drawn from the findings of this study. First, it suggests that some modifications should be made to the IELTS listening test to increase its validity. These changes may first involve the more frequent and diverse use of lexical bundles. It should also increase the informality typical of authentic conversations by adding fillers such as “like”. More hedging phrases such as “kind of” or “be going to” should be incorporated to reflect the natural caution

of real-life speakers. Also, attention to listeners should be enhanced through the use of collective pronouns “we” instead of “I” and some audience-centred phrases such as “if you look at”. Besides, to make the discourse more authentic, evaluative and direct obligatory expressions should be reduced. Noun phrases should be restricted to prevent the speech from being too written-like. However, it is worth mentioning that some of these changes, such as using more lexical bundles, may conflict with other priorities of test design necessary for the accurate assessment of candidates’ competency, such as the diversification of lexical and grammatical resources. Thus, there should be a balance between the various criteria of an assessment, as Buck (2001) suggests.

Second, this study helps confirm the importance of lexical bundles in facilitating multiple communicative purposes (Schmitt and Schmitt, 2020) and thus provide insights into a particular discourse. Without studying them, it would be hard to realise the general tendency of discursal functions of a particular discourse, or the social relationship between speakers and addressees, as shown in this study.

The findings also substantiate the results found in previous research on lexical bundle use in academic spoken discourse such as the popularity of VP-based phrases and infrequency of NP-based and PP-based phrases in conversational registers (Biber et al., 2004), the dominance of referential bundles in lecture speech (Chen and Chen, 2020). This helps to enrich the literature in lexical bundle use in academic spoken discourse, which is a relatively under-researched area. Finally, this study demonstrates a new approach to evaluating an assessment through corpus analysis of lexical bundles. Contrary to the commonness of lexical profile analysis in assessment evaluation which looks into the use of individual words (e.g. Read and Nation, 2006; Phung and

Ha, 2022), few studies examine testing materials, especially listening tests, using multiword sequences such as lexical bundles. As shown in this study, studies of lexical bundle, together with single word investigations, can provide researchers with a more complete picture of the authenticity of the language used in a test.

### **5.3. Limitations**

The following limitations are evident in this study. First, the large difference in corpus size between the IELTS corpus (128,924 words) and the ASE corpus (787,954) can influence cross-corpora comparability. Another factor which can contaminate research results is the different English varieties employed in the IELTS corpus, which uses British English, and the ASE corpus, which uses American English. Third, most decisions during lexical bundle extraction and classification were based on the researcher's personal judgement of the researcher, which may cause disagreements. For instance, selecting lexical bundles using frequency and dispersion thresholds normalised to the different corpus sizes can be a source of controversy, as argued in Chapter 3. Many bundles can fall into multiple functional groups (e.g. "that's a good") depending how the contexts are understood by the researcher. Relying on the subjective opinion of only one rater can lead to biased results.

### **5.4. Future research and development**

Based on the limitations mentioned, it is first recommended that the analysis should be made on corpora of similar sizes to minimise possible influences on cross-corpora comparability, as also suggested by Bestgen (2019). The texts collected for corpus compilation should also have comparable text lengths and use the same English variety. Second, another rater should be involved in the lexical bundle extraction process and inter-rater reliability should be measured to

enhance the consistency of data classification. Finally, future research can look into other categories of multiword items such as collocations, phrasal expressions for a more comprehensive view of the test's language use.

## References

- Altenberg, B., 1998. On the phraseology of spoken English: the evidence of recurrent wordcombinations. In Cowie, A. P., (Ed.), *Phraseology: theory, analysis and applications*. Oxford: Oxford University Press, 1998. pp. 101–122.
- Anthony, L., 2022. *AntConc* (Version 4.0.10) [Computer Software]. Tokyo, Japan: Waseda University. Available at: <https://www.laurenceanthony.net/software>.
- Aryadoust, V., 2012. Differential Item Functioning in While-Listening Performance Tests: The Case of the International English Language Testing System (IELTS) Listening Module. *International Journal of Listening*, 26 (1), pp. 40-60. Available at: [10.1080/10904018.2012.639649](https://doi.org/10.1080/10904018.2012.639649)
- Banerjee, S. & Pedersen, T., 2003. The Design, Implementation, and Use of the Ngram Statistics Package. *Lecture Notes in Computer Science*, 2588, pp. 370–381.
- Bestgen, Y., 2018. Evaluating the frequency threshold for selecting lexical bundles by means of an extension of the Fisher's exact test. *Corpora*, 13, pp. 205–228.
- Bestgen, Y., 2019. Comparing Lexical Bundles across Corpora of Different Sizes: The Zipfian Problem. *Journal of quantitative linguistics*. 27 (3), pp. 272-290. Available at: <https://doi.org/10.1080/09296174.2019.1566975>
- Biber, D., 1987. A Textual Comparison of British and American Writing. *American Speech*, 62(2), pp. 99–119. <https://doi.org/10.2307/455273>.
- Biber, D., 1993. Representativeness in corpus design. *Literary and Linguistic Computing*, 8(4), pp. 243–257.
- Biber, D., Johansson, S., Leech, G., Conrad, S., and Finegan, E., 1999. *Longman grammar of spoken and written English*. Essex: Pearson Education Limited.
- Biber, D., Conrad, S., and Cortes, V., 2004. If you look at . . . lexical bundles in academic lectures and textbooks. *Applied Linguistics*, 25, pp. 371–405. Available at: <https://academic.oup.com/applij/article-abstract/25/3/371/179465?redirectedFrom=fulltext>

- Biber, D. and Barbieri, F., 2007. Lexical bundles in university spoken and written registers. *English for Specific Purposes*, 26 (3), pp. 263-286. Available at: <https://www.sciencedirect.com/science/article/pii/S0889490606000366>
- Biber, D., and Reppen, R. (Eds), 2015. *The Cambridge Handbook of English Corpus Linguistics*. Cambridge: Cambridge University Press.
- Buck, G., 2001. *Assessing listening*. Cambridge: Cambridge University Press.
- Bybee, J. L., and Beckner, C. 2009. Usage-Based Theory. In Heine, B., and Narrog, H. (eds), *The Oxford Handbook of Linguistic Analysis*. Oxford: Oxford University Press, pp. 827-856.
- Byrd, P. and Coxhead, A., 2010. On the other hand: Lexical bundles in academic writing and in the teaching of EAP. *University of Sydney Papers in TESOL*, 5, pp. 31-64.
- Chaudron, C. and Richards, J., 1986. The effect of discourse markers on the comprehension of lectures. *Applied Linguistics*, 7 (2), pp. 113-127. Available at: <https://academic.oup.com/applij/article-abstract/7/2/113/163715?redirectedFrom=PDF>
- Chen, Y. H. and Baker, P., 2010. Lexical bundles in L1 and L2 academic writing. *Language Learning & Technology*, 14 (2) (2010), pp. 30-49. Available at: [https://scholarspace.manoa.hawaii.edu/bitstream/10125/44213/1/14\\_02\\_chenbaker.pdf](https://scholarspace.manoa.hawaii.edu/bitstream/10125/44213/1/14_02_chenbaker.pdf)
- Chen Y. H. and Baker, P., 2016. Investigating Criterial Discourse Features across Second Language Development: Lexical Bundles in Rated Learner Essays, CEFR B1, B2 and C1. *Applied Linguistics*, 37 (6), pp.849-880. Available at: <https://doi.org/10.1093/applin/amu065>
- Chen, Y. L. and Chen H. H., 2020. Analyzing the functions of lexical bundles in undergraduate academic lectures for pedagogical use. *English for specific purposes*, 58, pp. 127-137. Available at: <https://doi.org/10.1016/j.esp.2019.12.003>
- Conklin, K., and Schmitt, N., 2012. The processing of formulaic language. *Annual Review of Applied Linguistics*, 32, pp. 45–61.
- Cooke, V. J., 2017. *The Use of Lexical Bundles in Korean University Students' EFL Argumentative Writing: Learner Corpus Analysis* [MA TESOL dissertation, Nottingham Trent University]. MA TESOL.



- Cooper, T., 2013. Can IELTS writing scores predict university performance? Comparing the use of lexical bundles in IELTS writing tests and first-year academic writing. *Linguistics Plus*, Vol. 42, pp. 63-79. Available at: <https://journals.co.za/doi/abs/10.10520/EJC148752>
- Cortes, V., 2002. Lexical bundles in Freshman composition. In R. Reppen, S. M. Fitzmaurice & D. Biber (Eds.), *Using corpora to explore linguistic variation*. Amsterdam: John Benjamins Publishing Company, pp. 131–145.
- Cortes, V., 2004. Lexical bundles in published and student disciplinary writing: Examples from history and biology. *English for Specific Purposes*, 23, pp. 397–423.
- Cortes, V., 2008. A comparative analysis of lexical bundles in academic history writing in English and Spanish. *Corpora*, 3, pp. 43–57.
- Cortes, V., 2015. Situating lexical bundles in the formulaic language spectrum: Origins and functional analysis developments. In Cortes, V., and Csomay, E., (Eds.), *Corpus-based research in applied linguistics*. Amsterdam: John Benjamins, pp. 197–216.
- Dang, C. N. and Dang, T. N. Y., 2021. The Predictive Validity of the IELTS Test and Contribution of IELTS Preparation Courses to International Students' Subsequent Academic Study: Insights from Vietnamese International Students in the UK. *RELC Journal*, February 2021. doi:[10.1177/0033688220985533](https://doi.org/10.1177/0033688220985533).
- DeCarrico, J. and Nattinger, J., 1988. Lexical phrases for the comprehension of academic lectures. *English for specific purposes*, 7(2), 91-102. Available at: [https://doi.org/10.1016/0889-4906\(88\)90027-0](https://doi.org/10.1016/0889-4906(88)90027-0).
- Dechert, H., 1983. How a story is done in second language. In Faerch, C., and Kasper, G. (eds), *Strategies in interlanguage communication*. London: Longman, pp. 175-195.
- DeCock, S., 1998. A recurrent word combination approach to the study of formulae in the speech of native and non-native speakers of English. *International Journal of Corpus Linguistics*, 3, pp. 59–80.

- DeCock, S., 2000. Repetitive phrasal chunkiness and advanced EFL speech and writing. In Mair, C., Hundt, M., (Eds.), *Corpus Linguistics and Linguistic Theory*. Amsterdam: Rodopi B.V., pp. 51-68.
- Denham, P. A., and Oner, J. A., 1992. *IELTS Research Project. Validation study of listening sub-test* (IDP/IELTS commissioned report). Canberra, Australia: University of Canberra.
- Ellis, N. C., 1996. Sequencing in SLA: phonological memory, chunking and points of order. *Studies in Second Language Acquisition*. 18, pp. 91–126.
- Field, J., 2009. A cognitive validation of the lecture-listening component of the IELTS listening paper. In *The IELTS Partners: British Council, IDP: IELTS Australia and Cambridge English Language Assessment* (eds.), *IELTS research reports Vol. 9*. Available at: <https://www.ielts.org/for-researchers/research-reports/volume-09-report-1>
- Flowerdew, J., 1992. Student perceptions, problems and strategies in second language lecture comprehension. *RELC Journal*, 23 (2), pp. 60-80.
- Flowerdew, J., 1995. Research of relevance to second language comprehension – an overview. In J. Flowerdew (ed.), *Academic listening: Research perspectives*. Cambridge: Cambridge University Press, 1995, pp. 7-30.
- Granger, S., and Paquot, M., 2008. Disentangling the phraseological web. In Granger, S., and Meunier, F., (Eds.), *Phraseology: An interdisciplinary perspective*. Amsterdam: John Benjamins, pp. 27–50.
- Gray, B., 2016. Lexical Bundles. In Baker, P., and Egbert, J., (Eds.), *Triangulating methodological approaches in corpus linguistic research*. New York, NY: Routledge, pp. 33–55.
- Green, A., 2014. *Exploring language assessment and testing: Language in action*. New York: Taylor & Francis Group.
- Hunston, S., 2002. *Corpora in Applied Linguistics*. Cambridge: Cambridge University Press.

- Hyland, K., 2008. As can be seen: Lexical bundles and disciplinary variation. *English for Specific Purposes*, 27 (1), pp. 4-21. Available at:  
<https://www.sciencedirect.com/science/article/pii/S0889490607000233>
- Hyland, K., 2018. Academic lexical bundles How are they changing?. *International Journal of Corpus Linguistics*, 23 (4), pp. 383–407.
- James, K., 1977. Note-taking in lectures: problems and strategies. In Cowie, A. P., and Heaton, J. B., (eds.). *English for academic purposes*. UK: British Association of Applied Linguistics/SELMOUS, 1977, pp. 89-98.
- Kuiper, K., 1996. *Smooth talkers*. Hillsdale, NJ: Lawrence Erlbaum.
- McCall, W. A., 1922. *How to Measure in Education*. New York, NY: Macmillan.
- Mendelsohn, D., 2002. The lecture buddy project: An experiment in EAP listening comprehension. *TESL Canada Journal*, 20 (1), pp. 64-73.
- Messick, S., 1989. Validity. In Linn R. L., (ed.), *Educational Measurement*, 3<sup>rd</sup> ed. New York, NY: Macmillan, 1989, pp. 13-103.
- Morley, J., 2001. Aural comprehension instruction: Principles and practices. In Celce-Murcia, M. (ed.), *Teaching English as a second or foreign language*. 3<sup>rd</sup> ed. Boston: Heinle & Heinle, pp. 69-85.
- Nation, I. S. P., 2013. *Learning vocabulary in another language*. Cambridge: Cambridge University Press.
- Nattinger, J. R., and DeCarrico, J. S., 1992. *Lexical phrases and language teaching*. Oxford: Oxford University Press.
- Neely, E. and Cortes, V., 2009. A little bit about: Analyzing and teaching lexical bundles in academic lectures. *Language Value*, 1 (1), pp. 17-38.
- Nesi, H. and Basturkmen, H., 2006. Lexical bundles and discourse signalling in academic lectures. *International Journal of Corpus Linguistics*, 11 (3), pp. 283-304. Available at:  
<https://doi.org/10.1075/ijcl.11.3.04nes>

- Olsen, L. A., and Huckin, T. N., 1990. Point-driven understanding in engineering lecture comprehension. *English for specific purposes*, 9, pp. 33-47.
- Oshima, A., and Hogue, A., 2004. *Writing academic English*, 4th ed. Montreal, QC: Pearson.
- Ostman, J., 1981. *You know: A discourse functional approach*. Amsterdam: John Benjamins.
- Phung, D. H., and Ha, H. T., 2022. Vocabulary Demands of the IELTS Listening Test: An In-Depth Analysis. *SAGE Open*, 12(1). <https://doi.org/10.1177/21582440221079934>.
- Pinker, S., 1994. *The language instinct*. Harmondsworth: Penguin.
- Read, J. and Nation, P., 2006. An investigation of the lexical dimension of the IELTS speaking test. In the IELTS Partners: British Council, IDP: IELTS Australia and Cambridge English Language Assessment, (eds.), *IELTS research reports Vol 6*, 2006, pp. 207-231. Available at: [https://www.ielts.org/-/media/research-reports/ielts\\_rr\\_volume06\\_report7.ashx](https://www.ielts.org/-/media/research-reports/ielts_rr_volume06_report7.ashx)
- Sawaki, Y. and Nissan, S., 2009. Criterion-related validity of the toefl ibt listening section. In ETS, *ETS Research Report Series*, 2009, pp. i-82. Available at: <https://doi.org/10.1002/j.2333-8504.2009.tb02159.x>.
- Schmitt, N. 2000. *Vocabulary in Language Teaching*. Cambridge: Cambridge University Press.
- Schmitt, N. and Carter, 2004. Formulaic sequences in action: An introduction. In Schmitt, N. (ed.), *Formulaic Sequences: Acquisition, Processing And Use*. Amsterdam: John Benjamins B.V., pp. 1-22.
- Schmitt, N., Grandage, S., and Adolphs, S., 2004. Are corpus-driven recurrent clusters psycholinguistically valid?. In Schmitt, N. (ed.), *Formulaic Sequences: Acquisition, Processing And Use*. Amsterdam: John Benjamins B.V.
- Schmitt, N. and Schmitt, D., 2020. *Vocabulary in Language Teaching*. Cambridge: Cambridge University Press.
- Schoepp, K., and Garinger, D., 2016. IELTS and academic success in higher education: A UAE perspective. *International Journal of Applied Linguistics and English Literature*, 5(3), pp. 145–151.

Simpson, R. and Swales, J., 2001. *Corpus Linguistics in North America*. Ann Arbor: University of Michigan Press.

Simpson-Vlach, R. and Ellis, N. C., 2010. An Academic Formulas List: New Methods in Phraseology Research, *Applied Linguistics*, 31 (4), pp. 489-512. Available at: <https://doi.org/10.1093/applin/amp058>.

Sinclair, J.McH, 1991. *Corpus, concordance, collocation*. Oxford: Oxford University Press.

Sinclair, J., 2005. Corpus and text: Basic principles. In Wynne M., (ed.), *Developing Linguistic Corpora: A guide to good practice*. Oxford: Oxbow Books, pp. 1–16.

Underwood, G., Schmitt, N., and Galpin, A., 2004. The eyes have it: An eye movement study into the processing of formulaic sequences. In Schmitt, N. (ed.), *Formulaic Sequences: Acquisition, Processing And Use*. Amsterdam: John Benjamins B.V.

Vidakovic, I., Barker, F., 2010. Use of words and multi-word units in Skills for Life writing examinations. In IELTS Australia/British Council, *IELTS research reports 41*, pp. 7-14.

Vilkaitė, L., 2016. Formulaic language is not all the same: comparing the frequency of idiomatic phrases, collocations, lexical bundles, and phrasal verbs. *Taikomoji kalbotyra*, 8, pp. 28-54. Available at: <https://doi.org/10.15388/TK.2016.17505>.

Vo, S., 2019. Use of lexical features in non-native academic writing. *Journal of second language writing*, 44, pp. 1-12. Available at: <https://doi.org/10.1016/j.jslw.2018.11.002>

Weir, C. J., 2005. *Language testing and validation: An evidence-based approach*. New York: Palgrave MacMillan.

Williams, J., 2005. *Learning English for academic purposes*. Montreal, QC: ERPI.

Wolfson, S., 2022. Why do people, like, say, 'like' so much?. *The Guardian*. Available at: <https://www.theguardian.com/science/2022/may/15/why-do-people-like-say-like-so-much-in-praise-of-an-underappreciated-word>.

Wood, D., 2005. Lexical Clusters in an EAP Textbook Corpus. In Wood, D., (ed.). *Perspectives on Formulaic Language Acquisition and Communication*. London: Continuum International Publishing Group, 2005, pp. 88-106.

Woodrow, L., 2006. Academic success of international postgraduate education students and the role of English proficiency. *University of Sydney Papers in TESOL*. 1(1), pp. 51–70.

Wray, A., 2002. *Formulaic language and the lexicon*. Cambridge: Cambridge University Press.

Yen, D., and Kuzma, J., 2009. Higher IELTS score, higher academic performance? The validity of IELTS in predicting the academic performance of Chinese students. *Worcester Journal of Learning and Teaching*, 3, pp. 1–7.

Yorio, C., 1980. Conventionalised language forms and the development of communicative competence. *TESOL Quarterly*, 14 (4), pp. 433-442.

Zipf, G. K., 1935. *The psycho-biology of language*. Cambridge, MA: MIT Press.

## Appendices

### Appendix 1. Sources for the IELTS Section 3 sub-corpus

Sources	No. of texts
<i>Total</i>	<i>85</i>
IELTS Academic 1 (Cambridge English, 2015)	4
IELTS Academic 2 (Cambridge English, 2015)	4
IELTS Academic 3 (Cambridge English, 2015)	4
IELTS Academic 4 (Cambridge English, 2015)	4
IELTS Academic 5 (Cambridge English, 2015)	4
IELTS Academic 6 (Cambridge English, 2015)	4
IELTS Academic 7 (Cambridge English, 2015)	4
IELTS Academic 8 (Cambridge English, 2015)	4
IELTS Academic 9 (Cambridge English, 2015)	4
IELTS Academic 10 (Cambridge English, 2015)	4
IELTS Academic 11 (Cambridge English, 2016)	4
IELTS Academic 12 (Cambridge English, 2017)	4
IELTS Academic 13 (Cambridge English, 2018)	4
IELTS Academic 14 (Cambridge English, 2019)	4
IELTS Academic 15 (Cambridge English, 2020)	4
IELTS Academic 16 (Cambridge English, 2021)	4
IELTS Academic 17 (Cambridge English, 2021)	4
The Official Cambridge guide to IELTS for Academic and General training (Cambridge English, 2014)	8
Road to IELTS (British Council, 2012)	6
Official IELTS practice materials 2009 (ielts.org, 2009)	1
Official IELTS practice materials 2010 (ielts.org, 2010)	1
IELTS Practice Materials Academic Test (British Council, 2015)	1

## Appendix 2. Sources for the IELTS Section 4 sub-corpus

Sources	No. of texts
<i>Total</i>	<i>85</i>
IELTS Academic 1 (Cambridge English, 2015)	4
IELTS Academic 2 (Cambridge English, 2015)	4
IELTS Academic 3 (Cambridge English, 2015)	4
IELTS Academic 4 (Cambridge English, 2015)	4
IELTS Academic 5 (Cambridge English, 2015)	4
IELTS Academic 6 (Cambridge English, 2015)	4
IELTS Academic 7 (Cambridge English, 2015)	4
IELTS Academic 8 (Cambridge English, 2015)	4
IELTS Academic 9 (Cambridge English, 2015)	4
IELTS Academic 10 (Cambridge English, 2015)	4
IELTS Academic 11 (Cambridge English, 2016)	4
IELTS Academic 12 (Cambridge English, 2017)	4
IELTS Academic 13 (Cambridge English, 2018)	4
IELTS Academic 14 (Cambridge English, 2019)	4
IELTS Academic 15 (Cambridge English, 2020)	4
IELTS Academic 16 (Cambridge English, 2021)	4
IELTS Academic 17 (Cambridge English, 2021)	4
The Official Cambridge guide to IELTS for Academic and General training (Cambridge English, 2014)	8
Road to IELTS (British Council, 2012)	6
Official IELTS practice materials 2009 (ielts.org, 2009)	1
Official IELTS practice materials 2010 (ielts.org, 2010)	1
IELTS Practice Materials Academic Test (British Council, 2015)	1



### Appendix 3. Sources for the Interactive ASE sub-corpus

No.	Source	Link	Event name	Recording length	Word count
1	MICASE	<a href="#">ADV700JU023</a>	<i>Honors Advising</i>	52 min.	<b>9519</b>
2	MICASE	<a href="#">ADV700JU047</a>	<i>Academic Advising</i>	124 min.	<b>28160</b>
3	MICASE	<a href="#">DIS175JU081</a>	<i>Intro Biology Discussion Section</i>	59 min.	<b>7791</b>
4	MICASE	<a href="#">DIS495JU119</a>	<i>Intro to American Politics Discussion Section</i>	55 min.	<b>7751</b>
5	MICASE	<a href="#">INT425JG001</a>	<i>Graduate Student Research Interview 1</i>	34 min.	<b>5168</b>
6	MICASE	<a href="#">INT425JG002</a>	<i>Graduate Student Research Interview 2</i>	20 min.	<b>2963</b>
7	MICASE	<a href="#">INT175SF003</a>	<i>Interview with Botanist</i>	31 min.	<b>5159</b>
8	MICASE	<a href="#">LAB200JU018</a>	<i>Chemistry Lab</i>	47 min.	<b>8169</b>
9	MICASE	<a href="#">LAB175SU026</a>	<i>Biology of Birds Field Lab</i>	92 min.	<b>11769</b>
10	MICASE	<a href="#">LAB175SU032</a>	<i>Biology of Fishes Field Lab</i>	89 min.	<b>11370</b>
11	MICASE	<a href="#">LAB175SU033</a>	<i>Biology of Fishes Lab</i>	95 min.	<b>8153</b>
12	MICASE	<a href="#">LAB500SU044</a>	<i>Biopsychology Lab</i>	52 min.	<b>9455</b>
13	MICASE	<a href="#">LES385SU007</a>	<i>Number Theory Math Lecture</i>	36 min.	<b>4144</b>
14	MICASE	<a href="#">LES175SU031</a>	<i>Biology of Fishes Group Activity</i>	19 min.	<b>2866</b>
15	MICASE	<a href="#">LES215MU056</a>	<i>Intro Latin Lecture</i>	50 min.	<b>5883</b>
16	MICASE	<a href="#">LES320SU085</a>	<i>Visual Sources Lecture</i>	69 min.	<b>12526</b>
17	MICASE	<a href="#">LES565SU137</a>	<i>Sex, Gender, and the Body Lecture</i>	73 min.	<b>14629</b>
18	MICASE	<a href="#">LES220SU140</a>	<i>Ethics Issues in Journalism Lecture</i>	83 min.	<b>16291</b>
19	MICASE	<a href="#">MTG425JG004</a>	<i>Natural Resources Research Group Meeting</i>	83 min.	<b>9382</b>
20	MICASE	<a href="#">MTG400MX008</a>	<i>Immunology Lab Meeting</i>	60 min.	<b>9523</b>
21	MICASE	<a href="#">MTG999ST015</a>	<i>Forum for International Educators Meeting</i>	102 min.	<b>17323</b>
22	MICASE	<a href="#">MTG485SG142</a>	<i>Physics Research Group Meeting</i>	41 min.	<b>9076</b>
23	MICASE	<a href="#">OFC301MU021</a>	<i>English Composition Tutorial</i>	45 min.	<b>3586</b>
24	MICASE	<a href="#">OFC578SG037</a>	<i>Technical Communications Tutorial</i>	25 min.	<b>4178</b>
25	MICASE	<a href="#">OFC150MU042</a>	<i>Astronomy Peer Tutorial</i>	102 min.	<b>21798</b>
26	MICASE	<a href="#">OFC575MU046</a>	<i>Statistics Office Hours</i>	52 min.	<b>11265</b>
27	MICASE	<a href="#">OFC270MG048</a>	<i>Computer Science Office Hours</i> <i>Anthropology of American Cities Office Hours</i>	116 min.	<b>19977</b>
28	MICASE	<a href="#">OFC115SU060</a>	<i>Anthropology of American Cities Office Hours</i>	178 min.	<b>31268</b>
29	MICASE	<a href="#">OFC105SU068</a>	<i>American Culture Advising</i>	42 min.	<b>8511</b>
30	MICASE	<a href="#">OFC355SU094</a>	<i>Linguistics Independent Study Advising</i>	52 min.	<b>6943</b>
31	MICASE	<a href="#">OFC280SU109</a>	<i>Economics Office Hours</i>	92 min.	<b>14050</b>
32	MICASE	<a href="#">OFC195SU116</a>	<i>Heat and Mass Transfer Office Hours</i>	137 min.	<b>20603</b>
33	MICASE	<a href="#">OFC175JU145</a>	<i>Intro Biology Exam Review</i>	55 min.	<b>9014</b>
34	MICASE	<a href="#">OFC320SU153</a>	<i>Art History Office Hours</i>	66 min.	<b>9233</b>
35	MICASE	<a href="#">SEM475JU084</a>	<i>First Year Philosophy Seminar</i>	72 min.	<b>13906</b>
36	MICASE	<a href="#">SEM300MU100</a>	<i>English Composition Seminar</i>	125 min.	<b>21442</b>

37	MICASE	<a href="#">SGR385SU057</a>	<i>Math Study Group</i>	132 min.	<b>17753</b>
38	MICASE	<a href="#">SGR999MX115</a>	<i>Objectivism Student Group</i>	125 min.	<b>22416</b>
39	MICASE	<a href="#">SGR175SU123</a>	<i>Biochemistry Study Group</i>	109 min.	<b>17530</b>
40	MICASE	<a href="#">SGR200JU125</a>	<i>Organic Chemistry Study Group</i>	101 min.	<b>18124</b>
41	MICASE	<a href="#">SGR175MU126</a>	<i>Intro Biology Study Group Chemical Engineering Group Project</i>	103 min.	<b>24514</b>
42	MICASE	<a href="#">SGR195SU127</a>	<i>Meeting</i>	77 min.	<b>11289</b>
43	MICASE	<a href="#">SGR565SU144</a>	<i>American Family Group Project Meeting</i>	85 min.	<b>14116</b>
44	MICASE	<a href="#">SGR999SU146</a>	<i>Senior Thesis Study Group Chemistry Discussion Section Student</i>	64 min.	<b>15483</b>
45	MICASE	<a href="#">STP200JU019</a>	<i>Presentations</i>	51 min.	<b>7303</b>
46	MICASE	<a href="#">STP125JG050</a>	<i>Architecture Critiques</i>	123 min.	<b>24228</b>
47	MICASE	<a href="#">SVC999MX104</a>	<i>Media Union Service Encounters</i>	187 min.	<b>19072</b>
48	MICASE	<a href="#">SVC999MX148</a>	<i>Science Learning Center Service Encounters</i>	121 min.	<b>8613</b>

#### Appendix 4. Sources for the Monologic ASE sub-corpus

No.	Source	Link	Event name	Recording length	Word count
1	MICASE	<a href="#">COL999MX036</a>	<i>Provost Public Lecture</i>	61 min.	<b>9116</b>
2	MICASE	<a href="#">COL605MX039</a>	<i>Women's Studies Guest Lecture</i>	65 min.	<b>10370</b>
3	MICASE	<a href="#">COL385MU054</a>	<i>Public Math Colloquium</i>	51 min.	<b>7664</b>
4	MICASE	<a href="#">LEL500JU034</a>	<i>Intro Psychology Lecture</i>	47 min.	<b>7845</b>
5	MICASE	<a href="#">LEL500SU088</a>	<i>Drugs of Abuse Lecture</i>	68 min.	<b>11115</b>
6	MICASE	<a href="#">LEL485JU097</a>	<i>Intro to Physics Lecture</i>	49 min.	<b>7880</b>
7	MICASE	<a href="#">LEL200JU105</a>	<i>Inorganic Chemistry Lecture</i>	50 min.	<b>6918</b>
8	MICASE	<a href="#">LEL175SU106</a>	<i>Biology of Cancer Lecture</i>	70 min.	<b>11647</b>
9	MICASE	<a href="#">LEL320JU143</a>	<i>Renaissance to Modern Art History Lecture</i>	50 min.	<b>8332</b>
10	MICASE	<a href="#">LEL215SU150</a>	<i>Sports and Daily Life in Ancient Rome Lecture</i>	71 min.	<b>12958</b>
11	MICASE	<a href="#">LEL175JU154</a>	<i>Intro to Evolution Lecture</i>	98 min.	<b>12427</b>
12	MICASE	<a href="#">LES495JU063</a>	<i>Political Science Lecture</i>	96 min.	<b>15359</b>
13	MICASE	<a href="#">LES405JG078</a>	<i>Graduate Cellular Biotechnology Lecture</i>	83 min.	<b>13409</b>

**Appendix 5. Structural classification of lexical bundles across corpora<sup>6</sup>**

Structures	IELTS S3	Interactive ASE	IELTS S4	Monologic ASE
<b>1. Verb phrase expressions</b>				
<b>1.1.</b> <b>(connector+)</b> <b>personal</b> <b>pronoun +</b> <b>lexical verb</b> <b>phrase</b>	I don't think, I think it's, and I think it, I think that's, <i>I think I'll, I think we should, I think it was, you think of the</i> , we're going to, it's going to, you're going to, <i>we'll need to, you'll need to</i> , you don't need, you don't have, <i>we have to do, you'll have to, I'll have to, we'll have to, we didn't have</i> , I don't know, I didn't know, I've got a, we've got to, <i>I'd like to, I'll do that, I see what you mean</i>	I don't know, you don't know, we don't know, I didn't know, I don't think, you're going to, I'm going to, it's going to, we're going to, they're going to, that's going to, I was going to, you're not going, I'm not going, I think it's, I think that's, I think you're, I think I'm, you know what I, you know it's, you know you're, you see what I, you don't have, I don't have, they don't have, we don't have, it doesn't have, I don't want, you don't want, you might want to, you want to do, so you want to, I mean it's, I mean that's, I mean I don't, I mean I think, I mean if you, I don't understand, I don't care, I don't like, I don't remember, I don't see, you don't need, I have no idea, I'm trying to, I have a question, you're talking	I'm going to, <i>today I'm going</i> , we're going to, <i>today we're going, we've been looking</i> , I'd like to	we're going to, I'm going to, you're going to, it's going to, that's going to, they're going to, it's not going, and we're going, and you're going, so we're going, that we're going, I'm not going to, you don't have, it has to do, it has to be, I don't know, you want to make, so I want to, that I want to, I want to do, we don't want, I mean it's, I mean that's, I think it's, we'll talk about, I'll talk about, we're talking about, you look at the, you think about it, we need to think, of course you know, you know it's, I'd like to, I have to say, I told you that, we've got a, you've got to, you don't get, we're looking at, you're trying to, and you can see

<sup>6</sup> Lexical bundles in blue are those shared between IELTS and authentic speech.

		about, it doesn't matter, it has to be		
<b>1.2. (connector+) pronoun/noun phrase + be</b>	that's a good, that's what I, that's right it, it's a good, it's good to, it's hard to, well it's a, there's a lot, then there's the, but I'm not, I'm sure you	that's what I, that's what it, that's what you, that's a good, that's why I, that's kind of, okay so that's, yeah that's what, it's kind of, it's not a, it's like a, it's just like, it's not like, it's the same, it's in the, but it's not, so it's not, no it's not, so it's like, I'm not sure	N/A	that's one of, this is one of, it's a very, and this is the, so this is the, and this is a, there's a lot, and it's not, but it's not, he was able to, and that's what
<b>1.3. verb phrase with active verb</b>	let's think about, is going to be, have a look at, let's have a, let's look at	don't know if, don't even know, don't think so, was going to say, is going to be, are going to be, am going to be, are going to have, are going to do, is not going to, are not going to	am going to talk, going to talk about, going to look at, have been looking at, let's look at	going to talk about, is going to be, are going to be, are going to have, are going to do, going to do is, are going to see, are going to get, is not going to, are not going to, not going to be, going to be a, don't have to, have to do with, has to do with, let's look at, talk a little bit, is a lot of, and let's say
<b>1.4. Yes-no question fragments</b>	do you think it, did you think of, do we need to	do you know what, do you want to, is that what you, does that make sense	N/A	N/A
<b>1.5. Wh-question fragments</b>	what do you think, what did you think, why don't we, why don't you	why don't you, why don't we, what do you mean, what do you think	N/A	N/A
<b>1.6. Lexical bundles with wh-clause fragments</b>	N/A	don't know what, don't know how, what I'm saying, know what I'm, know what I mean, what I mean like, see what	N/A	what we're going, what I want to, what's going on, what you find is, what you get is

		I'm, what's going on, what you're saying		
<b>1.7. Lexical bundles with to-clause fragments</b>	need to think about, to look at the, need to look at, don't need to, don't have to, don't want to, to make sure that, to get hold of	to be able to, don't want to, don't need to	will be able to	to do with the, don't want to, want to make sure, to make sure that, not want to be, didn't want to, need to think about, to be able to, to keep in mind, more likely to be, to figure out what
<b>1.8. (Verb) + that-clause fragments</b>	N/A	don't know the, don't think I, don't think it, don't think that, don't think we	N/A	make sure that you, let's say that
<b>1.9. Adverbial clause fragments</b>	N/A	if you want to, if it's a, if you don't	N/A	if you want to, if you look at, if you think about, if there's a, when it comes to
<b>1.10. Copula be + noun phrase/adjective phrase</b>	is a good idea, be a good idea, would be a good	N/A	N/A	is one of the, to be in the
<b>2. Noun phrase expressions</b>				
<b>2.1. Noun phrase with of-phrase fragment</b>	that a lot of, the end of the, a bit of a, that sort of thing, a lot of work	the end of the	one of the most, a wide range of, the bottom of the, parts of the world	one of the things, and one of the, um one of the, a little bit of, the end of the, a whole bunch of, the rest of the
<b>2.2. Noun phrase with other post-modifier fragment</b>	a good idea to, a look at the, the best way to	N/A	N/A	a little bit about, a little bit more, little bit about the, the best way to, the ways in which, ways in which we

<b>2.3. Other noun phrases</b>	N/A	or something like that	N/A	one of the things
<b>3. Prepositional phrase expressions</b>				
<b>3.1. Preposition phrase with of-phrase fragment</b>	at the end of, in the middle of	at the end of	N/A	of the things that, at the end of, of the twentieth century, in terms of the, in terms of their, in the process of, in the course of, by the name of, in the case of, in the middle of
<b>3.2. Other prepositional phrases</b>	at the same time	at the same time, from here to. here	at the same time, on the other hand, all over the world	in the eighteenth nineties, in the nineteenth century, at the same time, with respect to the, on a regular basis, on the other hand, so in other words
<b>4. Others</b>	but I don't, and I don't, well I don't	I just don't, I don't even, but I don't, and I don't, yeah I don't, okay I don't, so I don't, like I don't, um I don't, well I don't, I don't really, so you don't, if you don't, no no no no	N/A	and I don't, and so on and, over and over again

### Appendix 6. Functional classification of lexical bundles across corpora<sup>7</sup>

Functions	IELTS Section 3	Interactive ASE	IELTS Section 4	Monologic ASE
<b>1. Stance expressions</b>				
<b>1.1. Epistemic</b>	I don't think, I think it's, and I think it, I think that's, I think I'll, I think it was, what do you think, you think of the, do you think it, what did you think, did you think of, I don't know, I didn't know, I'm not sure, I'm sure you, I see what you mean, why don't we, why don't you	I don't know, you don't know, we don't know, don't know if, don't know what, don't know how, don't know the, don't even know, I didn't know, I don't think, don't think I, don't think so, don't think it, don't think that, don't think we, I think it's, I think that's, I think you're, I think I'm, think it's a, I don't care, I don't understand, I don't like, I don't remember, I don't remember, why don't you, why don't we, I'm not sure, I have no idea, no no no no	N/A	I don't know, I think it's, I have to say
<b>1.2. Attitudinal/Modality</b>				
<b>1.2.1. Obligation/Directive</b>	need to think about, do we need to, we'll need to, you'll need to, need to look at, you don't need, don't need to, we have to do, don't have to, you'll have to, I'll have to, we'll have to,	don't have to, doesn't have to, not have to be, it has to be, going to have to, you might want to, you don't want, you don't need, don't need to	N/A	don't have to, it has to be, you've got to, we need to think, need to think about

<sup>7</sup> Lexical bundles in blue are those shared between IELTS and authentic speech.



	we've got to, I think we should			
<b>1.2.2. Intention/Prediction</b>	we're going to, it's going to, is going to be, you're going to, don't want to, to make sure that, to get hold of, I'd like to, I'll do that	you're going to, I'm going to, it's going to, we're going to, they're going to, that's going to, I was going to, is going to be, are going to be, am going to be, are going to have, are going to do, is not going to, are not going to, you're not going, I'm not going, am not going to, not going to be, going to be like, going to be a, are you going to// don't want to, I don't want, do you want to, if you want to, you want to do, so you want to// I'm trying to	I'm going to, we're going to, I'd like to	we're going to, I'm going to, you're going to, it's going to, that's going to, they're going to, is going to be, are going to be, are going to have, are going to do, going to do is, are going to see, are going to get, is not going to, it's not going, are not going to, not going to be, going to be a, I'm not going to, don't want to, if you want to, you want to make, so I want to, want to make sure, I want to do, we don't want, not want to be, didn't want to, to make sure that, make sure that you, to keep in mind, to figure out what, I'd like to, you're trying to
<b>1.2.3. Ability</b>	N/A	to be able to	will be able to	to be able to, he was able to, and you can see, more likely to be
<b>1.2.4. Evaluation</b>	that's a good, it's a good, is a good idea, be a good idea, would be a good, a good idea to, that's right it,	that's a good, it's the same, it doesn't matter	N/A	N/A

	yes that's right, it's good to, it's hard to			
<b>2. Referential expressions</b>				
<b>2.1. Identification/Focus</b>	I've got a, <b>that's what I</b> , the best way to, then there's the, well it's a	that's what I, that's what it, that's what you, that's why I, okay so that's, it's not a, it's just a, it's in the, but/so/no it's not, if it's a, what's going on, what you're saying, is that what you, yeah that's what	one of the most	what we're going, what I want to, that I want to, one of the things, of the things that, and/um one of the, that's one of, this is one of, is one of the, it's a very, and it's not, but it's not, and this is the, so this is the, and this is a, that there's a, so there's a, if there's a, we've got a, and that's what, that's part of, what's going on, what you find is, what you get is, when it comes to, the best way to
<b>2.2. Imprecision</b>	that sort of thing	it's kind of, it's like a, it's just like, it's not like, so it's like, or something like that, that's kind of	N/A	and so on and
<b>2.3. Specification of attributes</b>				
<b>2.3.1. Framing attributes</b>	N/A	N/A	N/A	in terms of the, in terms of their, in the process of, in the course of, by the name of, in the case of, the ways in which, ways in which we, with respect to the
<b>2.3.2. Quantity</b>	there's a lot, a lot of work, that a lot of, a bit of a	N/A	a wide range of	a little bit more, a little bit of, there's a lot, is a lot of, a

				whole bunch of, the rest of the
<b>2.4. Time/Place/Text reference</b>	at the end of, the end of the, at the same time, in the middle of	at the end of, the end of the, at the same time, from here to here	at the same time, the bottom of the, all over the world, parts of the world	at the end of, the end of the, in the middle of, in the eighteen nineties, of the twentieth century, in the nineteenth century, at the same time, over and over again, on a regular basis
<b>3. Discourse organisers</b>				
<b>3.1. Topic introduction/focus</b>	let's think about, to look at the, have a look at, let's have a, a look at the, let's look at	was going to say, what do you think, I have a question, you're talking about	am going to talk, going to talk about, going to look at, today I'm going, today we're going, have been looking at, we've been looking, let's look at	going to talk about, we'll talk about, I'll talk about, we're talking about, if you look at, you look at the, let's look at, we're looking at, if you think about, you think about it, a little bit about, talk a little bit, little bit about the
<b>3.2. Topic elaboration</b>	N/A	you know what I, you know it's, you know you're, what I'm saying, know what I'm, know what I mean, what I mean like, do you know what, you see what I, see what I'm, I mean it's, I mean that's, I mean I don't, I mean I think, I mean if you, what do you mean, does that make sense	on the other hand	to do with the, have to do with, has to do with, it has to do, I mean it's, I mean that's, of course you know, you know it's, on the other hand, so in other words, let's say that, and let's say
<b>4. Special conversational functions</b>				

<b>4.1. Politeness</b>	N/A	N/A	N/A	N/A
<b>4.2. Simple inquiry</b>	N/A	N/A	N/A	N/A
<b>4.3. Reporting</b>	N/A	N/A	N/A	I told you that
<b>5. Others</b>	but I don't, and I don't, well I don't, but I'm not, you don't have, you didn't have	I just don't, I don't even, I don't really but/and/yeah/okay/so/like/um/ well/no I don't, if/so you don't, you don't have, I don't have, they don't have, we don't have, it doesn't have, don't have a, don't have any	N/A	and we're going, and you're going, so we're going, that we're going, to be in the, you don't have, and I don't, you don't get