# The Influence of L1 on Processing of L2 Collocations in Tamil-English Bilingual Children

PhD Applied Linguistics (Taught Track & Thesis)

**Department of English Language and Applied Linguistics, School of Literature and Languages**

Roopa Katherene Leonard

**October 2020**

# Acknowledgements

It would not have been possible to complete this thesis without the encouragement and support of numerous people along the way.

Firstly, I would like to acknowledge my good fortune of having won the PhD supervisor lottery. Dr Michael Daller and Dr Holly Joseph have been incredible supervisors in every way and without them I wouldn't have made it very far with this thesis.

Michael, thank you for your invaluable support, timely input when I needed it, steady guidance, endless patience and encouragement during some of the hardest days of my PhD journey. Thank you for believing in me from the beginning and seeing me through until the end.

Holly, thank you for pushing me to do better when I thought I couldn't, for encouraging me and supporting me particularly during the trials of data collection, for always being available at short notice, for your infinite patience with all my many muddled mistakes, and for helping me through many steep learning curves. I always left your office feeling better and with a renewed sense of confidence.

I would also like to express my gratitude to Dr Jackie Laws who has been a constant source of support and advice from the time I started my PhD. Thank you, Jackie, for your genuine concern for my well-being, your kindness, your words of wisdom and always being available for a chat even after you left the University.

I would also like to thank other staff members from the Department of English Language and Applied Linguistics who have been supportive in various ways over the past four years: Dr Parvaneh Tavakoli, Dr Fraibet Aveledo, Dr Rodney Jones and Dr Clare Furneaux. Being surrounded by such a friendly and intellectually stimulating group of people who were always ready to offer advice on my work or a word of encouragement has greatly enriched my PhD experience.

I would like to thank Dr Daisy Powell and Dr Jeanine Treffers-Daller for all their support and kind words of encouragement, especially during the final phase of my PhD.

# Declaration of Original Authorship

I confirm that this is my own work and the use of material from other sources has been properly and fully acknowledged.

Roopa Katherene Leonard

# **Abstract**

The two studies in this thesis examine the activation of the L1 during the processing of L2 collocations. Models of bilingual lexical representation and access such as the Bilingual Activation Model (Dijkstra & Van Heuven, 2002) and the Multilink Model Dijkstra and Rekke (2010) posit that non-selective, cross-lingual activation is a feature of the bilingual lexicon. Following from increased research in formulaic language, this study uses online processing measures to investigate whether cross-lingual activation can be extended beyond single lexical items to collocations in bilingual children. A self-paced reading (Study 1) and an eye tracking (Study 2) experiment were conducted with Tamil-English bilingual children (age 8-11). In both studies, reading times on English collocations embedded in sentences were measured. All collocations were congruent or incongruent with collocations in Tamil. Study 1 with 58 participants was conducted in India and Study 2 with 80 participants was conducted in the UK and participants across the two studies varied substantially in their English and Tamil proficiency. All children also completed vocabulary tests (English & Tamil) and children in Study 2 also completed a questionnaire on their language background. Results clearly show that children rely on their vocabulary knowledge in L1 to aid their processing of collocations in L2 and that this L1 activation is immediate and can be captured in real time. An important finding from these studies is that the extent of L1 activation is dependent on L1 vocabulary as well as language exposure and dominance: the children in Study 1 showed a much bigger congruency effect than the children in Study 2. These findings indicate that L1 activation in bilingual children is not just for single words but goes beyond word level to collocations as well.

# Table of Contents

# List of Tables and Figures

# Chapter 1: Introduction

## 1.1 Background of the studies

Vocabulary plays a crucial role in language and reading development in children, and in this increasingly globalised world where monolinguals are increasingly a minority (Wei, 2020) there is a need to increase our understanding of how bilingual children process vocabulary. With increased interest in vocabulary research over the past few decades, we now have a much better understanding of how vocabulary processing takes place both in L1 (first language) and L2 (second language) (Schmitt, 2008). However, as noted by Hestetræet (2018), studies on child vocabulary acquisition in the L2 are relatively scarce. While theories of vocabulary development in adults do address the beginning stages of vocabulary acquisition and processing and this could be relevant to children, it is still necessary to carry out more research to understand the complexities of how exactly the process takes place in bilingual children. This thesis comprises two studies which were undertaken to further our understanding of how bilingual children process vocabulary—we know that developing a new vocabulary when already proficient in L1 is very different from developing L2 vocabulary while still acquiring the L1. Both studies were done with young Tamil-English bilingual children in the primary school age group. The first study was conducted in Chennai, India and the second study was carried out in the UK. This chapter will discuss the significance of these studies, give a brief introduction to bilingualism, and present an overview of both studies.

**1.2 Significance of the studies**

With an increasing number of children learning English as a second or additional language in India, the UK, and around the world, it is important to develop our understanding of how L2 <mark>vocabulary processing</mark> in children takes place, especially in relation to the influence of the first language on this process. As mentioned by Pearson, Hiebert, and Kamil (2017), vocabulary research in this area has until recently focused on single words. We often think of vocabulary in terms of individual words, but words actually co-occur frequently in systematic ways to form collocations and other multiword units i.e. formulaic language (formulaic language is further discussed in Section 2.9). Since we know that vocabulary is also acquired and stored in multiword units (Schmitt, 2010), not just as single words, it is essential that we expand our knowledge of the role of the first language (L1) in how bilingual children <mark>process formulaic language in the L2. Studying the processing of formulaic language</mark> can further our understanding of monolingual and bilingual lexicons are formed and how they differ. This knowledge can help in the design of vocabulary instruction material for children as well as inform vocabulary teaching practices in the classroom.

**1.3 Bilingualism**

Montrul (2008) defines bilingualism as knowledge of two languages, although not necessarily to the same degree; as Grosjean (1992) notes, this is contrary to the common belief that bilinguals must master two languages fluently. Since linguistic knowledge is complex and multidimensional at both

structural and processing levels, it is not usually the case that bilinguals develop unbalanced command of either language in one or more dimensions (Montrul, 2013). In his seminal

book *Studying Bilinguals* (2008), Grosjean presents his concept of a wholistic view of a bilingual person as a competent speaker-hearer who has developed competencies in both their languages according to their individual needs and the requirements of the environment—this is known as the Complementarity Principle. In this view, the bilingual uses both their languages together or individually, in different settings, for different purposes and with different people.

In terms of bilingual children, it is critical to investigate the relationship between their two developing linguistic systems (Grosjean & Li, 2013). Early researchers believed that a bilingual child developed one underlying linguistic system for both languages (Volterra & Taeschner, 1978), but a wealth of research has disproved that notion and has led researchers to conclude that bilinguals develop separate language systems or at least language systems that are significantly fused together (e.g. Döpke, 1998; Kim, 2009; Nicoladis, 2012; Paradis & Genesee, 1996; Pearson, Fernández, & Oller, 1993; Poulin-Dubois & Goodz, 2001). It is now well-established in the literature (e.g. Döpke, 2000; Hulk & Muller, 2001; Yip & Matthews, 2000) that both of these developing language systems influence each other—in particular that the two systems may influence each other when there is an actual or perceived overlap between them (Hulk & Müller, 2001). In the field of child L2 acquisition, this influence has been studied in a number of areas such as phonology (Bosch & Sebastián-Gallés, 2001, 2003), syntax (Yip & Matthews, 2000, 2007), speech rhythm (Mok, 2011) etc. The studies in this thesis aim to examine this influence of developing language systems on each other in the context of vocabulary acquisition.

**1.4 Focus of the Studies and Research Gap**

This section will briefly introduce the important concepts of this thesis: collocations, the role of the L1 in the processing of L2 collocations and the congruency effect, and it will also

introduce the research gap which this thesis aims to contribute towards. These concepts will be explored in greater detail in Chapter 2 with a review of the relevant literature.

Collocations are subset of formulaic language that can be broadly defined as words that occur together more frequently that would be expected by chance (Carrol & Conklin, 2020) e.g. *heavy wind*. Studies have shown that collocations, like other kinds of formulaic language, are processed more quickly than non-formulaic language (Bonk & Healy, 2005; Sonbul & Siyanova-Chanturia, 2015; Wolter & Gyllstad, 2011). As noted by Conklin and Carrol (2019), a number of studies have investigated how the correspondence between L1 and L2 contributes to online processing of formulaic language. In terms of collocations, when the L1 collocation has a translation equivalent in the L2 it is known as a congruent collocation, and when no translation equivalent exists it is an incongruent collocation. In this thesis it is hypothesised that if the L1 has a strong influence on the L2, the child will read congruent collocations faster than they read incongruent collocations and this is referred to as the congruency effect. As highlighted by Conklin and Carrol (2019), previous research has shown that if formulaic language shares form and meaning across both languages then learners perform better on comprehension and production tasks which could be indicative of faster processing. They also note that if either the form or meaning is not shared across both languages, it is more likely that processing tends to be slower. Since Tamil and English do not share forms because of their different scripts, the congruent collocations in this thesis are ones that share the same meaning across both languages. This thesis seeks to examine the congruency effect in collocations in bilingual children since previous research on L1 influence on L2 collocation processing has focused on adults. Additionally, this thesis also looks at children with varying degrees of proficiency and exposure in both languages.

**1.5 Study 1: An overview**

Study 1 was carried out in a primary school in the southern city of Chennai in the state of Tamil Nadu in India. English is the dominant language of higher education, business, the judiciary and also contributes significantly to the entertainment industry in India (Prabhu, 1984). With an increase in globalization, this preference and dominance for English has accelerated in most spheres of life. Despite the Tamil Nadu government's efforts to promote English education and development from the primary school level onwards, it has been beset with problems such as ineffective teacher training, lack of resources, large classroom sizes etc. (see Section 3.1 for further details). The school in this study is a private, low-fee English-medium school and caters to children from less privileged families (see Section 3.1 for more details).

In this context where children come from Tamil-dominant backgrounds and learn English at school, it is very likely that their knowledge of Tamil plays a role in their English vocabulary acquisition. This study focuses on a type of multiword phrase known as collocations (see Section 2.12), and it examined how children's knowledge of Tamil collocations influences their acquisition of English collocations.

**1.6 Study 2: An overview**

Study 2 was also carried out with Tamil-English bilingual children, but it was done in the UK with children who attend Tamil weekend schools (also known as complementary schools). The children in this study varied more widely in their linguistic profiles than the children in Study 1: some of them had higher levels of Tamil exposure compared to their peers, but most of them identified as English-dominant (see section 4.8.3). With the rapid increase of children

in UK schools from different language backgrounds (see Section 4.2), it is important to explore how these children acquire English vocabulary and to what extent their home or native language plays a role in this acquisition process.

**1.7 Outline of Research Questions**

The research questions in both studies in this thesis will look at the congruency effect in the processing of English collocations, i.e. whether the L1 influences the processing of L2 collocations and how this influence is associated or predicted by measures of vocabulary in both languages, as well as proficiency and exposure. Overall, it is predicted that children who are more proficient in Tamil are more likely to show an effect of congruency when they read English collocations. Since the sample groups are different for each of the studies, the predictions have minor differences and so the research questions for each study will be presented in the relevant chapter for each study.

**1.8 Summary**

This introductory chapter has presented the reasons these studies were carried out and has given overviews of the contexts of both studies. Chapter 2 will discuss the theoretical models behind these studies in detail, explain the concepts of formulaic language, collocations, and cross-linguistic influence and end with a review of previous studies done in this field. This will be followed by Chapter 3 which will outline the research questions, describe the design, methodology, procedure and results of Study 1 in detail and also present a brief discussion of the results. This thesis does not have a separate methodology chapter because different methodologies were used for Study 1 and Study 2, therefore each methodology will be described in its respective chapter. Chapter 4 will then list the research questions, present the design, methodology, procedure and results of Study 2, along with a brief discussion of the

results. The main discussion chapter is Chapter 5 which will discuss the results of both Study

1 and Study 2 in terms of the theoretical models, concepts of language exposure and

dominance and in the context of other studies done in this field. It will also list the limitations

of these studies and discuss the implications. Chapter 6 will conclude this thesis with a

summary of the differences between Study 1 and Study 2, present a recap of the findings and

end with the significance of the findings.

# Chapter 2: Theoretical Framework and Literature Review

## 2.1 Introduction

This chapter will lay out the background and rationale for the present studies (Study 1 and Study 2), in terms of theoretical understanding as well as research done in this field. It begins with a discussion of different aspects of language development in bilingual children. Since the present studies are concerned with the processing of collocations during reading, in the next section a discussion of the major models of vocabulary acquisition and processing in bilinguals will be presented. Next, the importance of formulaic language and collocations will be considered along with research that has been done in the field of the influence of L1 on the processing of collocations. The chapter will conclude with a summary of what we know and then present the research gaps that the present studies aim to fill.

## 2.2 Bilingual Children

According to McCardle and Hoff (2006), the number of children being raised in bilingual homes is growing, yet the language development of these children is not yet fully understood. A crucial issue for parents and educators is the rate of language development for children who grow up with two or more languages i.e. the question of whether the child will have "the same ages of onset for major language milestones" (Poulin-Dubois & Goodz, 2001:95). It is a common worry that bilingual children will show slower rates of language development when compared to monolingual children (Poulin-Dubois & Goodz, 2001). Most of these concerns stem from the popular but controversial concept of the 'balanced bilingual'—according to Müller and Pillunat (2008), if both languages develop at an equal pace, the child can be

considered a balanced bilingual. On the other hand, if one language develops faster or slower than the other, the child is considered an unbalanced bilingual. This concept of a balanced bilingual has been criticized by Grosjean as early as 1989 who states that the bilingual should not be viewed as two monolinguals in one, rather a bilingual's languages should be studied in tandem with each other, keeping in mind domain and usage. Research has shown that while there is no indication that a child's language development is affected in the long term by the addition of another language, the pace of language development may not be the same for both languages (Paradis & Genesee, 1996). However, these initial disparities in rate of development that are present during the early years are negligible and the child can be expected to develop each language steadily in time and gain command of both languages, although evidence shows that a bilingual's lexicon in each language is smaller overall but larger if both lexicons are added together (Genesee, 2002). Of course, quality and quantity of input and exposure are necessary for this process and will be discussed in detail later in this chapter.

Although the tendency to measure bilingual linguistic development in children against the monolingual standard is still present, there is a growing awareness that bilingual language development should be viewed in its own right and should be studied accordingly. In the present studies, the children began acquiring both languages (Tamil and English) in early childhood, but some of them are simultaneous bilinguals and others are sequential or successive bilinguals.

### 2.2.1 Simultaneous bilingualism

As noted by Serratrice (2018), bilingual children may have been exposed to two languages from birth (bilingual first language acquisition [BFLA]) in which case they tend to treat both languages as first languages—widely known as simultaneous bilingualism (Grosjean, 2010).

These simultaneous bilingual children are far less numerous than sequential bilinguals. Grosjean (2010) goes on to explain that simultaneous bilingualism can occur in children when each parent uses a different language with the child or when other caretakers such as staff at daycare or nannies use another language to talk to the child. Thus, these children receive input in two languages (possibly more) and hence during their early years, they acquire two languages. Even though the same developmental mechanism that applies to monolingual first language acquisition applies to these children, it manifests itself in different ways (Serratrice, 2013). Although there is some variability in the rate of language acquisition for bilingual children, the two groups appear to achieve the same milestones within the same age spans (Grosjean, 2010). In terms of vocabulary acquisition, Werker and Byers-Heinlein (2009) state that the two vocabularies of bilingual children seem to achieve typical vocabulary milestones, as long as they don't have minimal or cursory exposure to one of their languages—they need good exposure to both languages. Another study (Burns, Yoshida, Hill, & Werker, 2007) found that bilingual infants are able to efficiently establish phonetic representations in both of their languages on the same timescale as their monolingual peers. Two studies that looked at longitudinal vocabulary development in simultaneous bilingual children (Pearson, Fernandez, Lewedeg, & Oller, 1997; Pearson, Fernández, & Oller, 1993) found that bilingual children show the traditional lexical spurt—a sudden increase in vocabulary—either alternately, depending on the strength of each language and exposure, or the combined vocabulary of both languages. These researchers also report similarities between monolingual and bilingual children: sounds or sound groups that are easier to produce are the first ones produced by monolinguals and by bilinguals in both their languages, utterances gradually get longer, and the construction of these utterances changes from simple to more complex ones.

There are two main views on how simultaneous bilingual children develop their language system, although only one is commonly held today: proponents of one view claim that bilingual children develop a single unitary language system that is differentiated as they grow older, and other scholars favour the view that children develop each language as a separate system right from the beginning. Proponents of the first view cite evidence from much earlier studies by linguists like Leopold (1939), Volterra and Taeschner (1978) and Vihman (1985) who report examples of bilingual children mixing languages as evidence for the existence of one language system. However, as Genesee (2000) observes, these studies fail to provide context for their conclusions and while they do show instances of the children mixing languages, they do not provide sufficient evidence to show that children indiscriminately use both languages in all contexts of communication that would point to a unitary language system. There is a lot more evidence for the second view that bilingual children develop both their languages separately: numerous studies have found that infants less than a year old with bilingual mothers can discriminate between both languages (e.g. Bosch & Sebastián-Gallés, 2003; Byers Heinlein, Burns, & Werker, 2010). Meisel (2004) draws on a wealth of more recent research to substantiate his position that simultaneous bilingual children can differentiate the grammatical systems of both languages early on without much effort and that any mixing is just an instance of the child code-switching.

Some of the children in Study 2 are simultaneous bilinguals and it is important to take their language development into account as the influence of Tamil on their processing of English collocation is examined. This overview of simultaneous bilingualism reveals that simultaneous bilingual children need adequate exposure to both languages right from the beginning in order for them to have balanced language development. This notion will be explored with reference to language input, language exposure and language dominance later in this chapter.

*2.2.2 Sequential bilingualism*

In the case of early L2 acquisition (ESLA), children are exposed to a L2 after they have started acquiring their first language—this transition usually occurs when the children start childcare, nursery or school and it falls under sequential (successive) bilingualism (Grosjean, 2010). ESLA children typically hear a first language (L1) at home and then are exposed to the L2 (L2) when they attend regular childcare in the community or when they start school, unlike BFLA children who usually hear both languages in the home environment. This is different from the typical view of L2 acquisition which is typically characterized as the introduction of a L2 after a first has been developed. In her work that studied language acquisition in immigrant children, Wong Fillmore (1991) lists three components and three processes that enable children to acquire a L2 naturally. The three components are: (i) the children (learners) who have to learn a L2, (ii) the speakers of the language who will help them do so, and (iii) the social setting which usually refers to school and the community. The three processes are: the social aspect which involves the children observing what is spoken about in various social settings and asking the speakers for accommodations and adjustments when needed; (ii) the linguistic aspect which involves children using their existing knowledge of linguistic categories and structures to look for equivalents or differences in the L2; and (iii) the cognitive elements which means that using the L2 input they receive, children have to discover the rules and units of the language and integrate this knowledge into a grammar using analytical skills, memory, inferential skills etc.

McLaughlin (1995; 2013) notes that there are several factors that affect these language acquisition components and processes such as cultural, linguistic, and social differences as well as different education systems and individual attitudes and different cognitive abilities.

In a study done with 169 sequential bilinguals aged between four and seven who had one to five years of English exposure, Paradis (2011) found that internal variables such as chronological age, nonverbals intelligence and phonological short-term memory were better predictors of English proficiency than external variables like length of exposure, proportion of English use of home, richness of English use in afterschool activities and self-rated English proficiency of the mother. A study that examined L2 proficiency in Russian-Hebrew and English-Hebrew children (Armon-Lotem, Joffe, Abutbul-Oz, Altman, & Walters, 2014) found that L2 proficiency was very closely linked to attitudes towards the heritage language and L2 as well as ethnolinguistic identities—children with greater L2 proficiency lived in communities where there was a lot of mixed use of the heritage language and the L2 i.e. the Russian-Hebrew community and not the English-Hebrew community.

This overview of sequential bilingualism in children shows that there are many different variables which interact in the development of both languages in a sequentially bilingual child. The children in Study 1 and many of the children in Study 2 are sequential bilinguals and it is important to consider these factors when analysing how Tamil influences their processing of English collocations.

**2.3 Cross-linguistic influence**

In reference to L2 acquisition, cross-linguistic influence is understood as the influence of the knowledge of a previous language on the acquisition of a L2. However, cross-linguistic influence is different in both BFLA and ESLA children since they either acquire both languages simultaneously, or are still in relatively early stages of developing their first language when they are exposed to their L2. While there is no conclusive evidence on whether the knowledge of these languages is represented in a single or dual mental system, there is substantial evidence that however this knowledge is represented, the languages do not

develop in an autonomous fashion; instead, they influence each other with respect to different language features. Ortega (2009) observes that L1 influence on L2 acquisition is often subtle and selective and can have markedly different positive or negative consequences on L2 development depending on factors such as L1 background, different stages of L2 development, individual proficiency, different areas of the L2 etc.

A number of factors influence cross-linguistic transfer and one that has been studied extensively is L1-L2 differences and similarities (Ortega, 2009). Early studies (Long & Sato, 1984; Hyltenstam, 1977; Zobl, 1980) investigating cross-linguistic transfer in the context of these differences found that similarities do not always help in L2 acquisition and certain differences do not seem to impede L2 acquisition. However, more recent studies (Kim, 2009; Muller & Hulk, 2001; Nicoladis, 2012) in child bilingualism have provided more insights into how these differences and similarities influence language acquisition in children and have shown that they do influence cross-linguistic transfer, but are subject to directionality (the direction in which the transfer occurs) and grammatical/linguistic category. Directionality is heavily dependent on which of the bilingual's languages is more developed—it is more likely that crosslinguistic transfer happens from the stronger language to the weaker language. With regard to grammatical or linguistic categories, it is more likely that cross-linguistic transfer occurs when there is overlap or ambiguity between similar structures in both languages: a study by Nicoladis (2006) found that this was the case with adjective-noun strings in French-English bilingual children since English has only one order (adjective-noun) whereas in French, both adjective-noun and noun-adjective word orders are allowed. Another factor that has been found to influence crosslinguistic transfer is identified by Ortega (2009) is interlingual identification which is when the individual makes "the judgement that something in the native language and something in the target language are

similar (Odlin, 2003, p.45). This judgement can be a conscious, strategic choice to rely on the L1 to compensate for gaps in L2 knowledge, but it can also happen at a subconscious level.

Murphy (2014) distinguishes cross-linguistic influence into two kinds: qualitative cross-linguistic influence in the form of non-grammatical utterances and quantitative cross-linguistic influence which results in particular types of language errors that are typical in monolingual development, but at a rate higher than the typical monolingual rate (Paradis, Nicoladis, Crago & Genesee, 2011). From studies that have investigated the occurrence of both these kinds of cross-linguistic influence (Paradis et al., 2011; Yip and Matthews, 2000), it seems that bilingual children are more prone than adults to quantitative cross-linguistic influence i.e. they are more likely to use their knowledge of one language while producing constructions in their other language . This suggests that the kind of errors that a bilingual child makes are the same kinds that a monolingual child makes, only that a bilingual child makes more of these errors while developing linguistic competence in both languages. A recent case study by Babatsouli Nicoladis (2019) investigated cross-linguistic influence in a 4-year old bilingual child in the context of fixed expressions such as collocations and found that the child exhibited a degree of compositionality i.e. she showed evidence of cross-linguistic transfer in her use of collocations. This raises an interesting question as to whether this cross-linguistic influence is only evident in younger bilingual children or whether older bilingual children exhibit it too. The purpose of the current studies in this thesis is to determine whether the L1 is activated when bilingual children (7-10 years old) read and process collocation and the findings of this study will add to the growing knowledge base in this area.

**2.4 Language Dominance**

Within the field of child bilingualism research, there is a general consensus that a bilingual child is unlikely to develop the same set of skills, knowledge and competence in both languages (Murphy, 2014). According to Silva-Corvalan and Treffers-Daller (2016), even when children acquire two languages simultaneously from birth, they are very likely to demonstrate a higher level of proficiency in one of the languages. Meisel (2007) observes that "balanced bilingualism" in children is almost impossible to achieve and that one language will always be dominant. Language dominance, which can be explained as the "predominant of ambient languages in a given setting" (Meisel, 2007, p. 498) is a difficult construct to define since it is so closely tied up with language preference which refers to an individual inclination to use one language over the other. These two concepts are closely interrelated because if a language is dominant in the environment where the child lives, the child will have increased exposure to this language and will likely develop a preference for it in many cases. While it is true that factors like individual preference, frequency of use and complexity of language use are important factors in language dominance, Meisel (2007) suggests that the overarching feature determining language dominance is the context in which the child uses the language and goes on to list three context-related issues that contribute to language dominance: (i) the communicative and learning environment, (ii) the use bilinguals make of both their languages, and (iii) the individual development of linguistic knowledge. Therefore, a child's dominant language is likely to be the language that has a stronger presence in the child's environment. This explains why children with a minority language (i.e. a home language that is not used in wider society) often end up with a more unbalanced bilingual profile if they are not provided with adequate support for the home language. This idea of different languages being used in different contexts is known as language domains:

bilinguals distribute their languages across the domains of their lives and sometimes languages can be used across different domains (e.g. home, religious activities) and in some cases a language can be restricted to just one domain ( e.g. home or school) (Grosjean, 2010). The amount of time spent in each domain can affect language dominance.

Different criteria have been used to measure language dominance such as experience-based measures (age of exposure, length of exposure, language preference etc.) and performance-based measures (assessments of fluency, vocabulary knowledge etc.) which were the focus of a study done by Bedore, Pena, Sommers, Boerger, Resendiz, Greene, Bohman, and Gillam (2012). This study looked at language dominance in Spanish-English bilingual children in the United States and analysis of the results support the dynamic nature of language dominance—it was found to vary as a function of the point in time at which it was measured as well as the ways in which the child's language experience was measured. The factor that most reliably predicted language proficiency and dominance was "current language use" i.e. the children scored higher on performance assessments in the language they had the most experience using. The pattern of a bilingual child's language use can also be understood in terms of how they use their dominant and non-dominant languages. For example, in a study that investigated how children used code-switching to fill lexical or syntactic gaps, Nicoladis and Secco (2000) found that children more frequently used their dominant language than their non-dominant one to do so. In another study on cross-linguistic influence in bilingual children by Yip and Matthews (2000), it was concluded that the dominant language was more likely to influence the non-dominant one. This is because the child relies on the resources of the dominant language to compensate for lack of resources in the non-dominant language.

To conclude, language dominance in bilingual children depends on factors such as individual language proficiency, the status of the language in wider society, and most importantly, the exposure and experience the child has with each language. This last factor that depends on

input has received close attention and recent years and empirical evidence (Bedore et. al, 2012; Nicoladis & Secco, 2000) suggests that quantity and quality of input is the determining factor when it comes to determining the child's dominant language and consequently also determines the extent of unbalanced bilingual development.

## 2.5 Language Exposure

From the previous section, it is evident that language exposure is a crucial element in the development of bilingualism in children. There are different factors that contribute to the impact of exposure on the linguistic development of bilingual children. Armon-Lotem and Meir (2019) posit that a study of these factors will lead to a better understanding of the two competing views of language acquisition: a usage-based approach in which input and exposure are given a central role in language development (Bybee, 2006; Gathercole & Hoff, 2007) and the generative-nativist approach that accords importance to Universal Grammar in the development of language, thus placing much less importance on exposure and input (Pink & Bloom, 1990; White & White, 2003; White, 1989; Schachter, 1988 ).

Age of onset (AoO) is perhaps the most basic measure of exposure in early childhood bilingualism. For BFLA children (simultaneous bilinguals), age of onset is the same for both languages i.e. birth. For ESLA children (successive bilinguals), the age of onset for the first language is birth and the age of onset for the L2 is when exposure to that language begins. It is straightforward to determine age of onset for BFLA children and for ESLA children whose exposure to a L2 began when they migrated to a new country, but it can be complicated to determine AoO for children for whom the first language is a heritage language and the language of wider society is the L2 because it is often difficult to explore children's daily routines over a long time period to determine L2 exposure before preschool/school entry. Length of exposure (LoE) to a language measures the time from AoO to the child's

chronological age. For BFLA children, AoO and LoE are the same for both languages whereas for ESLA children, AoO and LoE are the same for the first language and differ for the L2. In addition to these basic measures, relative exposure and absolute exposure are more nuanced measures that take into account the amount of input received in each language (Armon-Lotem, De Jong, & Meir, 2018). Relative exposure refers to the approximate amount of exposure to each language in comparison to the other and is usually determined by questionnaires that assess how much each language is used in different scenarios or environments and how frequently they are used with different people such as family and friends. Absolute exposure is determined by measuring and recording child-directed speech, regardless of the exposure to the other language. It takes into account the fact that the density of exposure i.e. the number of words may not correlate with the percentage of exposure. Measures of absolute exposure also capture details of what kinds of lexical and syntactic constructions the child is exposed. The quantitative measures of language exposure introduced above will now be discussed in detail.

Length of exposure is the most simplistic measure and data on LoE is usually gathered in a questionnaire. There are several questionnaires used by researchers, either for parents or children and sometimes for both parents and children. Examples of these include the Alberta Language Environment Questionnaire (ALEQ) developed for use in the bilingual context of Canada (Paradis et. al, 2011) and the COST Action IS0804 Questionnaire for Parents of Bilingual Children (Tuller, 2015) developed for use in the multilingual European context. Paradis (2011) used the ALEQ in a study that looked at reading comprehension, vocabulary size, and morphosyntactic awareness of 169 sequential bilinguals between the ages of four and seven. This study was designed to measure the relative contribution of external and internal variables to a child's L2 proficiency and one of these external variables was LoE. According to the results, internal variables such as nonverbal intelligence and phonological

short-term memory had much higher predictive value for L2 vocabulary size and morphology than LoE. While LoE did predict L2 proficiency, it had low predictive value and regardless of LoE, older children had more developed L2 proficiency. In a study on the effects of language exposure on bilingual children, Armon-Lotem, Joffe, Abutbul-Oz, Altman and Walters (2014) also found that LoE had a low predictive power with regard to L2 proficiency scores; the results of their study showed that relative exposure and prestige of the heritage[1] language in society both had much higher predictive powers for L2 proficiency scores than LoE. Thus, it appears that while LoE has some predictive value for L2 skills and proficiency, other variables related to language exposure have bigger roles.

As an improvement on LoE, relative exposure is a quantitative measure that focuses on the amount of exposure the child receives to each language when compared to the other. It offers a more precise measure than LoE because it looks at language exposure in a proportionate fashion.  Different aspects of relative exposure such as home vs. school, exposure via media vs. child-directed speech are compared to determine how they contribute to linguistic skills and competencies in bilingual children. Gutierrez-Clellen and Kreiter (2003) compared grammatical competence to measures of relative exposure for a group of eight-year-olds in California who had Spanish as their home language (first language) and English as the language of society (L2). Out of all the measures of relative exposure, they found that only the relative exposure to Spanish at home predicted grammatical performance in Spanish. Exposure to English at home did not significantly predict children's linguistic abilities in English. In a large study with Spanish-English bilingual children, Peña, Bedore, and Byers-Heinlein (2018) compared the AoO for English and relative exposure to English and Spanish with their morphosyntactic and semantic knowledge. They found that AoO had high

---

[1] A heritage language is a language with which an individual has a historical and personal connection—this is what makes it salient and not necessarily the actual proficiency of individual speakers (Valdes, 1999)

predictive power for English performance although the effect decreased with age. The relative exposure to Spanish explained a large amount of the variance in English performance. In this study, the AoO and relative exposure were highly correlated with each other, making it difficult to determine the effect of the factors separately. In a study conducted in Montreal with French-English bilingual children aged five, Thordardottir, Grüter, and Paradis (2014) calculated relative exposure based on estimated proportion of exposure to each language since birth at both home and preschool. A single measure of time spent in an English or French dominant environment was computed by summing up their yearly proportions of exposure to each language. The exposure scores for each language were highly correlated with the vocabulary scores for that language and there was no effect of AoO. These studies show that relative exposure has effects on language proficiency, in particular vocabulary and morphosyntactic knowledge. However, it is important to remember that relative exposure is based on parental reports—this means that it may not accurately reflect the exact exposure the child receives to each language and like any self-reported measure, it does not include details on the specific aspects of language a child is exposed to. It is also possible to get similar relative exposure measures for two different scenarios: one in which the child received a lot of input in both languages and the other scenario where the child is not talked to much but the relative exposure for both languages could be the same.

In order to collect informative data on the nature of language input along with the amount of exposure, De Houwer and Bornstein (2003) designed the Language Input Diary (LID) as measure of absolute exposure. The LID is meant to capture detailed accounts of who a child interacts with and in which language. It is supposed to be filled in by the child's caregiver at a given time, in 30-minute blocks from 6 AM to 9 PM and is intended to be used regularly, not as a one-time assessment. In a study using LIDS with 90 Spanish-English bilingual toddlers, Place and Hoff (2016) were able to gather information on how much mixed input

the children received from their caregivers, how many speakers used each language with a child and which carers reported being native speakers of English or Spanish. While the relative exposure to Spanish and English was comparable, the LIDs showed that 94% of Spanish exposure was provided by native speakers while only about 40% of English exposure was provided by native speakers. This demonstrates that while relative exposure may be almost the same for each language, the quality of input can vary widely. While LIDs yield rich data, they are quite difficult to collect since they are time-consuming to fill in and require longitudinal commitment from the child's parents and caregivers. Researchers have developed detailed questionnaires (Thordardottir, Rothenberg, Rivard and Naves, 2006) in an attempt to capture some of the precise information captured by LIDs. Ideally, absolute exposure can also be measured directly by observing caregiver-child interactions; quantitative measures focus on the types of speech in child-directed speech (e.g. number of utterances, word tokens, syllables etc.) while qualitative measures focus on the richness of the input (e.g. number of word types, availability of syntactic structures etc.). In a study that compared the effects of measures of relative exposure to absolute exposure on bilingual children's language use (Marchman, Martinez, Hurtado, Gruter, and Fernald, 2017), researchers gathered information on 18 Spanish-English bilingual children. Information on relative exposure was obtained through a questionnaire and data on absolute exposure were extracted from natural interactions between the child and caregiver for at least eight hours per child. Results showed that there were moderate correlations between reported relative exposure and children's linguistic abilities in each language, but the absolute exposure measures had much higher predictive powers for children's linguistic abilities in each language, thus exemplifying the advantage of measuring absolute exposure. In a case study with a Japanese-English bilingual child, Nakamura and Quay (2012) found a high correlation between the number of word tokens and types in the child's input and the number of tokens and types the child knew in

terms of linguistic ability. From the results of these studies, it is evident that measures of absolute exposure are very useful in predicting bilingual development in children.

From the studies discussed in this section, we can see that relative exposure are absolute exposure offer more detailed insights into childhood bilingual development than AoO and LoE. However, both relative and absolute exposure have drawbacks—relative exposure measures rely on self-reporting and measures of absolute exposure can be difficult to collect since it requires a high level of commitment and cooperation from families. There are quite a few studies using relative exposure but only a few using absolute exposure. Studies using different measures of exposure will help provide a better understanding of how the linguistic competencies of bilingual children develop in each of their languages.

## 2.6 Input: Quality and Quantity

Access to the nature of input from absolute exposure measures allows researchers to make connections between the type of linguistic constructions a child is exposed to and bilingual language development. An issue that has come up in studies of absolute input is how the quantity of native and non-native input that bilingual children receive in the L2 influences language development. It must be noted that a complexity with distinguishing between native and non-native input is that these binary categories often do not reflect the quality of child-directed speech: while the "native" category can encompass a range of language proficiencies, the "non-native" category does not necessarily mean lower quality of input. Studies that have investigated the effect of input using this distinction have found conflicting results. In a study that looked at non-native English input in very young simultaneous bilinguals, Place and Hoff (2016) found that this input was not very supportive of English development. Yet in another study that looked at the effect of foreign domestic caregiver English input on older sequential bilinguals with Mandarin-speaking parents, Dulay, Tong

and McBride (2017) found that it had a positive effect on children's English development when compared to children in the same sample who received input in only Mandarin or Cantonese from their foreign domestic help, even though the English-speaking domestic helps were themselves L2 speakers. Other case studies (Chung, Liu, McBride, Wong & Lo, 2017; Place & Hoff, 2011) have shown that children often achieve high levels of proficiency in their L2s even when exposed to non-native input in that language from a parent or caregiver. Thus, the distinction between native and non-native input may not provide a clear picture of how quality of input influences language development and it is more beneficial to instead assess the types of linguistic structures they are exposed to in terms of quantity and quality.

Chan (2010) conducted a study that focused on the role of absolute input for the acquisition of double object constructions in Cantonese by Cantonese-English bilingual children when compared to their monolingual Cantonese peers. The study found that the bilingual children heard canonical double object constructions in Cantonese about half as frequently as the monolingual children which is probably why the bilingual children took longer than the monolingual ones to acquire these structures. Another study that exemplifies the importance of the specific type of input was done by Bialystok, Luk, Peets, and Yang (2010). They found that primary school children with different home languages but attending English schools performed comparably to their monolingual peers on an English comprehension test of "school" vocabulary, but there was a significant gap in performance between the monolinguals and the bilinguals in a similar comprehension test of "home" vocabulary. The reason for this is that monolingual children receive plenty of exposure to English at home while bilingual children probably receive exposure to home-related concepts in their home languages. Thus, for monolingual children the language of the home domain and school

domain is the same, but for bilingual children, the language of the school domain does not necessarily cross over into the home domain.

Altogether, there are very few studies that investigate how absolute exposure affects a child's linguistic development in the home language and the L2 and specifically how characteristics of the input affect linguistic development. It could be the case that bilingual children's language acquisition diverges from baseline patterns of monolingual language acquisition although as seen above, such divergences do not always take place. Further research is needed in this area for a more comprehensive and nuanced understanding of how the quantity and quality of input influences language development.

**2.7 The Bilingual Mental Lexicon**

The question of whether concepts are stored separately or together in the bilingual speaker's mind has been central to the study of bilingual language development from the earliest stages of bilingual research. Over the years, researchers have conceptualised the bilingual mental lexicon in a number of ways, with the aim of establishing whether a bilingual's two languages are stored in a single mental lexicon or separate ones. This section will examine important models that aim to explain the process of lexical storage, lexical access, and word recognition in bilinguals and then proceed to a discussion of how the bilingual mental lexicon is presently understood in terms of the latest research. Although all models have drawbacks and limitations, different features from various models can be used to explain the bilingual lexical network.

*2.7.1 Language Representation and Storage: Different Views*

In a systematic account of bilingual models of language representation, Marini and Fabbro (2007) identify Weinreich's 1953 model as one of the earliest attempts to explain language representation in the bilingual mind. Weinreich based his seminal model on Saussure's distinction between the signifier (the word) and the signified (the concept). According to Weinreich (2010), there are three possible ways in which the word form and word meaning can be represented in the bilingual mind: the coordinate system, the compound system and the subordinate system. In the coordinate system, the bilingual individual has two separate stores for words in each language (L1 and L2) and two corresponding separate stores for word meanings. In the compound system, there are two separate stores for the words of each language (L1 and L2), but both of them share the same store for concepts/meanings; however, in the subordinate system, the L1 word is linked to the concept/meaning store and so the L2 word has to access meaning through the L1. Weinreich proposed that none of these systems are necessarily mutually exclusive and they can co-exist in the bilingual mind, dependent on factors such as individual proficiency, word type, word frequency etc. He suggested that for a given individual it would be possible that for certain words compound or coordinate relationships would form while for other words, subordinate relationships would develop. In contrast, Kolers (1963) proposed that bilingual lexical storage consists of two separate systems that function independently, and any interaction occurs only through translation. However, Kolers' conception of the separated bilingual lexicon did not gain much traction and has been contested by a majority of researchers who favour the concept of an integrated lexicon.

This view of a completely integrated lexicon has had supporters over the past few decades. Notably, Cook (1992) and Cook and Cook (1993) draw from a wide range of bilingual research that, according to them, provides overwhelming support to the view of an integrated mental lexicon. Some of the salient points from the evidence they cite are as follows: (i) a study by Caramazza and Brones (1979) strongly indicates that reaction time to a word in one language is related to the frequency of its cognate in another known language; (ii) during the processing of interlingual homographs (e.g. English/French *coin*), bilinguals access the meanings in both languages, not just in the target language (Beauvillain & Grainger, 1987); and (iii) translation performance between two known languages is influenced by morphemic similarities between them (Cristoffanini, Kirsner & Milech, 1986). The phenomena of code-switching and code-mixing are also seen as evidence to support the concept of an integrated lexicon; Muller-Lance (2003) observes that the frequency of switching between a bilingual's known languages strongly favours the view of an integrated lexicon.

A closer examination of the evidence from studies of bilingual language development has led to a position that has found favour with an increasing number of researchers—as proposed by Kroll and Tokowicz (2005) and supported by other language researchers (Brysbaert & Dijkstra, 2006; De Angelis & Dewaele, 2009), there is no reason to believe that the way languages are represented in the mental lexicon is the same for all features of a lexical item e.g. orthography, phonology, semantics and syntax. According to this position, these different features and aspects of words may be organised in different ways depending on variables such as language learning environment, age of acquisition, aptitude, typology and individual differences. Cook (2002) later on adopts this view, stating that neither complete separation or integration is likely and it is highly plausible that different relationships exist in the mental lexicon for different features of language such as vocabulary, syntax and pragmatic functions. In a synthesis of studies that look at differing views of the bilingual lexicon, Singleton (2003)

concludes that there is sufficient evidence to argue against total integration, but there is overwhelming evidence for a "very high degree of cross-lexical connectivity and interaction" (176).

### 2.7.2 Lexical Access

Following on from and closely intertwined with the issue of lexical storage, another issue that has been at the forefront of bilingual language development research is lexical access i.e. how lexical items are retrieved from the bilingual mental lexicon during language processing. Hofer (2015) identifies the two major positions that language researchers have taken on this issue. One position is the language-selective or non-specific view of lexical access which predicts that only lexical items in the target language are activated in the selection and retrieval process i.e. lexical items in the non-target language do not compete for selection with items in target language. However, as discussed above, there is little evidence to support the idea of completely separate lexical stores, and so consequently there is not much backing for this selective view of lexical access.

Presently, there is plenty of support for the non-selective view of lexical access which predicts that activation during lexical selection is not restricted to the response or target language since scholars now agree that the lexical stores in bilinguals overlap at least partially (French & Jacquet, 2004). This view of non-selective access has its foundation in the view of an integrated or at least overlapping (partially integrated) mental lexicon that was discussed in the previous section.

### *2.7.3 Models of Lexical Representation, Storage and Access*

To gain a clearer understanding of how the theoretical understanding of bilingual language processing has developed over the years, seminal models of bilingual lexical representation and access will be discussed in this section. The discussion will provide an overview of the salient features and mechanisms of each model, empirical evidence in favour of certain features and also present the shortcomings in the context of empirical research.

### *2.7.3.1 Revised Hierarchical Model*

Using Weinreich's distinct explanation of bilingual language representation as a foundation, researchers developed models of language representation that came to be known as hierarchical models because they consist of three components or nodes. These models with three components differ from previous conceptions of the bilingual mental lexicon because they assume a separation between the form and meaning of a word. Initially proposed by Kroll and Stewart in 1994, the Revised Hierarchical Model (RHM) is the earliest and most recognised hierarchical model of bilingual concept and word representation. Experimental findings that bilinguals are faster at translating words from L2 (L2) into L1 (native language/first language) formed the basis for the creation of this model (Kroll & Stewart, 1994). Although it was originally created as a model of word production, it was later modified to include word acquisition (Kroll, Van Hell, Tokowicz & Green, 2010) and a primary presupposition of this model is that the L1 lexicon is separate from the L2 lexicon. As seen in Fig 2.1, the RHM assumes that the L1 mental lexicon is larger than the L2 mental

lexicon. The link between the L1 lexicon and concepts is seen as bidirectional and very strong since a person initially acquires concepts in their native language.



*Figure 2.1* Revised Hierarchical Model (Kroll and Stewart, 1994, p.26)

As a person acquires their L2, the L2 lexicon is built up and a connection is established between the L1 lexicon and the L2 lexicon as represented by the solid directional line between the L2 and L1 lexicons. There is an opposite directional arrow between the L1 and L2 lexicons, but this connection is seen as weaker than the former, since bilinguals usually acquire the translation of the L2 based on the L1 lexicon and not vice versa. Similarly, the connection between the L2 lexicon and concepts is also seen as weaker. However, according to Kroll and Stewart (1994) these weaker links become stronger as bilinguals achieve higher levels of fluency and proficiency in their L2. Furthermore, if the L2 is used more frequently and continuously, it means that the link between the L2 lexicon and concepts becomes stronger. Therefore, the frequency with which bilinguals need to access the L1 translation while processing words in the L2 presumably decreases as they advance in L2 proficiency.

The primary strength of the RHM is that it shows that in the bilingual lexicon, lexical access and conceptual access change based on individual proficiency. Since its conception, the RHM has dominated the conversation of bilingual language representation: Brysbaert and Duyck (2010) identify a few major contributions that the RHM has made to this field. Firstly, it separated lexical and conceptual representations by supporting distinct representations when a task involves activating word forms. In contrast, when a task involves accessing conceptual information, it supports a shared representation. This helped explain the early, conflicting evidence with regard to the issue of shared or separate representation in terms of form and meaning. While the RHM was not the first model to conceptualize this distinction, it was the first to clearly state the implications of the hierarchical model on how language is represented for bilingual processing. Secondly, since the RHM distinguishes between the L1 and L2 lexicon it gave credibility to studies that found minimal interference from one language during the processing of another i.e. lending support to the idea of selective access. Thirdly, the RHM offers an explanation as to how representation changes as a function of developing proficiency. It suggests that in the early stages of L2 acquisition, the L1 acts as a mediator between L2 words and their conceptual representation. As L2 proficiency increases, links between L2 words and conceptual representation, gradually lessening the need for the L1 to act as a mediator—this means that for a proficient L2 user, the links between the L2 lexicon and the conceptual store are as strong as, or almost as strong as, the corresponding L1 links.

Although several experimental findings have been able to support certain features of the RHM since its inception, it does have its drawbacks. For example, it is now widely accepted that when an L2 word is acquired it could it be stored as a separate lexical and conceptual entity as opposed to sharing a conceptual representation with an L1 word (Brysbaert & Duyck, 2010). Additionally, translation priming studies have shown that it is possible for bilinguals to activate the L2 lexical representation of an item without accessing the L1.

Brysbaert and Duyck (2010) go on to identify and explain some of the major issues with the RHM. Firstly, more recent research that looked at cross-linguistic transfer (e.g. Van Heuven, Schriefers, Dijkstra & Hagoort, 2008; Van Heuven, Dijkstra & Grainger, 1998; Thierry and Wu, 2007) has found that word representations in different languages compete in a manner very similar to word representations in the same language—this indicates a strong likelihood of the lexicon being a single, integrated one instead of the separate ones stipulated by the RHM. Secondly, according to the RHM there are strong links between L1 and L2 words thus predicting a strong translation effect from L2 words to their L1 targets. There is plenty of evidence for L1 to L2 translation priming effects ( Gollan, Forster & Frost, 1997; Jiang & Forster, 2001), but it is barely any evidence that this translation priming exists especially if both languages have different orthographies. Thirdly, although the simplicity of the RHM is appealing, in reality many translations are not simple one-to-one mappings since words can have different translations based on the context. Additionally, words and their translations have synonyms thus implying that the lexicon contains a network of intricate connections rather than a set of one-to-one mappings. Lastly, the RHM is based on the conception that semantic information in the lexicon is language-independent which is contradicted by evidence from studies of bilingual memory representation. For example, studies by Sahlin, Harding, and Seamon (2005) and Marian and Neisser (2000) document the importance of language-specific cues for activation of memory traces in the corresponding language.

As seen in this section, the results of recent empirical research suggest that the conceptualisation of the bilingual mental lexicon must move beyond the RHM in the light of more recent research in bilingual word recognition and processing. Brysbaert and Duyck (2010) recommended that an advantageous way to do this would be to incorporate an additional language into an existing model of monolingual word processing.

*2.7.3.2 Bilingual Interactive Activation Model*

Proposed by Dijkstra and Van Heuven in 1998, the Bilingual Interactive Activation Model (BIA) borrows its basic form from the Interactive Activation model (McClelland & Rumelhart, 1981) and proposes a bottom-up manner of word recognition, from letter features to letters to words. It assumes that both of a bilingual's languages are automatically active during word recognition and processing and that the bilingual lexicon is an integrated one, instead of a separate L1 lexicon and L2 lexicon, thus supporting the non-selective view of language representation and access. Unlike monolingual models of word recognition, the BIA model requires a mechanism for a word to be selected correctly in the intended language: this is the primary way that the BIA differs from hierarchical models is the introduction of a fourth level of node—the language nodes. These language nodes suppress words in the non-target language, not in the earliest stages of activation but in the later stages of word selection. Van Heuven and Thomas (2005) note that an important aspect of the BIA model is that all the nodes are interconnected at the word level—this is called lateral inhibition and it means that words from both languages can inhibit or have an effect on the other's activation.

Dijkstra and Van Heuven (2002) explain the mechanism of the BIA model as follows: when a string of letters is presented to a bilingual individual, this input excites features at each letter position and this in turn excites letters that contain these features, at the same time inhibiting letters for which these features are absent. Next, these activated letters excite words in both languages in which the letter is present at that particular position and all other words, irrespective of language, are inhibited. Following this, activated words nodes belonging to the same language send activation signals to the corresponding language nodes while at the same time, the activated language node sends back inhibitory signals to all word nodes in the other

language. In this manner, the language nodes are activated by the words in the languages they represent and they also inhibit activated words in the other language.

The BIA model has been able to account for a range of empirical findings and a few of them will be listed here. First, the model predicts that neighbourhood density effects[2] take place between languages during word identification and recognition. Studies that investigated the effects of cross-language manipulation (Van Heuven, Dijkstra, Ton & Grainger, 1998; Dijkstra, Ton, Van Heuven & Grainger, 1998) found that during the presentation of a target word, neighbours from both languages are activated. Another phenomenon the BIA model can account for is masked orthographic priming across languages as seen in a study by Bijeljac-Babic, Biardeau and Grainger (1997). In this study, French-English bilinguals were presented with prime and target words from both languages and results indicated that lexical knowledge from both languages had an effect on target word recognition. Both findings— the effects of the neighbourhood density effect and the effect of masked orthographic priming— support the concept of non-selective access and an integrated mental lexicon.

---

[2] In neighbourhood studies, a target word's orthographic "neighbour" is defined as any word that differs from the target words by a single letter (as long as they have the same length and letter position) (Coltheart, Davelaar, Jonasson, & Besner, 1977). Neighbourhood density refers to the number of neighbours a word has.

*Figure 2.2* The Bilingual Interactive Activation model (Dijikstra and van Heuven, 1998, p.200)

However, although the BIA in its original form can account for certain empirical findings, it has its limitations. Perhaps the most obvious limitation is that it assumes that both languages under consideration share orthographic features i.e. it does not account for bilingual mental lexicons in which each language has different orthographic representations and it does not account for phonological and semantic features. Additionally, the representational and the functional features of the language nodes are confounded and there is not much room to account for how linguistic and non-linguistic factors influence word recognition in bilinguals.

*2.7.3.3 Bilingual Interactive Activation Plus Model*

In an effort to incorporate additional representations and processing components in the form of a revised model, Dijkstra and Van Heuven proposed an adapted model, the BIA+ model, in 2002. They describe the BIA model as being "nested in the BIA+ model" (Dijsktra & Van Heuven, 2002, p.181), meaning that the simulations relating to orthographic word recognition are still valid, but the revised model includes lexical, semantic, and phonological features as functions of the language nodes. Dijkstra and Van Heuven (2002) list out the ways the BIA+ model differs from the BIA model: (i) representation and processing of orthographic, phonological, and semantic codes; (ii) representation of interlingual homographs and cognates; (iii) linguistic context effects; (iv) non-linguistic context effects; (v) relationship between word identification and task demands; (vi) stimulus-response binding in lexical decision; and (vii) stimulus-response binding in language switching. To illustrate in detail how the BIA+ model differs from the BIA model, a few of these points will be discussed here.

Like the BIA model, the BIA+ model presupposes that the lexicon is integrated and lexical access is non-selective, but it extends these presuppositions from orthographic representations to phonological and semantic representations. In the process of word recognition, the first stages are identical to the BIA model. With regard to orthography, it follows logically that that is no cross-language activation across language pairs that do not share orthography at all i.e. orthographic activation is language specific. In the process of lexical activation, orthographic representations are activated just before phonological and semantic representations (Ferrand & Grainger, 1993). Based on this, in the BIA+ model once orthographic codes are activated, they begin to activate phonological and semantic representations. The exact moment of activation depends on whether they belong to the L1 or

L2, subjective frequencies, and other factors which imply that L1 codes are likely to be

activated before L2 activation codes. This assumption is termed the "temporal delay

assumption" (Dijsktra & Van Heuven, 2002, p. 183) and it has two consequences: (i) cross-

linguistic effects are expected to be larger from the L1 to the L2 than the opposite direction,

and (ii) if task demands allow responses to faster codes, such as orthographic codes, will

mean that it is possible that semantic and phonological codes can have no effect during

activation.



*Figure 2.3* BIA + Model (Dijkstra & Van Heuven, 2002, p.182)

In the BIA model, the language nodes have two distinct functions—linguistic and non-

linguistic—which operate at the same level of processing. The BIA+ model was revised

based on evidence that suggests the linguistic and non-linguistic functions of language nodes

should be assigned to different levels of processing since they operate in different ways. Specifically, the linguistic function of the language nodes is restricted to language membership representations and they no longer function as language filters. Studies (e.g. Dijkstra, Ton, Timmermans, & Schriefers, 2000) indicate that language information becomes available later on in the bilingual word recognition process (for isolated words), after the stage of word selection and thus language identification does not affect the word selection process. Dijsktra and Van Heuven (2002) suggest that this could be due to how language nodes and words are mapped to each other in the bilingual mental lexicon: each language node is connected to thousands of words, but each word is connected to only one language node. Therefore, if the amount of activation is constant at different levels, the feedback from the word level to the language node will be much larger than the feedforward from the language node to the word level.

Since the BIA+ model makes a distinction between the word identification system and the task-based decision system, it also makes the assumption that linguistic context directly affects word identification and non-linguistic factors affects aspects of the task-based decision system. Linguistic factors in this context are defined as effects associated with the lexical, syntactic and semantic features of the test and non-linguistic factors are defined as effects arising from features such as task instruction and task demands. Most studies examining the mechanism of bilingual word recognition have focused on words in isolation but there are a few that have looked at it in sentence context and Dijsktra and Van Heuven (2002) cite a few of them that support this view of the effects of linguistic factors. For example, in an eye-tracking study with Spanish-English bilinguals, Altarriba, Kroll, Sholl, and Rayner (1996) found that target word recognition interacted with linguistic sentence

context and word frequency interacted with semantic constraint[3], suggesting that lexical characteristics play an important role in word recognition. In a more recent study by Titone, Libben, Mercier, Whitford, and Pivneva (2011) that investigated the effects of semantic constraint on bilingual lexical access, significant effects of cognate facilitation were found for sentences with low semantic constraint.

In summary, experimental studies have shown that the basic architecture of the BIA+ Model captures bilingual word recognition fairly well. A primary strength of this model is that it is clear, specific and empirically testable predictions have made it possible for research to further explore the complexities of non-selective access and an integrated bilingual mental lexicon. However, there is still plenty of room for research to develop a better understanding of the roles of the phonological and semantic representations in the BIA+ word representation system as well as how these representations interact with linguistic and non-linguistic context effects.

*2.7.3.4 Multilink Model*

Drawing from the strengths of the Revised Hierarchical Model (RHM) and the Bilingual Interactive Activation Model (BIA), the Multilink Model was developed by Dijkstra and Rekke (2010) with a particular focus on translation. This localist-connectionist model seeks to explain how word translation happens in bilinguals with different proficiencies during the stages of recognition, retrieval, and production. There are two key features of this model based on the RHM and BIA models: (i) it recognizes that in the bilingual mind, the size of the lexicon in each language varies based on proficiency level and exposure and for L1 and L2, so the links between form and meaning might be different and (ii) lexical access is non-

---

[3] Semantic constraint refers to the degree to which a sentence is semantically biased towards a target word (Titone, Libben, Mercier, Whitford, & Pivneva, 2011).  E.g. 'In the church nativity scene, the *manger* was broken' (high constraint) and 'He built the *manger* using delicate wood' (low constraint).

selective, with differences between phonological, orthographic, and semantic representations. While Brenders, van Hell and Dijkstra (2011) were running simulations during the development of the model, they found that although early bilinguals may have longer processing times, they ultimately have the advantage of the possibility of reaching high levels of proficiency in both languages. Finally, the Multilink model is based on the presupposition that orthographic nodes are responsible for the activation of semantic nodes. Dijkstra and Rekke (2010) go on to explain that this semantic activation is dependent on individual proficiency and word frequency. It takes into consideration the word association links between languages; it also accounts for priming effects at the activation stage for words that are strongly associated with each other. These features of the Multilink model make it relevant to the studies in this thesis since collocations are words that are strongly associated with each other and the congruency effect that will be explored is similar to a priming effect.

*2.7.3.5 Distributed Feature Model*

The Distributed Feature Model (DFM) proposed by de Groot in 1992 is based on the observations that words in one language do not always have precise translation equivalents in other languages and it is well known that often translations only capture approximate meanings (Hofer, 2015). For example, Altarriba (2003) studied emotion words in Spanish-English bilinguals and found that emotion words in Spanish are usually contextualized to a specific situation or episode which is not the case in English. The DFM offers a detailed explanation of how words are stored at the conceptual level based on word type. The main assumption of the DFM is that concrete words and cognates share more features across languages which makes them easier to translate, but abstract words and non-cognates share fewer features thus making them more difficult to translate. The model in Fig. 2.4 shows connections from L1 and L2 words to their conceptual nodes, each of which represents a

word meaning. As depicted in Fig. 2.4, concrete words (e.g. *father* and *padre*) can completely or almost completely overlap (in case they are used in different senses in each language) while abstract words such as *advice* in English and *consejo* in Spanish (*consejo* means advice but can also be translated as council, board, tip) have much less overlap and are more likely to contain language-specific information.



*Figure 2.4* The Distributed Feature Model (de Groot, 1992, p.400)

Evidence to support this model can be found in studies that have shown that bilinguals recognise and translate concrete words faster than abstract words, much like how monolinguals process concrete words faster (De Groot, 1992; Van Hell and de Groot, 1998). Other studies (Altaribba, 2003; Van Hell and de Groot, 1998) suggest that grammatical category also plays a role: nouns behave like concrete words while verbs are more similar to abstract words (Hereida & Cieślicka, 2014).

The main strength of this model is its usefulness in explaining the nuances in differences and similarities of word meanings across languages. Hereida and Cieślicka (2014) give the example of the English word "love" and its Spanish equivalent "amor". In English, although "love" is polysemous, it can be used for people, animals and inanimate objects but in

Spanish, "amor" can be used only for people. These distinctions are possible even for concrete words: "pelota" in Spanish has a high degree of overlap with its English equivalent "ball", but less overlap with an alternative translation "balon"—this is because the Spanish concept "balon" refers to a large, heavy ball which is usually a specialised sport ball, thus all "balons" are balls but not all balls are "balons". The DFM effectively captures and portrays these subtle distinctions that are present at the conceptual level.

Although the DFM takes cross-linguistic features and differences into account, it has a few significant drawbacks some of which are listed by Pavlenko (2009). Firstly, it does not include a developmental component that would adequately explain the acquisition and representation of partial translation equivalents. Secondly, it equates interlingual connection with the degree of shared meaning, but research indicates that other factors affect interlingual connection such as level of activation of both languages, context of acquisition, context of use, similarity of word forms etc. In particular, level of proficiency in both languages is seen as an important factor—while unbalanced for unbalanced bilinguals, connections between translation equivalents may be weaker regardless of shared meaning, for balanced or expert bilinguals, it is possible that even partial equivalents may have strong connections. Finally, there are studies that contradict the main premise of this model and indicate that even concrete words may have partial or completely different conceptual representations (Ameel, Storms, Malt & Sloman, 2005; Malt & Sloman, 2003).

*2.7.3.6 The Shared Asymmetrical Model*

Dong, Gui, and MacWhinney (2005) developed a model that outlines a more dynamic model of L2 vocabulary learning and bilingual performance. In the Shared Asymmetrical Model

(SAM), the L1 and L2 lexicons are linked to each other, to separate L1 and L2 concepts and to common conceptual elements (see Fig. 2.5).This model is based on experimental findings that the authors have used to formulate two tendencies of L2 learners. The first is the convergence tendency: it states that "conceptual differences between a pair of translation equivalents tend to converge in the minds of L2 learners. The more advanced the L2 is, the greater co-effects the two languages produce on the conceptual representations of the two languages" (Dong et al, 2005, p.232). This means that L2 learners are dependent on their L1 at the earliest stages of learning an L2 word, but gradually the conceptual language differences become smaller to the L2 learner and the conceptual representation for the corresponding L1 word is influenced by the L2 conceptual system. The second is the separatist tendency in which L2 learners exhibit the "tendency to maintain the L1 conceptual system in the representation of the L1 word and to adopt the L1 conceptual system in the representation of the L2 word" (Dong et al, 2005, p.233). Both these tendencies are shown to some extent in the process of L2 vocabulary acquisition.



*Figure 2.5* The Shared Asymmetrical Model (Dong et al, 2005, p.233)

L1 = L1 lexical item names, L2 = L2 lexical item names, L1 elements = L1 concepts, L2 elements = L2 concepts and common elements = common concepts.

In this model, common elements are the conceptual elements that have translation equivalents in both L1 and the L2, whereas L1 elements and L2 elements constitute language-specific and culture-specific concepts. In general, the number of common elements is far greater than L1-specific and L2-specific conceptual elements and this is represented in the size of the corresponding circles in Figure 2.5. According to the dynamics of L2 vocabulary acquisition represented in this model, the link between lexical item names (L1 and L2 in Figure 2.5) and the common elements is stronger than the link between lexical item names and language-specific elements (L1 elements and L2 elements), represented by the corresponding thickness of lines in the diagram. Additionally, the link between L1 lexical item names (L1) and common elements is stronger than the link between L2 lexical item names (L2), as shown by the solid and dotted line in the diagram. With the progression of L2 proficiency, the initial link between L2 (lexical item names) and L1 elements (concepts) slowly weakens as the link between L2 (lexical item names) and L2 elements (concepts) grows stronger.

Thus, this model highlights the key role that bilingual lexical memory plays in L2 vocabulary acquisition and makes the claim that conceptual convergence is vital to understanding how the bilingual mental lexicon develops over time.

*2.7.3.7 The Modified Hierarchical Model*

The Modified Hierarchical Model (MHM) proposed by Pavlenko (2009) seeks to act as a transitional model, retaining the strengths of earlier models while exploring new questions and positing new hypotheses. It draws on the developmental progression from lexical to conceptual representation which is central to the RHM and also includes the concept of partial and shared representations which is key in the DFM. Additionally, it seeks to clarify

the nature of conceptual representation proposed in the SAM. The features of the MHM that develop and goes beyond these earlier models will be explained here.

(i) Organization of the conceptual store

In the MHM, conceptual representation is understood to be either fully shared, partially overlapping or fully language-specific, unlike in the RHM which assumes a unified conceptual store. In Figure 2.6, L1 and L2 represent language-specific concepts and language specific-parts of partially overlapping concepts. This recognition of language-specific conceptual storage is a key factor that differentiates the MHM from the DFM.



*Figure 2.6* The Modified Hierarchical Model (Pavlenko, 2009, p.147)

Pavlenko (2009) explains that this language-specific storage has certain implications for bilingual processing: if some linguistic categories are language-specific and culture-specific,

only one of a bilingual's languages would have the required word forms or lexical representations. Thus, to use language-specific lexical concepts of one language in another language, bilinguals usually resort to practices such as codeswitching and loan translation (Pavlenko, 2003; Pavlenko & Driagina, 2007).

In this view, the activation process is context-dependent: it is a two-way interaction between the mind and the environment in which linguistic and social contexts have effects on the conceptualizer by activating concepts and frames linked to one language while inhibiting other concepts and making them less accessible. Studies in cross-cultural psychology have demonstrated evidence of the context-dependent nature of bilingual cognition (Hong, Morris, Chiu, & Benet-Martinez, 2000; Ross, Xun, &Wilson, 2002).

(ii) Conceptual transfer

The second distinguishing feature of this model is the importance given to the phenomenon of conceptual transfer which is based on the differentiation between conceptual and semantic levels of representation. Distinguishing between these two levels of representation allows us to differentiate between the sources of transfer for lexical items and provides more information on the role they play in L2 vocabulary acquisition. While semantic transfer has been extensively looked at in previous models, conceptual transfer can be described as the use of L2 words in accordance with L1 linguistic categories (L1 conceptual transfer) and the use of L1 words according to L2 linguistic categories (L2 conceptual transfer) (Pavlenko, 2009). This leads to the final distinguishing feature of the MHM model which is conceptual restructuring.

(iii) Conceptual restructuring in L2 learning

The MHM model is particularly relevant to L2 learning because it views the main goal of L2 learning as conceptual restructuring and the creation of L1-target-like linguistic categories,

while the RHM regards L2 learning as the development of direct links between L2 words and existing conceptual categories. The MHM incorporates the concepts of implicit knowledge and explicit knowledge which are widely accepted in language acquisition studies but are relatively new with regard to models of the bilingual lexicon. Implicit knowledge is knowledge that learners are unaware of but can be inferred from their performance (Ellis, 2005; Hultsjin, 2007). Explicit knowledge is knowledge that learners possess metalinguistic knowledge of what they have learned and can usually verbalize it when required (Ellis, 2005; Suzuki & DeKeyser, 2017). Studies have shown that L2 learners' performance while drawing on the two types of knowledge can differ: while acquiring vocabulary in L2 classrooms, learners can complete tasks by drawing on word use rules from their explicit knowledge when they are not under time pressure (DeKeyser, 2008; Nazari, 2013). However, in time-constrained spontaneous communication, they find it more difficult to do this and often make the mistake of overgeneralizing L2 words or using the wrong word in a given context (DeKeyser, 2008; Geeslin, 2003). The MHM contends that to achieve target-like word use abilities in the L2, learners have to move from explicit vocabulary knowledge to implicit vocabulary knowledge. This model contributes to the understanding of the bilingual mental lexicon in two primary ways: firstly, the differentiation between semantic and conceptual transfer and the integration of explicit and implicit knowledge in L2 vocabulary acquisition helps us to understand learner errors from a new perspective. Secondly, it enables us to understand why negative conceptual transfer is harder to correct than negative semantic transfer, since conceptual transfer involves conceptual restructuring.

**2.8 Summary of models of the bilingual lexicon**

In this section, the key models that have attempted to explore how the bilingual mental lexicon is represented and accessed were discussed in some detail. It is evident from empirical research that there are complex interactions between the different levels of the bilingual lexicon and how they are organized. Although significant progress has been made in the understanding of bilingual lexical representation and access and its dynamic mechanisms, methodological constraints have to be overcome to present a clearer picture and to refine our understanding of the finer aspects of the bilingual mental lexicon. The studies in this thesis aim to examine how the influence of L1 collocational processing in the L2 can be explained by these models since they provide a theoretical model for our current understanding of lexical acquisition, storage and processing. Although the focus of these models is on single lexical items, they take word associations into consideration as well. With the field of lexical studies growing to acknowledge that formulaic language is an important part of how language is acquired, it is important to look at models of lexical storage and processing in this context. The next section will review the literature on formulaic languages and collocations in detail.

**2.9 Formulaic language**

To understand the nature of collocations, it is necessary to first look at the group of language elements to which they belong: formulaic sequences. In recent years, there has been a significant increase in research on the acquisition, processing and teaching of multiword

sequences and formulaic language within the fields of language acquisition and pedagogy. Although it resists a strict definition, formulaic language can be described as an umbrella term for the different types of multiword units found in written and spoken discourse. Henriksen (2013) provides an overview of the basic types of formulaic sequences with examples: idioms (*when life gives you lemons make lemonade*), figurative expressions (*to freeze to the spot*), pragmatic formulas (*have a nice day*), discourse markers (*let me see now*), lexicalized sentence stems (*this means that…*), and collocations (*rough crossing, remotely clear*).

Research indicates that at least one-third of language is composed of formulaic elements (Conklin & Schmitt, 2008). Schmitt and Carter (2004) list some of the ways that formulaic language can be used: to state a common maxim (*look before you leap*), to express a concept or idea (*raining cats and dogs, a bitter disappointment*), common discourse markers (*nice weather today, it's good to see you*) etc. Pioneering work by Pawley and Syder (1983), Nattinger and DeCarrico (1992), and Lewis (1993) all served to draw the attention of language teachers and researchers to the frequency of formulaic language and its importance for language acquisition (Henriksen, 2013). Schmitt (2010) emphasises the importance of focusing on formulaic language in vocabulary research and lists the following reasons for doing so: (i) although the exact estimates vary, it is agreed that both written and spoken discourse contain large percentages of formulaic language; (ii) the existence of a large percentage of formulaic language in discourse implies that proficient language users have access to a large number of formulaic sequences; (iii) formulaic language is a varied phenomenon that serves different communicative purposes as mentioned above; and (iv) the use of formulaic language enables fluency—apart from the belief that formulaic language helps compensate for online cognitive demands, there is growing evidence that formulaic language is processed more quickly than non-formulaic language, although research shows

that L2 learners have difficulties acquiring it (Conklin and Schmitt, 2008; Underwood, Schmitt, and Galpin, 2004).

### 2.9.1 Identifying Formulaic Language

Due to the diverse and widespread nature of formulaic language, defining and identifying it is a particularly difficult task and is perhaps the most significant obstacle in formulaic language research. Although vocabulary researchers have proposed various comprehensive and suitable definitions, the following definition given by Wray (2005:12) is suitable for the purpose of defining formulaic language as accurately as possible:

*"a word or word string, whether incomplete or including gaps for inserted variable items, that is processed like a morpheme, that is, without any recourse to any form-meaning matching of any subparts it may have."*

Moon (1997) proposes three criteria for distinguishing formulaic language or multiword units from other kinds of word strings: institutionalisation, fixedness, and non-compositionality. Institutionalisation refers to the degree to which a word string is considered conventional in a language by a language community i.e. how often it recurs in language use. Fixedness is the degree to which the order of words in a multiword sequence is fixed; for example, *it's raining cats and dogs* is frozen as it is and cannot be changed. Fixedness also takes into account if a formulaic sequence varies in its components to give different meanings—for example, *another kettle of fish* and *a different kettle of fish* are accepted variations with a difference in meaning, but *on another hand* is not an accepted variant of *on the other hand*. Non-compositionality refers to the degree to which a formulaic sequence cannot be interpreted according to its individual words i.e. the meaning of the whole sequence is not the sum of the meanings of its constituents. For example, break a leg (which means good luck) has nothing

to do with actually breaking a leg. These three variables occur in different combinations and in different degrees in all formulaic sequences.

## *2.9.2 Approaches to identifying formulaic language*

Based on different research purposes and the kind of data available, there have been various approaches to identifying formulaic language. Although there is no standard framework for this purpose, Schmitt (2010) observes that it is possible to broadly split these approaches into four major trends which will be discussed in this section.

The most common approach to identifying formulaic language is through corpus statistics and this is usually done by identifying sequences that commonly recur in a chosen corpus based on a set of previously identified frequency ranges—this process has the very convenient advantage of easily being automated. As what is frequently recurring can vary drastically across different subsections of language users (even for the same language), this approach is particularly useful when studying different groups of language users and learners. With advances in corpus linguistics over recent years, there are numerous corpora now available such as the following: the International Learner Corpora of English (ICLE) is a growing corpora of essays by learners from over 16 language backgrounds; the widely cited British National Corpus (BNC) consists of a 100 million words of text samples drawn from a range of genres, curated to represent a sample of spoken and written British English at the time of its creation; and the Corpus of Contemporary American English (COCA) contains more than 560 million words, with an equal divide between academic and fictional texts, spoken language, as well as newspapers and magazines. Using the frequency calculators available for these corpora, it is possible to identify how many times word strings of different lengths (two-word, three-word etc.) appears in each corpus. The word strings identified by this method do not have to be structurally complete e.g. *once upon a time* would be

recognised as a formulaic sequence. This frequency measure alone can be used to determine which words are collocates of each other. For example, in the BNC, the word string *firm foundation* occurs 32 times whereas *strong foundation* appears only 8 times. From this, it can be assumed that *firm foundation* is more likely to be a collocation than *strong foundation*. However, there are two significant issues with using only simple frequency to identify formulaic language: firstly, function words are the most frequently recurring words in both spoken and written discourse so word strings like *at the* and *in the* show up as having extremely high frequencies, even though they co-occur simply as a matter of structural and lexical chance and not as formulaic language. Secondly, some collocations have a very low frequency but are highly restricted and therefore quite salient e.g. In order to circumvent these issues, researchers use strength of association measures which will be discussed in Section 2.12.1.

Another approach to identifying formulaic sequences adopted by some researchers is to look at them in terms of how transparent or substitutable the constituents are. For example, the verb *drive* collocates only with some modes of transport such as *car, bus, truck* but not with others such as *motorbike* or *bicycle.* The main problem with this approach is that it cannot be automated and is therefore labour-intensive—it cannot be used to analyse large corpora and is useful only for small-scale studies. Additionally, it also has an inescapable element of subjectivity and evaluations of formulaic sequences would require multiple assessors to be valid and reliable.

Two slightly lesser-known approaches to the identification of formulaic language are the acquisition approach and the psycholinguistic approach. In the acquisition approach, the focus is on the lexis that children repeat in the process of language acquisition. In the psycholinguistic approach, it is assumed that formulaic language is stored in chunks in the mental lexicon. Bybee (2002) attributes phonetic reduction in collocations with a high degree

of association to the automatization that results from the repetition of sequence, which then results in the fossilization of these changes as multiword chunks. However, the problem with this approach is that defining a formulaic sequence by what is stored in the individual lexicon is tricky due variations in formulaic sequences. Although in theory formulaic language is fixed, in reality there are variations due to factors such as cultural differences, varying lexical components and substitutability of verbs.

### 2.9.3 Acquisition of Formulaic Language

While it is known that formulaic language is essential for fluent communication, it may also aid the process of further language learning. The results of studies in L1 acquisition (Peters, 1983; Wray, 2000) have led researchers to propose that for children, the acquisition of these unanalysed chunks of language provide the basis for language development and facilitate learning of language components and grammar. For example, a child uses phrases like *I wanna go home* even though such phrases may be beyond their current linguistic and grammatical knowledge. The child can then break the phrase down into its constituents and gain an understanding of how syntax works in language—this can be done with other phrases as well if they have useful recombining components (Wray, 2000).

The acquisition of formulaic language seems to be a rather slow process for L2 learners, and when it comes to active usage, even highly advanced L2 learners have been shown to rely on a narrower set of formulaic sequences than their native counterparts (Durrant & Schmitt, 2010). Studies in L2 acquisition have shown that the acquisition of formulaic sequences is an area in which L2 learners tend to lag behind native speakers, even in comparison to other linguistic aspects (Kuiper, Columbus, & Schmitt, 2009; Conklin & Schmitt, 2008). Other studies have shown that only very proficient L2 learners—often those who have been immersed in their L2 community—achieve a level of formulaic language competency that

resembles that of native speakers. Yamashita and Jiang (2010) found that with an increase in language proficiency, L2 learners' knowledge and usage of formulaic sequences increased. Irujo (1993) suggested that L2 learners receive less input of a particular kind of formulaic language, idioms, and this leads to a lack of idioms in learner output. In later studies on collocations, Durrant and Schmitt (2009) found that frequency is an important factor in collocation acquisition— in this study, L2 learners produced frequent collocations but not infrequent ones.

Schmitt (2010) observes that the field of formulaic acquisition is turning towards pattern-based models and of language acquisition and construction grammar in an attempt to understand how exactly formulaic language is acquired by language learners. These theories posit that individuals learn language implicitly based on their ability to extract patterns from the input they receive. For example, the differences subject-verb agreement for different persons (*I want, you want, he wants, they want*) are acquired by the individual observing and extracting the patterns from language input. Of course, at later stages the individual may be able to provide a grammatical rule for verb conjugations, but the initial acquisition is based on patterns of language use rather than on explicit knowledge of language rules. This pattern-based model also explains how morphemes combine to form different words e.g. *un-deni-able, un-desire-able.* It also can account for collocational association found in language use e.g. it is by statistical learning through input rather than by explicit instruction that a proficient language user knows that *blonde* collocates with *hair* and not with *paint*, even though this association can be described as arbitrary and there is no semantically logical explanation for this.

*2.9.4 Processing of Formulaic Language*

Formulaic language has been shown to have processing advantages, primarily because it makes use of a relatively abundant resource (long-term memory) to compensate for working memory which is a limited resource. In this way, the mind uses long-term memory to store prefabricated chunks of language which can be readily used in language production instead overloading the working memory, which would otherwise be required to recall different lexical and syntactical rules during the process of language production (Conklin & Schmitt, 2008). Although this explanation was first postulated as an assertion, recent research using different methodologies (e.g. Arnon & Snider, 2010; Conklin & Schmitt, 2008) is now converging to support this claim that formulaic language has a processing speed advantage (Conklin & Schmitt, 2012). The earliest studies that investigated this advantage looked into speech production—the first salient one was done by Dechert in 1983. Dechert analysed the spoken output of a German speaker of English and found that some parts of her output were fluent and he labelled these parts "islands of reliability" i.e. formulaic language. Based on this finding, Dechert contended that these islands of reliability are essential for the smooth planning and execution of speech production. The results of other studies that examined the role of formulaic language in speech production support this claim: Bannard and Matthews (2008) found that monolingual young children showed a greater sensitivity to high frequency formulaic sequences than similar sequences of a lower frequency and studies by Kuiper (2000; 2004) showed that people who are required to produce fluent speech under time pressure (e.g. sports commentators) rely heavily on formulaic language to make this possible.

## 2.10 Processing of Figurative Language

Many of the studies that examined the mental processing of formulaic sequences in native speakers have focused on formulaic sequences in the form of metaphors, idioms, proverbs and other figurative language. Gibbs, Bogdanovich, Sykes and Barr (1997) conducted a study with idioms and control phrases in a story context and determined that the idioms were read more quickly than the control phrases. More recent studies using eye-tracking, such as Underwood, Schmitt, and Galpin (2004), have found that for formulaic sequences, both native speakers and proficient non-native speakers fixated less on the terminal words than for non-formulaic control phrases because by the time they got to the terminal word of the formulaic phrase it was highly predictable and they did not need to spend as much time on it—this is evidence of a processing advantage for formulaic language. Conklin and Schmitt (2008) tested the reading times of formulaic sequences (both idiomatic and literal) for native speakers and L2 speakers—even though the reading times of the L2 speakers was longer than those of the native speakers, both groups read the formulaic sequences more quickly than they read the control strings which were non-formulaic sequences. From these findings and the findings of other similar studies, Conklin and Schmitt concluded that connectionist accounts of language acquisition are a probable explanation for how the words of formulaic sequences become associated with each other due to repeated input. These connectionist/emergent theories of language acquisition have their roots in pioneering work on neural networks and computational models done by neuroscientists and computer scientists in the 1940s and 1950s (Gasser, 1990). These accounts suggest that the extraction of patterns from language input is essential for the process of language acquisition. Since formulaic language is relatively frequent and salient in written and spoken discourse, it is highly probable that formulaic sequences are included in these extracted linguistic patterns.

Effectively, each time a person is exposed to a formulaic sequence it reinforces existing knowledge of that sequence that exists in the mental lexicon.

However, it must be noted that in Siyanova-Chanturia, Conklin and Schmitt's (2011) eye-tracking study, there was no evidence that the non-native speakers processed idioms faster than the matched literal sequences—in fact, the reading times for the idiomatic sequences were longer than those for the literal sequences, even though the idioms chosen were common ones that were likely to be known to the non-native speakers. Citing these conflicting results in their review of research on the processing of formulaic language, Conklin and Schmitt (2012) suggest that idioms may not be best suited for drawing conclusions on the processing of formulaic sequence— this is because idioms do not occur very frequently in language use and language learners and children may not have much exposure to them. Additionally, idioms have varying degrees of transparency and are often ambiguous because they could have both a figurative and literal sense. Considering these factors, other research on other kinds of formulaic sequences could give us a more comprehensive idea of how formulaic language is processed.

**2.11 Processing of other formulaic language**

In a comprehensive overview of the processing of formulaic language, Carrol and Conklin (2020) analysed recent research and found that regardless of differences in the properties of different kinds of formulaic language, all types of formulaic language are processed more quickly than their non-formulaic counterparts. In a set of two studies, Bod (2000; 2001) examined the processing of frequent three-word (*I love it*) SVO sentences and low-frequency control sentences (*I keep it*), all matched for lexical frequency and complexity, in native speakers and found that the participants responded more quickly to the high-frequency

sentences, thus showing that frequency is an important factor in the processing of word strings and is not limited to single lexical items.

Other important factors that have been shown to have an effect on individual word recognition and processing are frequency and word length. Tremblay and Baayen (2010) investigated processing times for four-word formulaic sequences using phrase recall and electrophysiological measures and found that the participants showed faster processing times for the sequences with higher frequencies. In their eye-tracking study Siyanova-Chanturia, Conklin and van Heuven (2011) examined the effect of frequency on the processing speeds for three-word binomial strings (*kith and kin*) embedded in sentences for native and non-native speakers. Using mixed effects modelling, they found that across a range of proficiency levels, both groups had shorter processing times for more frequent formulaic sequences when compared to the less frequent ones—this is in line with the overwhelming evidence found in the literature in this field.

## 2.12   Role of L1 in the Processing of Formulaic Language

As explained by Van Lecker Sidtis (2015) language users need to know complex details of form, meaning, and use of formulaic language to use it correctly and efficiently. For L2 learners, this process of knowing and using formulaic language correctly is complicated by the fact that they have access to an existing repertoire of formulaic language in their L1 which may partially overlap with L2 formulaic language in various ways (Carrol & Conklin, 2019). For example, Carrol and Conklin provide examples of English and French idioms which overlap to different degrees: (i) almost full overlap is observed between the English idiom "to throw money out of the window" and the French idiom "jeter l'argent par les fenêtres" is translated as "to throw money out of the windows" and both versions have the

same meaning; (ii) the English and French idioms for expensive things overlap in meaning but take on different forms—in English it is "costs an arm and a leg" and in French you would say *coûter les yeux de la tête* which is translated as "costs the eyes in your head"; (iii) an instance when English and French use different phrases to express the same idea the English expression "to feel blue" which means feeling sad or depressed and the French equivalent is *avoir le cafard* which is literally translated as "to have the cockroach". Thus from these examples, we can see that there are different degrees of overlap of formulaic sequences between different languages and navigating these complex differences and similarities can be challenging for L2 learners to navigate. A number of studies have been conducted to investigate how exactly the L1 influences the processing of L2 formulaic language and these studies will be reviewed in this section.

In a study designed to examine the processing of English idioms and transliterated Chinese idioms in English by Chinese-English bilinguals, Carrol and Conklin (2014) used self-paced reading to measure the response times of these bilinguals as well as native English speakers (control group) to the aforementioned types of idioms. The participants were presented with a combination of four sets of items: English idioms, English control phrases, translated Chinese idioms in English and translated Chinese control phrases. The English idioms were chosen from the *Oxford Learner's Dictionary of Idioms* (Warren, 1994) and the Chinese idioms were chosen from the *Dictionary of 1000 Chinese Idioms* (Lin & Leonard, 2012). For the Chinese idioms, only those which had a literal translation that is plausible in English and with an identical word order in English were selected. The chosen Chinese idioms all had a monosyllabic final word that had a translation equivalent in English. The control items for both English and Chinese idioms were created by replacing the final word with a plausible alternative, with half of the final words being nonwords. The participants were required to complete a lexical decision task with the first part of the collocation acting as the prime and

the last word of the collocation acting as the target. The participants were first presented with the primes and then with the target word—they were required to press "Yes" if they thought the word was real and "No" if they thought was not real. The reaction times were recorded by the self-paced reading software. The Chinese participants were also given a vocabulary test after the LDT. Analysis of the results showed that both native speakers of English and the Chinese participants responded significantly more quickly to the target words of the idioms in their own languages (the transliterated idioms in the case of the Chinese speakers) than the final word of the control phrases. The Chinese speakers had significantly faster reaction times to the translated Chinese idioms than the English idioms, even though the Chinese idioms were presented in English. With regard to the English idioms, the Chinese speakers did not show any significant differences between the English idioms and the control phrases—in this study, the Chinese-English bilinguals did not show any processing advantage for formulaic sequences in the L2 possibly because the Chinese speakers were not familiar with the English idioms. Thus the results of this study showed that for these learners, the L1 was activated even when they encountered L1 transliterated idioms in an unfamiliar form in the L2.

Carrol and Conklin (2015) followed up the previous study by investigating the processing of translated Chinese idioms to determine whether L1 idioms show a priming effect when encountered in context in the L2 by Chinese learners of English. The researchers used eye-tracking as a tool to tap into automatic processes that take place during reading. The first experiment in this study examined whether the local lexical context provided by an translated idiom in the L2 was enough to facilitate priming for the final word. The Chinese idioms were selected from the *Dictionary of 1000 Chinese Idioms* (Lin & Leonard, 2012)—all the selected idioms had final characters with a translation equivalent in English and were judged to be highly familiar in the original Chinese from by native Mandarin speakers. The English idioms were chosen from the *Oxford Learner's Dictionary of English Idioms* (Warren, 1994) and a

set of these idioms were judged to be highly familiar by native English speakers. Similar to the previous study, the chosen idioms were used for the stimuli along with control idioms in both English and Chinese in which the last word of the idiom was replaced with a plausible alternative. In this study, the idioms were embedded in sentences that supported the figurative meaning of the idioms. The participants were required to read each sentence on a computer screen and press a button to go to the next sentence. Half of the sentences were followed by a yes/no question to encourage them to pay attention and to check comprehension. The reading times were measured using an eye-tracker. The reading times for the final word of each idiom were the focus of the analysis, based on the assumption that if the participants was familiar with the idiom and had stored it as a whole unit the processing time for the last word would be considerably shorter. Results showed that unlike the native speakers, the L2 Chinese learners did not read the final word of the English idioms any faster than the English control phrases. For the translated Chinese idioms, results showed that the Chinese learners read the last word significantly faster than the last word of the control items. This effect was clearest in the early measure of first fixation, suggesting that the L1 influence is strongest in the early stages of processing, although the effect on late reading measures such as total time were also significant. These results led the researchers to conclude that there is evidence of L1 influence even L2 speakers encounter L1 idioms in an L2 translation context.

A study by Carrol et. al (2016) investigated the influence of the L1 on the processing of idioms in advanced L2 Swedish learners of English using eye-tracking. English native speakers and the L2 Swedish learners were presented with English idioms, translated Swedish idioms, English-Swedish congruent idioms and matched control phrases. This allowed them to compare the influence of Swedish on English idioms when there was an equivalent in both languages and when there wasn't. The English idioms were chosen from a variety of idiom dictionaries and previous stimuli lists used by the authors, and the same

procedure was followed for the selection of the Swedish idioms. From the list of English idioms, native speakers of Swedish determined which idioms had equivalents in English and these were selected as the congruent idioms. Short sentence contexts were then created for each of the idioms. Participants were presented with a set of these sentences on a computer screen while their reading times were measured with an eye tracker. They were required to press a button to go to the next sentence and half of the sentences were followed by yes/no questions to ensure the comprehension was taking place. The reading times for the entire idiom as well as the final word of each idiom were recorded. Swedish participants were also required to take a receptive vocabulary test to provide an idea of their vocabulary sizes. The results of the vocabulary test showed that the Swedish participants were a fairly homogenous group of advanced English users. With regard to the eye movement measures, the researchers compared the eye movement measurements for the final word of each phrase and for the phrase as a whole. Results showed that the Swedish participants exhibited an overall pattern of reading all three groups of idioms significantly faster than the control phrases. They read the translated Swedish-only and the congruent idioms faster than they read the English idioms although the differences in reading times were not significant, and they were also more likely to skip the final word of these idioms than the English-only ones. This finding is interesting: the lack of significant difference in reading times between English-only idioms and the other types is most likely because the Swedish participants in this study were advanced users of English and already have access to a repertoire of English idioms in their mental lexicons. Although most literature in the field states that it is difficult for L2 learners to achieve nativelike processing of formulaic language, it appears that L2 users with advanced levels of proficiency are able to do this.  As expected, for the native English speakers there was no difference in the processing of congruent and English collocations, but they had significantly slower reading times for the translated Swedish collocations. In terms

of L1 influence on the processing of formulaic language in the L2, there was no evidence that congruency had a faciliatory effect on the congruent idioms over the Swedish-only ones: there was no significant difference between their reading times for Swedish-only idioms and the congruent ones. Thus, this study shows that with advancing proficiency at high levels, there is less likely to be a congruency effect i.e. L1 influence in the processing of L2 formulaic language.

Although research in this area generally points to L1 activation during the processing of L2 formulaic sequences, there are a few exceptions. For example, Cieślicka and Heredia (2013) conducted a study in which they compared Spanish-English bilinguals reading times for congruent idioms (English-Spanish equivalents) as well as for idioms that express the same idea but take on completely different forms. They found that the participants had more difficulties with the congruent idioms than with the different idioms and took longer to read them as well. They attributed this to L1 activation but concluded that the L1 equivalent had been activated and then had to be supressed, thus slowing down the processing. Interpreting these findings in light of the majority of other studies that have found that the L1 has a faciliatory effect on the processing of L2 formulaic language is not straightforward but the authors point out that there were several factors that affected their results including the varied language dominance of the participants and the differences in opacity and transparency. Overall, the existing evidence indicates that until L2 users reach advanced levels of proficiency, the L1 has a significant influence on the processing of L2 formulaic language. This is most likely in due to the quick activation of L1 equivalents while processing during L2 reading (Conklin & Carrol, 2019), which speeds up processing in the L2.

**2.13 Collocations**

A simple but rather broad definition of a collocation is a frequently recurring two-to-three word syntagmatic unit (Henriksen, 2013). Hunston (2002) describes it as the tendency of two words to co-occur or the tendency of one word to attract another. Wray (2002) observes a critical difference between collocations and other kinds of formulaic language such as idioms: collocations are more "fluid" in nature, whereas as idioms are fixed e.g. *to make a mountain out of a molehill* is a fixed expression, while though *rain* often collocates with *heavy*, it is still associated with other words.  Defining the exact nature of a collocation is a critical issue in the area of collocation studies and there is a range of differing opinions. Researchers have proposed different criteria and structures for identifying and classifying collocations. Overarching all these different classifications, there is now a general consensus that there are two main approaches to defining collocations: the frequency-based approach and phraseological approach. In the frequency-based approach, which is more widely used, collocations are defined as two or more words that frequently occur together. Nurmukhamedov (2015) notes that researchers who adopt the phraseological approach use native-speaker judgement, degrees of restriction, and corpora to identify collocations (e.g Nesselhauf, 2003) with a focus on the semantic and lexical restrictions of each word in the combination, whereas researchers who endorse the frequency-based approach rely only on corpus evidence (Henriksen, 2013; Shin & Nation, 2008; Webb, Newton & Chang, 2013). In a recent study on collocational links, Wolter and Gyllstad (2013) merge both these approaches to give a comprehensive definition of collocations as follows: "A collocation is a sequence consisting of two or more words which co-occur more frequently than chance would predict based on the frequency of occurrence of the individual constituent words" (Wolter and Gyllstad, 2013, p.434).  This definition is well rounded and provides an inclusive and precise insight into the nature of collocations. Based on this, most researchers now adopt

both these approaches: first identifying frequently occurring word combinations and then eliminating those that do not fit their analysis criteria (Henriksen, 2013).

### *2.13.1 Collocational Strength: Strength of Association Measures*

To overcome the problems of using simple frequency as a measure of collocational strength, one of the most commonly used measures of strength of association for collocations is the Mutual Information (MI) score. Schmitt (2010:130) describes MI score as "a measure of how much one word tells us about the other". This means that when a word pair has a high MI score, if one of the words is present, there is a high probability that the other member of the pair is nearby. Hunston (2002) explains the MI score as a score that compares the actual co-occurrence of two words (observed occurrence) with their expected co-occurrence if the words in the corpus appeared in a random order. In a given corpus, the MI score is based on the number of times the words occur together versus the number of times the words occurred separately, and this can be done automatically in most corpora. For example, although the frequency of *cloven hooves* is low in the BNC, the MI score is high (16.1) because when the word *cloven* appears, it his highly likely it will be followed by *hooves*. On the other hand, the MI score for *rainy day* (8.8) is about half of that of the *cloven hooves* score even though it occurs much more frequently in the corpus. This is because even though the word rainy is quite often followed by day, it is also likely that it could be followed by other words such as weather, morning, afternoon, season etc. In terms of computation, the MI score for each collocation is the Observed occurrence divided by the Expected, converted to a base-2 logarithm (Hunston, 2002). As mentioned, in most cases this calculation can be done automatically in the chosen corpora.

In collocational research, there is a consensus that the threshold for statistical significance is an MI score of 3 (Hunston, 2002). However, using only MI scores as a measure of collocational strength is not advisable since it is more useful for ranking collocations based on their collocational strength. To remedy this, the t-score is often used in conjunction with the MI score to give a more robust assessment of collocational strength. While the MI score accounts for strength of association, the t-score takes frequencies into account and gives the confidence with which we can assert there is an association between the two words i.e. that it's not just a matter of chance. In terms of computation, the t-score for each collocation is calculated by subtracting the Expected occurrence from the Observed occurrence and dividing the result by the standard deviation (Hunston, 2002). In general, these are the two most commonly used measures of collocational strength—while t-score is a better measure for frequent collocations (*every day*), the MI score is useful for collocations with stronger links, but which appear less frequently in the corpus (*cloven hooves*).

### 2.13.2 Processing of Collocations

Even with the recent surge in research on the processing of collocations, the underlying mechanism of this process remains only partially understood. In one sense, the processing of collocations is different from that of other formulaic language because unlike other formulaic language (idioms, metaphors, etc.), they usually do not pose comprehension difficulties even if some degree of figurative interpretation is required (Nesselhauf, 2005; Henriksen, 2013). For many collocations, if the meanings of the constituent words are known, it is not difficult to deduce the meaning of the whole collocation. Wolter and Gyllstad (2014) suggest that this could affect the way and native speakers and L2 learners process collocations: for L2

learners, this could mean that they tend to process collocations word-by-word rather than recognising them as a whole unit.

This view was most notably argued by Wray (2000, 2005) who contended for the processing of collocations, native speakers rely mainly on their knowledge of the meaning assigned to the whole collocational unit while L2 learners rely mainly, perhaps even exclusively, on their knowledge of the meaning of each individual word. Wray (2002) illustrates this argument with the collocation *major catastrophe.* In this example, Wray explains her position that native speakers would process this as one unit with a singular meaning and do not process the meanings of *major* and *catastrophe* separately—an idiomatic way of processing the collocation. In contrast, L2 learners would decompose the collocation into individual words during processing: this means that they would break down major catastrophe into two words meaning big and disaster and store them separately. Consequently, when they need to express this idea in future, they would have no access to the collocation major catastrophe that they had previously encountered but would instead word pairing with a similar meaning as equally plausible for use in that situation.

However, with advances in research on collocational processing, this view has come under some criticism because it is now believed that not all L2 learners and users take this "break down" approach to processing collocations—there are other factors that have been found to influence how L2 learners process collocations. Wolter and Gyllstad (2013) found that advanced L2 learners showed sensitivity to the frequencies of L2 collocations—this suggests that advanced L2 learners do notice recurrent patterns in the language they encounter and the distributional frequency does play a role in how they process collocations. This finding was supported by the results of another study—Durrant (2014) found that the processing of collocations in L2 learners has a positive correlation with collocational frequency. In a study that examined the rate of acquisition of collocations, Gonzalez-Fernandez and Schmitt (2015)

found that collocational frequency was an accurate predictor of collocation acquisition for learners at different proficiency levels. There is almost always a correlation between the processing of the collocations and the proficiency level of the learners—with an increase in proficiency level, learners are more likely to acquire and process collocations quickly. For example, in an eye-tracking study, Siyanova-Chanturia, Conklin and van Heuven (2011) found that native speakers and advanced L2 learners were sensitive to collocational frequency, but the intermediate L2 learners did not. These studies suggest that unlike Wray's (2000) argument, with gains in proficiency, L2 learners acquire a significant collocational processing advantage. Additionally, in a study on language exposure and collocational processing, Durrant and Schmitt (2010) found a learning effect for collocations and this led them to conclude that any differences between L1 and L2 collocational processing are probably due to a lack of exposure.

## 2.14 Influence of L1 on the Processing of Collocations

Collocations can vary quite considerably from language to language (Wolter & Gyllstad, 2011), and it has been observed that this variation is because of arbitrariness (Pawley & Syder, 1983) i.e. there is often no logical or grammatical explanation for why certain words collocate with each other and others do not. With this variation from language to language, the influence of the L1 on L2 collocational acquisition and processing is of particular interest.

Early research that looked at the L2 mental lexicon mainly accounted for collocational responses as only one kind of response on word association tests (e.g. Meara 1982; Wolter 2001, 2002). These studies, among others, focused primarily on interlexical links between the L1 and L2 mental lexicons and very little attention was given to intralexical links in the L2 lexicon. However, the influence of the L1 on the acquisition of collocations (intralexical

links) is an area that has been attracting a great deal of scholarly attention over the last decade. Although there are different explanations for the influence of the L1 on L2 collocation acquisition, the understanding of how this influence works is still in early stages. This influence of L1 is usually examined by studying the differences in the production, reception, and processing of congruent collocations (collocations that have a direct equivalent in the learners' native language) and incongruent collocations (collocations that do not have a direct equivalent in the learners' native language).

Quite a few studies on the influence of the L1 on the production of L2 collocations have shown that L2 learners tend to rely on a narrower set of collocations than their native speaker counterparts and the errors that they make can usually be attributed to L1 influence. For example, Nesselhauf (2003) studied the production of collocations and looked at the effect of the L1 on wrongly produced collocations. The researcher analysed English essays of 32 German-speaking final year university students. The selected essays were non-technical and argumentative and had an average length of 500 words. All verb-noun collocations were extracted from the essays and sorted into correct and incorrect collocations based on native speaker judgement as well as input from dictionaries such as *The BBI Dictionary of English Word Combinations* and the *Oxford Dictionary of Current Idiomatic English*. About a quarter of the collocations were identified as incorrect (quite evenly distributed over the essays) and the collocations were further sorted into nine categories based on the type of mistake that rendered the collocation incorrect. The most frequent error was the wrong choice of verb and the lest frequent error was the syntactic structure of the collocation. An interesting finding in this study was that not only did the L1 seemed to play a role in every type of incorrect collocation that the participants produced, its influence in each type of mistake is of remarkably similar strength. After translating the incorrect collocations into the L1 (German), the researcher found that the participants had made verb, noun, preposition, article and usage

errors at least partly due to L1 influence; it also revealed that regardless of the degree of restriction of a collocation (i.e. how tightly bound together the words are in a collocation), incongruent collocations are more difficult for learners to produce. Analysis of the correct collocations also revealed that a significantly high number of them were congruent with German collocations, thus lending support to the notion that the L1 can have a faciliatory effect on how collocations are processed and produced.

Yamashita and Jiang (2010) presented verb-noun and adjective-noun collocations to Japanese EFL learners, and native speakers of English and asked them to judge each collocation on whether it would acceptable in English. The researchers found that while there was no difference in processing times or error rates between congruent and incongruent collocations for the native speakers (as expected), for both ESL and EFL Japanese learners the error rates were higher for the incongruent phrases. Additionally, the EFL learners had slower reaction times for the incongruent phrases but this was not the case for the ESL learners. They also used a cloze test with the Japanese speakers to determine the difference in levels of L2 proficiency between the EFL and ESL learners—there was a significant difference in mean scores between both groups, with the ESL learners outperforming the EFL learners. This could be because although they had approximately the same length of English formal education, the ESL learners had a lot more exposure to English outside the classroom. For the main experiment involving the phrase acceptability task, the participants were presented with the collocations on a computer screen using a type of self-paced reading programme to measure their processing times. Firstly, the results of the phrase acceptability task suggest that both L1 knowledge and L2 exposure play a role in the processing of collocations i.e. collocations that also exist in the L1 are processed more easily in the L2, possibly because less cognitive effort is required. Secondly, these results show that incongruent collocations remain difficult for L2 learners even at higher levels of proficiency since acquiring

incongruent collocations requires a very high level of exposure. Finally, for the ESL learners there was a difference in error rates but no difference in processing times—this could be because once L2 collocations represented in the lexicon, they are processed independently of the L1. This led the researchers to conclude that L2 collocations are processed independently of the L1 lexicon only at the later stages of language acquisition i.e. when the learners become more proficient. This means that the influence of the L1 is greater during the initial stages of language acquisition and gradually subsides in the later stages.

Wolter and Gyllstad (2011, 2013) investigated collocational priming on congruent and incongruent collocations native speakers of English and Swedish learners of English and found similar results to Yamashita and Jiang's study. In their first study, they used a lexical decision task and an additional test of receptive collocational knowledge to assess the participants' reaction times to congruent and incongruent collocations as well as a group of unrelated baseline items. In a typical lexical decision task (LDT), the participant is first presented with an initial word (the prime) followed by an additional word (the target). The participant is required to identify whether or not the target word represents a word in the specified language by pressing a key. The software used was the same used by Yamashita and Jiang in their 2010 study, which was a type of self-paced reading software. The items used in the LDT were congruent and incongruent collocations with the first word of the collocation acting as the prime and the second word of the collocation acting as the target. As mentioned previously, there was also a set of unrelated phrases to act as baseline items. As in a typical LDT task, the participants were asked to press "Yes" if they thought the target word was a real English word and "No" if they thought it was not a real word. After completing the LDT, the Swedish learners of English were required to complete a pencil-and-paper yes/no 100-item collocation test which consisted of the collocations used in the LDT as well as a set of distractors—this was the test used to determine their receptive knowledge of the

collocations they had encountered. They found that the reaction times for the congruent collocations and the incongruent collocations was almost the same but significantly more for the unrelated items. This result, although it was expected, supported the assumption that both congruent and incongruent conditions were equivalent and also that words do indeed prime their collocates. For the Swedish learners of English, the processing times for the target word for the incongruent collocations were significantly more than the same measure for the congruent collocations. A notable finding for these Swedish learners was that although the lexical priming was indeed more for congruent collocations, they were less likely to recognize the L2-only collocations as collocational—this was seen in the results of the test of receptive collocational knowledge that they took. From the results of these two studies, Wolter and Gyllstad concluded that the advantage for congruent collocations may be due to a lexical priming effect: the knowledge of the collocation in the L1 primes their knowledge of the equivalent L2 collocation, i.e. the congruent collocation, thus reducing the processing time for such collocations. Thus, the L1 appears to be providing easier access to L2 collocations which have an equivalent in the L1, which is not possible for L2-only collocations. In the 2013 study, Wolter and Gyllstad investigated the effects of frequency on the processing of congruent and incongruent collocations in native English speakers (control group) and advanced Swedish learners of English (experimental group). The focus of the study was to explore how high proficiency learners of English would process congruent collocations, incongruent collocations and noncollocational items, ensuring that items in all three categories had matched frequencies. They also took into account the frequencies of the L1 translated forms for the congruent collocations. An acceptability judgment task chosen for this experiment was administered using the same self-paced reading programme as the 2011 study. In this task, the participants were presented with one item at a time on a computer screen and asked to press "Yes" if they thought it was a phrase commonly used in English

and "No" if they thought it wasn't. This was a slightly modified version of a typical acceptability judgement task in which participants are usually asked to determine if the items are acceptable or not. The Swedish learners of English also took a receptive vocabulary test in order to provide a general idea of their overall English proficiency. As in the 2011 study, the Swedish learners in this study had significantly shorter reaction times for the congruent collocations in comparison with the incongruent collocations and they also produced significantly more errors on the incongruent items when compared to the congruent ones. In terms of frequency, for both native speakers and the Swedish learners of English the biggest predictor of reaction times was the frequency of the collocations in relation to the English corpus. Wolter and Gyllstad cite other studies (e.g. Durrant & Schmitt, 2010) which also found that high-proficiency language learners are sensitive to frequency effects just like native speakers. The findings from this second study by Wolter and Gyllstad (2013) suggest that advanced L2 learners are sensitive to frequency effects not only at the word level but also at the collocational level. Additionally, it must be noted that collocational frequency was the biggest predictor of reaction times, not word frequency. This indicates that both native speakers and the Swedish learners of English processed the collocations holistically as single units and not as separate words.

Wolter and Yamashita (2018) followed up on these previous studies by conducting a study that examined the effects of word frequency, collocational frequency, L1-L2 congruency and language proficiency on L2 collocational processing in Japanese speakers of English (intermediate and advanced) and native English speakers. The focus of this study was to modify the tasks so that the participants would have to pay attention to the meaning of the collocations and not just the form and also to look more closely at the role of word frequency and collocational frequency in collocational processing. Similar to Wolter and Gyllstad (2013), an acceptability judgement task was used in this study and participants were asked to

determine whether each collocation was commonly used in English or not using the same self-paced reading programme used in the previous studies. Four types of items were used in this study: (i) Japanese-English congruent collocations, (ii) English-only (incongruent collocations), (iii) Japanese-only (incongruent collocations), and (iv) baseline items. Results showed that congruent collocations were processed significantly faster than incongruent collocations by both groups of Japanese speakers of English. This was in line with the previous studies with learners with different native languages and with different methodologies, thus strengthening the conclusion that the L1 influences processing of collocations in the L2. With respect to frequency, the results showed that the native English speakers and the advanced Japanese speakers of English showed a greater sensitivity to collocational frequency than the intermediate Japanese speakers of English. Interestingly, further analysis of the word frequency and collocational frequency with both groups of Japanese learners showed that with increased proficiency there was a shift away from reliance to word frequency to reliance on collocational frequency. This indicates that with increasing proficiency, L2 learners move towards nativelike collocational processing. However, it must be noted that even the most advanced L2 Japanese learners in this study relied more heavily on word frequency than the native English speakers which supports the conclusions of previous studies that even advanced L2 learners have difficulties with acquiring and processing L2 collocations.

Taking the results of these studies together, it is clear that across different contexts and with different groups of learners the L1 influences L2 collocational processing i.e. there is a congruency effect in which the L1 is activated during L2 collocational processing. Congruent collocations are recognized and processed more accurately than incongruent collocations i.e. L1-only or L2-only collocations. Also as some of these studies have shown, once a collocation is registered in the L2 lexicon its L2 frequency influences its processing. The

main methods used to measure the congruency effect in these studies was the measuring of processing times using LDTs and acceptability judgement tasks, Additionally, most of the studies used extra tests to measure vocabulary and proficiency size in order to see if these measures could be used to analyse the extent to which L1 influences L2 collocational processing. With regard to previous research, the single mode is similar to the lexical decision task in terms of presentation mode i.e. each word of the collocation presented one after the other. Thus, the measures and methods used in previous studies have informed the design of the studies in this thesis.

## 2.15 Orthographic Differences between Tamil and English

Since both studies in this thesis look at how Tamil influences the process of reading of English collocations, it is necessary to have a basic understanding of how the orthographic features of both languages may play a role in this process. There is growing evidence that script typology may play an important role in cross-linguistic transfer between L1 and L2 (Geva & Siegel, 2000; Geva, Wade-Woolley, & Shany, 1993). This section will briefly discuss the orthographies of Tamil and English and how they differ from each other.

Tamil belongs to the South Dravidian group of languages which is a language family consisting of 26 languages that are native to the Indian subcontinent. It is primarily spoken in the South Indian state of Tamil Nadu, which is where Study 2 was conducted, but it is also an official language in Singapore and Malaysia where there are substantial Tamil communities. It is spoken widely by Tamil diasporic communities in countries like the United Kingdom, the United States, Sri Lanka, Mauritius, Fiji etc. As noted by Bhuvaneshwari and Padakannaya (2013), Tamil is an alphasyllabic script: the Tamil script is *akshara*-based

which typically represents language at the orthographic or syllabic level. Each *akshara* is either a vowel or a consonant-vowel combination. Salomon (2000) observes that alphasyllabic languages are neither strictly alphabetic because each character does not represent a single sound unit or phoneme, and neither are they strictly syllabic because each character does not stand for a single and indivisible sound unit. According to Nag and Narayanan (2019) Tamil orthography is partially transparent, and certain current and historical influences have brought some opacity into it. It is more opaque than other languages in the same language group. Additionally, Tamil is a highly diglossic language—Schiffman (1998) established that Tamil has a distinct "High" variety used in literary writing and formal communication and a "Low" variety that is standardized and used for informal communication and everyday speech. Although there is a scarcity of research on how Tamil diglossia affects Tamil literacy (Nag & Narayanan, 2019), it is possible this distinction between written Tamil and spoken Tamil affect how children acquire Tamil collocations.

The English language belongs to the Indo-European family and is written using the Modern Latin Alphabet. English orthography is considered to be almost totally opaque because both decoding and encoding require additional mappings apart from direct mapping of pronunciation to spelling (Kahn-Horwitz, Schwartz & Share, 2011). As Kahn-Horwitz et al. (2011) note, acquiring symbol-sound correspondence of the English alphabet is only the first step in English literacy acquisition—fluent reading requires recognition of orthographic patterns that make words graphemically complex as well as acquisition of words that have varying degrees of phonemic regularity e.g *one* where only /n/ is phonemically regular. Frost (2005) states that the multiple vowel system in English (15 in total) with fewer graphemes to represent them is what makes English orthography opaque. As Kahn-Horwitz et al. (2011) explain in their study regarding the acquisition of English orthography, this acquisition poses

an extraordinary challenge for a young English L1 learner and even more so for a young English L2 learner.

From the above descriptions of both languages, we can see that Tamil and English are orthographically far apart because they use different scripts and vary widely in their levels of opacity. Also, the diglossic nature of Tamil plays a role in its acquisition and should be taken into account when considering how knowledge of Tamil collocations may influence processing of English collocations for Tamil-English bilingual children.

## 2.15 Conclusion: Research Gaps

From the discussion on the theoretical framework for the present studies as well as an overview of the empirical research that has been done in the area of the influence of the L1 on L2 collocation processing, it is clear that although significant progress has been made over the past few decades, there is still plenty of work that needs to be done to understand how the L1 influences L2 collocational processing and formulaic language in general. To begin with, most of the work that has investigated L1 influence in collocation processing has been done with adult learners (e.g. Wolter & Gyllstad, 2011, 2013; Wolter & Yamashita, 2018; Yamashita & Jiang, 2010) and there is a lack of evidence as to how this process takes place in bilingual children. From the discussion on childhood bilingualism, it is apparent that bilingual children differ from bilingual adults in several aspects of how they acquire and therefore process their languages and it is possible that this will show differences in how they process collocations and the influence of their L1 on this processing. Additionally, these studies have also taken into account vocabulary levels and proficiency levels of the participants and it appears that with advanced proficiency in the L2, L1 influence of formulaic language

diminishes (Carrol et. al, 2016). The studies in this thesis will also take into account the vocabulary levels and proficiency levels of the participants to see when this holds true for bilingual children as well. The studies in this thesis will employ self-paced reading and eye-tracking, two methodologies that have been used extensively in investigating the processing of formulaic language. Secondly, the primary focus of most of models of bilingual lexical representation and access is on how single words or lexical items are processed. Since recent research has demonstrated that multiword items, such as collocations, are often processed as one lexical unit, it would be interesting to explore how bilingual children process them and how this can expand our understanding of the bilingual mental lexicon. Finally, these studies aim to add to our understanding of the extent of L1 activation during online L2 processing in languages that are orthographically different. The two studies presented in this thesis aim to contribute to filling these gaps and develop our knowledge of how the L1 influences L2 collocational processing in bilingual children.

# Chapter 3: Study 1

Based on the review of literature which has covered research done in the field of processing of collocations, Study 1 was designed in order to investigate the influence of the L1 on L2 collocational processing in children. This study was carried out in a primary school in Chennai, India with 58 Tamil-English bilingual children. This chapter reports the context, design, methodology and results of Study 1 and ends with a brief discussion of the results which will be explored in greater detail in Chapter 5.

## 3.1 Linguistic Context in Chennai, India

In the highly multilingual setting of India, English has the constitutional status as one of the two official languages and serves as a lingua franca (Ayyar, 1993), particularly in large cities such as Chennai. According to Prabhu (1984), English is the dominant language of higher-level administration, large-scale business and commerce, the judiciary and forms also a considerable portion of creative and artistic output. It is viewed as the language of opportunity and it is seen as a desirable and important tool for upward social and economic mobility—consequently, there is a widespread desire to be educated in English from school-level onwards which has led to English being the dominant medium of instruction especially in wealthier urban contexts.

In Chennai, English-medium schools range from schools that follow an international curriculum that cater to high-income families, elite private schools that follow the state curriculum (Tamil Nadu State Board) or a national curriculum (Central Board of Secondary Education) to low-income state and private schools that cater to children from less privileged families. Schools in this last category are only English-medium in a formal sense—although

the whole curriculum is meant to be taught in English, classroom transactions are mostly in the regional language since teachers themselves struggle with communicating effectively in English. In most of these schools, class sizes range from 40 to 60 and virtually all teachers of English have learnt English in the same education system (Prabhu, 1984). Typically, children who attend these low-income private and public schools have very minimal or no contact with English outside the school setting. As noted by Ponnuchammy (2012), the state government's efforts to promote English at these schools have been met with criticism due to poor school facilities, lack of adequate and effective training for teachers, and teachers' inefficiency at communicating in English. Gargesh (2006) observes that after approximately 12 years of education in state or low-income private schools (English-medium), most of these students in Tamil Nadu are unable to cope with the academic demands of tertiary education in colleges and universities which all have English as the medium instruction. The school that participated in this study belongs to this category and caters to children from low-income families in Chennai. Since these children are not very proficient in English and their main source of English exposure is from an L2 speaker who is not very proficient (their teacher), based on the literature review in Chapter 2 it is likely that these children's processing of English collocations will be influenced by their Tamil.

## 3.2 Research Questions and Hypotheses

1. Is the L1 activated when Tamil-English bilingual children process collocations in the L2?

2. Does congruency have a priming effect during the reading of collocations for Tamil-English bilingual children?

3. What is the relationship between the proficiency levels/vocabulary size of the children and the time they take process congruent and incongruent collocations?

Using these research questions, the following hypotheses have been formulated:

Congruent collocations will be processed more quickly than the incongruent ones which is an indication that Tamil (L1) will be activated during the processing stage of English (L2) collocations for Tamil-English bilingual children.

Congruency will have a priming effect (when the first word of the collocation prepares the reader for the second word) on the second word of collocations during the processing stage of collocations for Tamil-English bilingual children.

Children with a bigger vocabulary size in Tamil and a smaller one in English will show a larger congruency effect i.e. they will show a larger difference in the reading times between congruent and incongruent collocations. Children with higher levels of English proficiency will show a smaller congruency effect or none at all, i.e. they are likely to have very small differences between the reading times of congruent and incongruent collocations, or no difference at all.

**3.3 Design**

This study followed a 2x2 within-participants design with congruency and presentation mode as the independent, between-participants variables. In this study, congruent collocations refer to collocations that can be translated word-for-word from English to Tamil while still retaining the exact same meaning i.e. English collocations that have an exact equivalent in Tamil; while incongruent collocations are English collocations that do not have equivalents in Tamil. There were two independent variables: congruency and presentation mode. Each of the independent variables had two levels: congruency had congruent and incongruent version and presentation mode had the single mode and the chunk mode. In the single mode, the sentences were presented to the participants word-by-word, so that they were required to

press a button to go on to the next word in the sentence, thus the collocations were read word-by-word as well. In the chunk mode, the participants were presented with each sentence in two-word units (chunks) so the participants read each collocation as a single unit. This was done so that the reading times for collocations as a whole could be recorded, as well as each individual word to investigate whether there was a priming effect (when the first word of the collocation prepares the reader for the second word). Thus, each participant was exposed to four experimental reading conditions: single congruent, single incongruent, chunk congruent and chunk incongruent. The dependent variables were the reading times for the four conditions, with scores of the measures of vocabulary and proficiency in both languages as covariates (these measures will be described below) X-lex tests (English and Tamil), the scores of the C-test in English as the covariates. The reading times for the collocations were measured in milliseconds

## 3.4 Participants

The participants for this study were fifty-eight 5[th] grade students (34 girls, 24 boys) who were all 9-10 years old from an English-medium primary school in Chennai, India. This age group was chosen because a certain degree of reading ability and English proficiency was required for participation in the self-paced reading experiment so that they could complete the self-paced reading task. All of them had Tamil as their L1 and English as their L2 and did not speak any other languages. The medium of instruction was English, but they also learned Tamil as a separate subject. All the students had 6-7 years of English-medium education, from the time they started school. The school chosen for this study catered to children from low-income families because although it is a private school, it is a low-tier private school with relatively cheaper fees and would be an ideal choice for parents who want their children to be educated in an English-medium school, but cannot afford other pricier, well-resourced

private English-medium schools. So although the medium of instruction in their school was English, they had very little access to English outside the school environment and their preferred language of communication was Tamil. However, since most of their subject books were in English, they had more exposure to written English than written Tamil.

**3.5 Ethics**

Prior to the data collection phase, this project was subject to ethical review by the School of Literature and Language's Ethics and Research Committee and was granted ethical approval according to the University's ethical guidelines. Permission was obtained from the headteacher to conduct the study and then information sheets and consent sheets were sent to the parents via the class teachers prior to the commencement of the data collection. As the children who participated in this study were below the age of consent (all the children were either 9 or 10 years old), information sheets and consent forms were given to the parents in Tamil, as this was deemed to be more appropriate than English (see Appendix 1). Opt-in consent was used and parents were required to sign the form and return it to the class teacher if they agreed to allow their child to participate in the study. The researcher gave an short talk to each class to explain the components of the study and what each child would be required to do. As the study took place in India, there was no requirement for a DBS check or any such equivalent. The data obtained from the study was anonymised and securely stored as per University guidelines.

**3.6 Preparation of Materials**

As the primary aim of the current study was to study the influence of Tamil on the processing of English collocations, suitable stimuli with both congruent and incongruent collocations

had to be prepared in order to measure the reading times. This section will describe how the collocations were selected and how the stimuli were prepared using these stimuli.

### 3.6.1 Selection of Collocations

Since the definition of collocations is inclusive of various kinds of word combinations, a critical part of every collocation study is determining and defining the kind of collocations to include in the study. Most researchers use a combination of collocation lists, word frequency lists, collocation dictionaries, and, when appropriate, learner and native speaker corpora in the selection process. One of the most common measures of collocational strength is Mutual Information (MI)—it compares the actual co-occurrence of two words in a corpus with their expected co-occurrence if the words in the corpus appeared in a random order (Hunston, 2002). According to Hunston (2002) the MI score of a collocation is computed by dividing the value of the observed co-occurrence of the words divided by their expected co-occurrence and then converting it to a base-2 logarithm. It is a measure of how strongly two words seem to associate in a corpus, based on the independent relative frequencies of the two words (see Section 2.12.1 for more details). To illustrate this principle, Hunston (2002) goes on to give an example from the Bank English corpus of the words *baleful* and *unwavering*, which by themselves are not particularly frequent words, while *gaze* is more frequent. *Baleful* and *gaze* co-occur only 6 times in the corpus but the MI score is high; since *baleful* is a low-frequency word, even the low co-occurrence is significant.

To determine how to select the collocations and create the stimuli needed for the experiment, the researcher looked at similar studies to gain an understanding of how collocations are generally selected and used as stimuli. Since they were looking at how quickly native and non-native speakers react to frequent and infrequent collocation, Siyanova and Schmitt (2008) extracted adjective-noun collocations from essays by both native and non-native

speakers and used the British National Corpus (BNC) to determine the frequency and Mutual Information of each collocation. Underwood et al. (2004) included a range of formulaic sequences in their study, including lexical phrases and transparent metaphors and idioms. For the lexical phrases, the researchers consulted a list from a previous study done by one of the researchers and for the idioms, they consulted the *Oxford Learner's Dictionary of English Idioms* (1994). They used two corpora for the frequency analysis: the BNC and the CANCODE (Cambridge and Nottingham Corpus of Discourse in English). Since this study employed eye-tracking they had to use additional criteria with regard to word length, the number of function words used in the sentences, and the predictability of the sentences.

Since the current study focused on testing the difference in reaction times between congruent and incongruent collocations, the collocations were extracted from the textbooks of the children who participated in the study. The collocations consisted of both adjective+noun and verb+noun collocations. This was done for two reasons: (i) these Tamil-English bilingual children were young learners and did not have much exposure to English outside the classroom and (ii) it was done to ensure that the learners had already had some form of exposure to the selected collocations. This decision is supported by the findings of Northbrook and Conklin (2019) who measured the sensitivity of beginner Japanese learners of English to lexical bundles and found that these learners had processing advantages for the lexical bundles extracted from their learning materials rather than the ones from the corpus based on frequency. Thus, children in their beginning stages of learning English are more likely to be familiar with collocations from their textbooks rather than ones drawn from a corpus. After all the collocations were extracted from the textbooks, the frequency and Mutual Information of each collocation was determined using the BNC. The *Oxford Collocation Dictionary* and *Collins CoBuild English Dictionary* were consulted to verify that all the collocations are used in standard English. Using these tools, collocations were selected

based on the following criteria: (i) they had to be on the list of the 5000 most frequent words in English according to the BNC, (ii) the MI (Mutual Information) score had to be above 3, which has been determined as the minimum level for collocates to be recognisably associated with each other (Hunston, 2002) .The selected collocations were translated into Tamil by the researcher (the native language of the learners) to check whether they are congruent or incongruent. The translation was cross-checked by two other adult native speakers of Tamil. Finally, an equal number of incongruent and congruent collocations were selected to be used in the experiment. The list of chosen collocations are included in Appendix 5.

### 3.6.2 Preparation of Stimuli

Using the extracted collocations, the researcher created the stimuli to be used in the experiment. The stimuli consisted of eleven mini stories: each story was three to four sentences long and contained a mixture of four to five congruent and incongruent collocations with up to two collocations per sentence. The words in each collocation had a minimum of three words and a maximum of eight words. The collocations were controlled to match the criteria outlines in the previous section. The stories were given to 5 adult native English speakers and they were asked to rate them on their readability and difficulty, and the necessary modifications were made before finalising the mini stories. As per the design, the first five stories were presented to the learners in single mode and the latter six were presented in the chunk mode and this was the same for all participants (see Appendix 6). The stories presented in both modes were comparable in terms of length and number of collocations.

**3.7 Research Instruments**

This section will describe the research instruments chosen for this study, the vocabulary and proficiency tests and the self-paced reading task. It will also discuss the validity and reliability of these tools, as well as their suitability for the purpose of this study.

*3.7.1 X-lex test*

Developed by Meara and Milton for L2 learners (2003) the X-lex 5000 test measures receptive vocabulary size and it consists of both real words and pseudo-words. It focuses on vocabulary breadth and is based on the Yes/No test format wherein the participants are asked to tick the words they know. Although it was originally designed for EFL students at the college level, it has successfully been used with children (Daller and Ongun, 2017; Milton, 2006). It contains the 5000 most frequent English words in the BNC (British National Corpus) as its corpus and these words are drawn from British newspapers, periodicals, academic books, popular fiction etc so are only representative of written language. Each test consists of 100 words drawn from the first five frequency bands (K1, K2, K3, K4, and K5), with 20 words from each of these bands. An example of a word from the high frequency band is "from" and "defence" is a word from the low frequency band. The test also included 20 pseudo-words which were designed to be phonologically plausible but actually do not exist e.g. "fremby". For every correct word, the participants receive 50 points and lose 250 points for every pseudoword that they incorrectly mark as known. Thus, the pseudo-words act as a kind of correction formula to control for guessing. For the X-lex English test, the existing test was used (Meara & Milton, 2003) (see Appendix 2). For the X-lex Tamil test, the researcher

developed a comparable test (see Appendix 3) using the Tamil corpora available on the Sketch Engine website based on how a similar test was developed for Turkish (Daller & Ongun, 2017). Similar to the English test, 20 words from each of the first five frequency bands were chosen for the test and 20 pseudo-words were also included. Similar to the English pseudo-words, these Tamil pseudo-words were created to appear phonologically plausible but do not actually exist. Participants scored 50 points for every real word they marked and lost 250 points for each pseudoword. Thus, the maximum score would be 5000 (all real words marked as known but no pseudo-words) and -5000 would be the minimum score (all pseudo-words marked as known but no real words).

### *Validity and Reliability of the X-lex test*

Research on the Yes/No format of receptive vocabulary tests, such as the X-lex test, suggests that the general proficiency levels of the test-takers have an influence on the test scores for items belonging to different frequency bands and on their tendency to mark pseudo-words (Milton, 2010; O' Dell, Read & McCarthy, 2000). It has been reported as reliable and valid for screening and placement purposes and also as a measure of average vocabulary size (Milton, 2010). Due to the fact that it is easy to administer and score and also allows for a sampling of large number of items it has been used widely in vocabulary testing in L2 research (Harsch & Hartig, 2015). In a study that measured the relationship between X-lex scores and reading and listening skills, Harsch and Hartig (2015) found a significant positive correlation between the X-lex scores and the reading proficiency levels of the learners although this correlation is slightly lower than the one found for a C-test (see Section 3.7.2). Nevertheless, the findings of this study showed a significant correlation for the X-lex test as well. Based on the literature that confirms the reliability and the validity of the X-lex test, as well as its ease of use, it was chosen to test the learners' vocabulary size in this study (see Appendices 2 and 3). This Tamil X-lex test was piloted with native speakers of Tamil (adults)

to ensure that all the words were well-known and that the pseudo-words were sufficiently plausible. The validity was not determined the same way as for the X-lex English test.

### 3.7.2 C-test

The C-test is an integrative written test of general language proficiency for L2 learners that was developed as a response to the shortcomings of the cloze test by Raatz and Klein-Braley in 1981 (Klein-Braley, 1984). Klein-Braley and Raatz (1981) explained it was designed as an alternative to the cloze test from a theoretical psychometric perspective. According to them, it is different from the cloze test in two principal ways: (i) deletions are made at the word-level instead of at text-level and (ii) instead of one long text, it consists of 4-6 short texts. It is based on the assumption that learners with higher proficiency levels will be able to draw on their automated language skills to fill in more of the blanks, also known as the principle of reduced redundancies (Klein-Braley,1997).

A typical C-test consists of four or five authentic, short texts which are coherent in themselves (Harsch & Hartig, 2016). Each text covers a different topic in order to control for knowledge bias. In each of these texts, apart from the first line and the last line, half of every second word is deleted. The number of missing letters is not indicated. If the word has an odd number of letters, there is a choice about deleting either the smaller or the larger "half" is deleted, as long as it is followed consistently for the whole test. One-letter words, proper nouns, and numbers are left untouched. The task of the test-taker is to fill in the remainder of the half-deleted words. In most cases, there is only one right answer for each blank. Each right answer gets one point and there is no negative marking. For this study, a C-test was developed based on topics from the learners' textbooks. There were 50 blanks and so the final scores were marked out of 50. The maximum score would be 50 (all blanks completed correctly) and the minimum score would be zero (no blanks filled in correctly).

Eckes and Grotjahn (2006) and Grotjahn (2002) report the C-test to be easy to administer, objective and a reliable measure of language proficiency. Overall, the C-test has received accolades for its ease of use, simple scoring method and efficiency. However, Babaii and Ansary (2001) identify poor item discrimination and unclear construct validity as potential problems with the C-test which can lead to ceiling effects. To control for ceiling effects, Sigott (1995) recommend deleting two-thirds of the word instead of one half or leaving only the first letter of the deleted words, or alternatively deleting the first half of the word instead of the second half (left-hand deletion instead of right-hand deletion). However, the original form of the C-test is the form that has been widely used and administered as a test of proficiency and placement and this study follows the original format.

*Validity and Reliability of the C-test*

The C-test is recognised in language testing because it elicits a range of linguistic knowledge and language skills. It is acknowledged as a reliable test of general language proficiency since in a relatively short time, it allows for the sampling of a comparatively large number of items and has been used successfully as a screening and placement mechanism (Eckes & Grotjahn, 2006; Hastings, 2002). Since its introduction by Raatz and Klein-Braley in 1981, various studies have been carried out to validate the C-test as a test of general language proficiency and have reported positive results with the C-test being a strong predictor of general language proficiency (Babaii & Ansary, 2001; Dornyei & Katona, 1992). It has also found to be partly influenced by individual test-taker characteristics such as closeness of L1 and L2 and characteristics of the L2 background of the learner such as input and exposure. In a study designed to test the predictive power of C-test scores on reading ability, Harsch and Hartig (2015) found that the C-test had high correlations with both reading and listening skills. The researchers concluded that this could be because the C-test draws on a number of skills that learners need to access for the reading process. Additionally, a recent study by

Müller and Daller (2019) found that the five C-test sub-scores (for each of the texts in the C-test) yielded a value for Cronbach's alpha of 0.885 with regard to measuring general proficiency—although the C-test does not differentiate between different elements of language proficiency, it correlates to them all.  Based on the findings in the literature that confirm its validity and reliability and the simplicity of the test format , the C-test  was chosen as a test of language proficiency for this study (see Appendix 4 for the C-tests used in this study).

### 3.7.3 Self-Paced Reading Experiment

Self-paced reading, the methodology chosen for this study, is a well-established technique used to study linguistic processing and measure reading times of individual words and word units.  It was first developed by psycholinguists in the 1970s (Mitchell & Green, 1978) to allow researchers in the field of cognitive psychology to measure language comprehension and processing in real time with tasks that are as close to the natural reading process as possible (Jegerski, 2014). It is based on the assumption that the participants' reading time of each word or word unit indicates their knowledge of linguistic phenomena and their degree of reading ability (Marsden, Thompson & Plonsky, 2018). The primary advantage of self-paced reading is its relative ease of administration, implementation and data analyses and for these reasons, its popularity has persisted over time. Studies that use self-paced reading are based on the premise of a proposition by Just and Carpenter (1980) which posits that time taken to process a word is reflected in the amount of time taken to read a word. Although subsequent empirical work has revealed that the relationship between reading times is more complex than this, it is still a basic assumption that reading times provide a good measure of how processing takes place i.e. longer reading times denote a level of reading difficulty and shorter reading times indicate that faster processing has been facilitated. Juffs and Harrington

(1995) were the first to use self-paced reading to investigate the differences in grammatical processing between native and non-native speakers. Over the past few decades, self-paced reading has been used to investigate different dimensions of language processing such as processing of structurally complex sentences (Grodner & Gibson, 2005), processing of temporally ambiguous sentence structures (Garnsey, Pearlmutter, Myers, & Lotocky, 1997), the way pronouns and reflexives are interpreted ( Fedele & Kaiser, 2012), the effects of syntactic priming (Traxler & Tooley, 2008) etc.

Marsden, Thompson and Plonsky (2018) define self-paced reading as an online computer-assisted research tool in which participants read sentences which are broken into words or chunks, at a pace they determine by pressing a key to go on to the next word or chunk. The computer records each press of the key and the time in between the pressing of the key is measured as the reading time for each unit. Jegerski (2014) provides an overview of the different possible formats for self-paced reading: the main ones are the cumulative and noncumulative formats. In the cumulative format, as each stimulus unit is revealed to the participant one by one, but the previous stimuli remain on the screen as each new stimulus is revealed until finally the entire sentence appears on the screen. In the noncumulative format, only one stimulus unit is visible at a given time i.e. when each new stimulus segment appears, the previous one disappears from the screen. The problem with the cumulative display format is that participants read several stimuli units at a time then finally read them all at once, which can make it difficult to measure reading times for each unit accurately. In light of this, most self-paced reading studies use the noncumulative display format (also known as the moving windows format) and this is the format used in the present study. Jergerski (2014) goes on to describe the procedure and elements (cue and stimulus unit) of a typical self-paced reading experiment. The first element is the cue which is usually an asterisk or a similar symbol which appears at the location on the screen where the first stimulus unit will appear

subsequently. This cue encourages the participant to look directly at the location of the first word of the stimulus when it appears instead of spending time initially gazing at another location, which adds time to the reading time of the first stimulus, thereby making it an imprecise measure. After the cue, the first stimulus unit appears on the screen. Since self-paced reading is mainly used to measure sentence-level processing, each stimuli item is usually one sentence long and is divided into stimulus units which are divided into word-by-word or phrase-by-phrase units, according to the regions of interest of the specific experiment. Each of these stimulus units corresponds to a separate data point in the form of a reading time recorded in milliseconds (ms). The word-by-word format yields more precise data, but the phrase-by-phrase format is more like natural reading and may eliminate some of the more unnatural effects that can be present in self-paced reading. However, it should be noted that for both the word-by-word format and the phrase-by-phrase format the reading time also encompasses the time needed to plan and execute the action of pressing the space bar which is shortcoming of self-paced reading. Both formats have been widely used in self-paced reading research.

Self-paced reading was also the first online measure to be used in language processing experiments with L2 learners (Jegerski, 2014). It has been used in a range of linguistic processing studies—from studying the effects of language switching on reading comprehension (Bultena, Dijikstra & van Hell, 2015) to examining the processing of past tense morphology (Pilatsikas & Marinis, 2013). There are also a number of studies that have used self-paced reading tasks to study the processing of formulaic sequences: Kim and Kim (2012) used it to study the effects of frequency on multiword processing in L2 learners and native speakers, Schmitt and Underwood (2004) used it to analyse the processing of component words in formulaic sequences, and Siyanova and Schmitt (2008) used it to study the processing of L2 collocations. It has also been used in lexical studies to investigate the

extent to which L1 influences L2 processing, learning or representations based on the assumption shared with this present study i.e. the speed of processing will be affected for lexical items that are similar to L1 lexical items when compared to lexical items that do not share similarities with L1 lexical items (Marsden et. al, 2018).

In the present study, PsychoPy (Peirce, 2007) v1.8, an open-source application was used to design and run the self-paced reading task. This application was designed to handle the display of stimulus and to measure the timing of this display. Built using the programming language Python, it allows the user to either use the existing stimuli or to import their own. The first five mini stories were entered word-by-word so that they would be presented to the participants in the single mode in the noncumulative format. The remaining six stories were entered in two-word chunks so that they would be presented to the participants in the chunk mode. The stories were presented in a different random order for each participant. PyschoPy automatically recorded the reading times for each word and chunk in a separate file for each participant. Similar to the studies cited above, the participants were presented with the stimuli on a computer monitor and were required to press the space bar to proceed to the next screen.

Thus, self-paced reading is a well-established research tool that has recently been used in studies that have measure reading times for collocations and formulaic sequences as mentioned earlier in this section (Schmitt and Underwood, 2004). PsychoPy v.1.8.2 was the open-source application used to design and run this self-paced reading task (Peirce, 2007). For any experiment that measures time, precision of measurement is of crucial importance and PsychoPy syncs with the system clock which is accurate to milliseconds. For display of the stimuli, PsychoPy uses a robust double buffering mechanism that allows for smooth transition from one frame to the next (Peirce, 2009). Based on these technical details, self-paced reading was determined to be a suitable tool for this study and PsychoPy was selected as an appropriate application to use.

**3.8 Procedure**

At the beginning of the experimental phase, the researcher administered the X-lex tests in English and Tamil, as well as the C-test, to the participants, always in the same order. The participants took the tests in their classrooms—the X-lex tests were given in the morning and the C-test was given in the afternoon. For the X-lex tests, the researcher instructed the students to mark all the words they understood and knew how to use. Following test instructions, the students were not told that the test contained pseudo-words. For the C-test, the researcher explained the format of the test to the students using the example on the test paper.

The self-paced reading experiment took place during individual sessions in a room at the students' school using a laptop. Each child came in individually to take part in the experiment. For each participant, the researcher demonstrated how the task was expected to be completed using an example mini story. Before starting the task, each participant was given the opportunity to practise using the practice trial mini story. After completing the trial, the participants went on to read the 11 mini stories in a different random order for each person. Each session lasted approximately 30-40 minutes based depending on how long each child took to read the mini stories.

**3.9 Data Processing**

PyschoPy recorded the reading time for each stimulus unit on a by-participant basis in the form of an individual Excel sheet for each participant. Although PsychoPy recorded the reading times for each word and unit during the self-paced reading, only the reading times for

the collocations were required for data analysis. PsychoPy generated an output file which the researcher then converted to an Excel file. This Excel file had the reading times for each word or word unit in a by participant format. The researcher extracted the reading times for the whole collocations in the chunk mode and the reading times for each word of the collocations in the single mode. The means of the reading times were computed by participant for the F1 analysis and by item for the F2 analysis.

First, the means for each condition were computed (see Table 3.1). As is usually the case with reading time data, the data were positively skewed and therefore not normally distributed, i.e. for all the variables $p < .05$ in the Shapiro-Wilk's test. Based on accepted practice in the field, any reading times that were 2.5 deviations away from the mean for each of the four conditions were excluded to eliminate outliers (Conklin & Pellicer-Sanchez, 2016). Three sets of data were excluded from the single condition and three sets were excluded from the chunk condition i.e. data from three participants because they were outliers. This cleaned set of data (see Appendix 7) was determined to be normally distributed, for all conditions ($p > .05$ in the Shapiro-Wilk's test) and the analyses were performed on these data. All analyses were conducted by participants ($F1$) and by items ($F2$).

Table 3.1

*Mean reading times for congruent and incongruent, and single and chunked collocations (ms). Standard deviations in parentheses.*

|  | Congruent | Incongruent |
|---|---|---|
| **Single (Word 1+ Word 2)** | 2322 *(807)* | 2631 *(1204)* |
| **Chunked** | 1535 *(443)* | 1742 *(470)* |

Table 3.2:

*Means and Standard Deviations for X-lex test scores (out of 5000) and C-test scores (out of 50)*

|  | **Mean** | **SD** |
|---|---|---|
| **X-lex English** | 2742 | 957 |
| **X-lex Tamil** | 2161 | 857 |
| **C-Test** | 32.39 | 7.54 |

### *3.9.1 Vocabulary size tests*



*Figure 3.1* Distribution of vocabulary size scores

A boxplot was plotted with the results of the vocabulary size test scores in English and Tamil (see Appendix 9 for the full list of scores). The maximum raw score for each test is 5000. From the distribution of scores in the boxplot, it is evident that the students have overall bigger vocabularies in English and that they vary quite widely in their Tamil vocabulary scores ($M = 2289$, $SD = 309$), but not so much for English ($M = 2976$, $SD = 795$). The

difference between the means for English vocabulary and Tamil vocabulary was significant: $t(57) = 3.39$, $p < 0.05$.

### 3.9.2 Effects of Congruency and Test Scores

A two-way repeated measures ANCOVA was run to determine the effects of congruency and presentation mode on mean reading times while also accounting for participants' vocabulary knowledge in English and Tamil and for their proficiency in English. The combined word length (Word 1+Word 2) was also included as a covariate to determine whether the word length affected the reading times. There was a main effect of congruency with longer reading times on the incongruent than congruent collocations (as seen in Fig. 1), $F1(1,48) = 50.93$, $p <.001$; $F2(1,42) = 16.36$, $p =.016$. There was a main effect of presentation mode with longer reading times for single mode than for chunk mode $F1(1,48) = 48.21$, $p = .034$; $F2(1,42) = 20.31$, $p =.008$. The interactions between congruency and presentation mode were not significant $F1(1,48) = 16.48$, $p = .101$; $F2(1,42) = 12.78$, $p =.119$.

In terms of the vocabulary and proficiency scores, there were no significant covariate effects (all $F$s < 1; all $p$s > .1) on reading times, that is, vocabulary knowledge in English and Tamil as well as English proficiency did not affect how long the children took to read the collocations. Only the F1 values are reported here as the vocabulary scores are by participant. The results showed that there were no statistically significant interactions between congruency and X-lex English Test score, $F1(1) = .029$, $p = .866$; X-lex Tamil Test score $F1(1) = .666$, $p = .419$; or C Test score, $F1(1) = .122$, $p = .729$. There were also no significant interactions between the presentation mode and the X-lex English Test, all $F$s < 1; all $p$s > .1. With regard to the word length, there was no significant covariate effect on the reading times ($F < 1$; $p > .1$) i.e. word length did not have an influence on the reading times.

*3.9.3 Congruency Effect on Word 2 in the Single Mode*

Although vocabulary scores and proficiency scores did not affect reading times, further analyses were carried out to determine if they predicted the size of the congruency effect. Since the children were likely to show a congruency effect on Word 2 after they had already encountered Word 1, analyses were run to investigate the congruency effect on Word 2 of each collocation (see Appendix 8 for mean reading times per word). This was done in line with previous research on formulaic language which investigated the priming effect either on the last word of formulaic sequences or on the second word of collocations e.g. Carrol and Conklin (2015), Carrol and Conklin (2014), Wolter and Gyllstad (2011, 2013). From the results of their studies, Wolter and Gyllstad found that the second word of congruent collocations were primed by the first words, unlike the incongruent collocations. Based on this, it was decided to investigate the congruency or priming effect on Word 2. To assess the congruency effect (priming effect) on Word 2 of the collocations, standardized reading times (Z-scores) were calculated for the reading times of Word 2 of the collocations in single mode. The difference between the standardized reading times for Word 2 of the incongruent collocations and the standardized reading times for Word 2 of the congruent collocations was calculated as a new variable: Word 2 Difference. This difference showed the extent to which each individual read Word 2 of the congruent collocations faster than they read Word 2 of the incongruent collocations. Preliminary analyses showed the values to be normally distributed as assessed by Shapiro-Wilk's test ($p > .05$). A Pearson's correlation test was run to look at the correlation between the Tamil and English vocabulary scores and Word 2 Difference. The correlation between Tamil vocabulary and Word 2 Difference as well as the correlation between English vocabulary and Tamil vocabulary were not statistically significant. There was a statistically highly significant, negative correlation between English vocabulary scores

and Word 2 Difference $r(58) = -.84$, $p < .001$. A linear regression established that an increase in English vocabulary scores predicted a decrease in Word 2 Difference F $(1,57) = 139.51$ and English vocabulary scores accounted for 71% of the variability in the Word 2 Difference values.



*Figure 3.2*: Reading times (in milliseconds) for Word 1 and Word 2

### 3.9.4 Correlation between Reading Times of Congruent Collocations and Test Scores

A Pearson's correlation was run to examine the relationship between the test scores and the reading times of the collocations. There was a positive correlation between the scores of the X-lex English test and the X-lex Tamil tests, $r(56) = .487$, $p < .001$. There was a positive correlation between the scores of the X-lex English tests and the C-test, $r(56) = .564$, $p < .001$. However, there was no correlation between the reading times of the collocations and the test scores and this will be explored in the discussion.

**3.10 Discussion**

*3.10.1 Hypothesis 1: Congruency*

The first hypothesis was concerned with the role that congruency plays in the processing of collocations. The most important finding in this study is that the effect of congruency was significant—both F1 and F2 analyses for the ANOVA as well as the ANCOVA had significant results, meaning that the children read congruent collocations faster than they read incongruent ones—this supports the first hypothesis. This congruency effect indicates that they drew on their L1 while reading L2 collocations, in both presentation modes, i.e. when they were presented with the collocations as single words as well as in the chunk mode.

From a theoretical standpoint, it appears clear that the non-selective lexical activation is not limited to single words but also takes place in the processing stage of L2 collocations. It can be assumed that the L1 provides quicker access to L2 collocations which have an L1 equivalent than those L2 collocations that have no L1 equivalent—this recalls the lexical activation process previously discussed. As seen in Chapter 2, there are two major conceptions of the bilingual lexicon: the Revised Hierarchical Model (Kroll and Stewart, 1994) views the L1 and L2 lexica as two separate entities but the more recent Bilingual Activation Model and Multilink model posit the theory of one integrated lexicon in which both L1 and L2 lexical items are stored. Regardless of whether the L1 and L2 exist as two separate lexica or are combined in one, it is possible to examine activation independent of storage condition as has been done in similar studies (Siyanova & Schmitt, 2008). The finding that L2 collocations are processed more quickly if they have a L1 equivalent indicates that the L1 is not suppressed when L2 collocations are activated; instead it is a process of non-selective activation wherein both L1 and L2 are activated even though the learner is

presented with input only in the L2. Exploring this non-selective dual activation further in the context of these results, it is plausible that when the child is presented with the L2 word, the L1 translation is activated along with the L1 collocates and these L1 collocates activate the L2 collocates thus allowing for quicker processing times. The results of the present study support the assumption common to the all the models of bilingual lexical representation and access discussed in Chapter 2 that this non-selective, cross-lingual activation is not limited to single lexical items and takes it one step forward by extending it to collocations.

Additionally, this effect of congruency can also be explained by Jiang's model of L2 lexical representation and development (2000) which was discussed in Chapter 2. This model posits that the acquisition of L1 words with an L2 equivalent is quicker because the learner already has access to the required semantic and syntactic information at lemma level and thus only needs to acquire the phonological and orthographic information at the lexeme level. However, when a word does not have an L2 equivalent, the learner has to acquire the L2 information at both lexeme and lemma level—this process requires more effort and time. Although Jiang's model describes the acquisition of single lexical items, it can be extended to explain collocational acquisition as has been done in previous studies (Wolter & Gyllstad, 2011). Based on the congruency effect found in this study, it can be assumed that collocational knowledge is part of the L1 sematic and syntactic information that the learner retains and uses while acquiring L2 vocabulary.

As discussed in the literature review, earlier studies by Yamashita and Jiang (2010), Wolter and Gyllstad (2011) and Siyanova and Schmitt (2008) all found that with an increase in L2 proficiency, the L1 influence on collocation acquisition decreases (see Section 2.13). Although the participants in this study have all been learning English for 5-7 years and also study in a school where the medium of instruction is English, they receive negligible

exposure to English outside school. Therefore, their number of years of English education is not necessarily indicative of their English proficiency. Their main language of communication with friends and family, as well as media consumption, is Tamil. This dominance of Tamil in the everyday lives of these children could explain the proficiency level of the children and consequently the congruency effect, despite the number of years of English education. Quality of English input is also a factor to be considered since these children are taught by L2 speakers, most of whom are not proficient in English due to the lack of training.

### 3.10.2 Hypothesis 2: Priming Effect due to Congruency

From the results of the analysis on the reading times for Word 1 and Word 2 of both congruent and incongruent collocations, there appears to be a priming effect on the second word of the congruent collocations. As discussed in the literature review, in a similar study that investigated the influence of L1 intralexical knowledge on L2 collocational knowledge Wolter and Gyllstad (2011) gave an explanation of what occurs when a bilingual is presented with an L2 word which mirrors the non-selective, cross-lingual activation under the RHM discussed in Section 2.7. According to them, assuming that upon activation of an L2 word its known collocates are activated, along with its L1 translation and its known L1 collocates, there are four possible scenarios when the learner encounters a second word in any given collocation: (i) the word has been primed as an L1 collocate but not as an L2 collocate; (ii) the word has been primed as an L2 collocate but not as an L1 collocate; (iii) the word has been primed as a collocate in both L1 and the L2; and (iv) the word has not been primed as an as an L2 collocate because the learner does not yet recognize it as an L2 collocate. Out of these four scenarios, only the last two can be accounted for by the data obtained in this study as follows: since the reading times for Word 2 were shorter than those for Word 1 for the

congruent collocations, this implies that that upon reading the first word of the congruent collocation the children recalled the Tamil translation and its collocates and this helped them to predict the next word for the English version and therefore spend less time processing Word 2. This finding is explained by the third scenario—the knowledge of the collocation in Tamil (L1) has enabled the priming effect for the second word in the English collocation (L2). Interestingly, the children spent more time on Word 2 than on Word 1 while reading the incongruent collocations. A possible explanation is that the L1 was activated for Word 1 along with its collocates and so the children expected a different word for Word 2, leading to longer reading times for Word 2. This is explained by the fourth scenario of the framework which implies that without the aid of the L1, the learner does not have sufficient knowledge of L2 collocates to be able to predict Word 2 in incongruent collocations. The results of the present study support the results of previous studies that have found a processing advantage for congruent collocations (see Section 2.14) e.g. Yamashita and Jiang (2010), Wolter and Gyllstad (2011, 2013), and Wolter and Yamashita (2018)—these studies all show that across contexts and with different groups of learners, there is a priming effect on the second word of collocations which means that congruent collocations are processed more quickly than their incongruent counterparts. This will be further explored in Chapter 5.

**Presentation Mode**

Similar to the findings for congruency, the results showed that the learners had shorter reading times for the collocations in the chunk mode than the collocations in the single word mode. However, this finding could also be partly attributed to the methodology used for the self-paced reading experiment—for the single word mode, the participants read the first word of the collocation on one screen and had to press a button to go to the next screen to read the

second word of the collocation whereas for the chunk mode, they were able to read the whole collocation on one screen. Once this has been taken into account however, it should be noted that the difference between the mean reaction time for the chunk mode and the single mode was over 750 ms for both congruent and incongruent collocations— this suggests that to some extent at least, the children process collocations more quickly when presented with them as units than as single words. However, these findings are only suggestive since self-paced reading cannot give a clearer picture: the most reliable and accurate method of measuring this assumption would be through eye-tracking. With regard to previous research in the field, the chunk mode is similar to the lexical decision tasks used by Wolter and Gyllstad (2011) and Carrol and Conklin (2014) since they also used self-paced reading and presented one word of the collocation at a time, although in this study the participants were not required to decide whether the second word of the collocation was acceptable or not. However, these studies do not provide any comparison between collocations presented as a whole and collocations presented in chunks.

### 3.10.3 Hypothesis 3: Vocabulary and Proficiency Tests

There was a positive correlation between the X-lex English scores and the X-lex Tamil scores as well as between the X-lex English scores and the C-test scores, i.e. the children who scored well on the X-lex English test also scored well on the X-lex Tamil, with the same being true for those who received lower scores. The test scores from both vocabulary tests, the X-lex English and the X-lex Tamil, did not correlate to the reading times under any of the four conditions and did not correlate to the congruency effect either.

### *3.10.4 X-lex Tests*

Firstly, as mentioned in section 3.4.1, the X-lex tests draw from a list of the 5000 most frequent words in English. However, as shown in previous studies, the words which constitute the most frequent words in a child's lexicon may not correspond to the overall most frequent words in English. Words found in children's storybooks, textbooks and the average classroom need not be representative of the words that occur frequently in general usage. For the Tamil X-lex test, the words were taken from a corpora of the 5000 most frequently occurring words in written Tamil from the website Sketch Engine since there is no comparable corpora for spoken Tamil. While the frequency discrepancy discussed above could have also explained the lack of correlation between the test scores and the reading times, for the Tamil test there is an additional factor that should be taken into consideration. As previously discussed, Tamil is diglossic in nature—there is a wide gap between the spoken and written forms of the language (see Section 2.14). The children in this study are more familiar with the spoken form of the language (low Tamil) because they encounter it and use it far more in daily life than they encounter the written form in school (high Tamil). This could also explain why the mean score of the Tamil X-lex is much lower than the mean score of the English X-lex; they scored higher on the English X-lex because English is not diglossic.

Secondly, the Yes/No format and the control system of the test might not have been suitable for young learners, especially for learners with lower proficiency levels who are more prone to guesswork which could affect the test results. Although the correction formula has been widely proven to control for guessing (Eyckmans, 2004; Fitzpatrick & Clenton, 2010), it is possible that when taken by younger and less proficient learners this correction formula could work negatively. Learners in this category could have the tendency to guess as many words

as possible which could raise the number of false-alarms (number of marked pseudo-words), thus adding to the number of negative points the participant receives. This would mean that the overall score is less likely to accurately reflect the vocabulary size of the test-taker.

Thirdly, learner attitude during test taking could have affected the test results. The attitude of the participants is an important factor in any test but even more so with children. It is possible that some of the children did not take the test seriously as they saw it as a fun activity rather than something they should give their full attention to, which is something the researcher observed during the data collection. It is also possible that the children were not used to the test format and this affected their ability to answer according to their knowledge.

Additionally, the method of administration of the test could have affected the test score—in other studies with children the test was orally administered (Daller & Ongun, 2018) in order to prevent confusion with unfamiliar spellings which is particularly pertinent for diglossic languages. In the present study, the children were given copies of the test and asked to circle the words—this could have affected the overall scores.

Finally, it is also possible that this test did not accurately capture the children's vocabulary knowledge since it relies heavily on children's self-reporting of their vocabulary knowledge rather than directly testing it. The best way to test vocabulary knowledge as accurately as possible is with multiple tests—including tests on vocabulary breadth, vocabulary depth, productive knowledge, receptive knowledge—which directly test understanding of words (Bogaards & Laufer, 2004).

**3.10.5 C-test**

Similar to the results of the vocabulary tests, the scores of the C-test did not correlate with any of the reading times either. There was a correlation between the X-lex English score and the C-test scores: the children who scored well on the X-lex English test also did well on the

C-test with the same being true for those who received average and low scores, showing a correspondence between vocabulary size and language proficiency (which is expected when taking more than one language measure from the same individual), even though this correspondence was independent of reading times.

Although the C-test has widely been recognized as a reliable and measure of language proficiency for L2 learning (Hastings, 2002), Eckes and Grotjahn (2006) point out that it is possible for test takers with high levels of reading comprehension to perform poorly on it due to lack of productive skills such as writing and speaking. As discussed in the vocabulary test section, it is also possible that learner attitude affected the performance of the children in different ways. It could also be the case that some of the children did not fully understand the concept of the C-test—this can be seen in a few of the test papers in which they wrote complete words in the blanks instead of completing the given words.

**3.11 Additional Factors**

Of course, there are other factors that influence reading times that were not taken into account for this study since the focus was on the effects of congruency. The most significant of these factors is perhaps frequency of occurrence of the collocations and the individual words. Research in vocabulary acquisition has established that higher frequency words are generally acquired before lower frequency words and it is possible that the same pattern applies to acquisition of collocations. However, Schmitt and Underwood (2004) found that although frequency of occurrence of formulaic language had a significant effect on the processing abilities of native speakers, this was not true for non-native speakers because their experience of language might differ in terms of how often they are exposed to formulaic language. Another factor in this study was that the stimuli collocations were not randomly chosen from a corpus—they were extracted from the textbooks of the learners so it can be safely assumed

that the learners had encountered them, regardless of their frequency. Nevertheless, exposure via textbook may not account for all forms of collocational exposure and an analysis including the frequency of each of the collocations at a later stage could provide further insight into the interaction between congruency and frequency in this context.

Another factor that could influence the reading times of collocations is the input. As already stated the collocations in this study were extracted from the textbooks of the learners so it was assumed that all the participants had been exposed to all the collocations at least once. We know that words need to be encountered many times in order that their meaning is learned and the context of those exposures is important (Nation, 2017). However, to control for input and exposure, while also testing for congruency, it would be interesting to conduct a study in which learners are given a controlled number of exposures to novel congruent and novel incongruent collocations and then measure whether congruency and number of exposures have a significant interaction. A study by Durrant and Schmitt (2010) revealed that the number of exposures does have an influence on adult learners' retention of collocations—testing this in children with the additional factor of congruency could provide some valuable insights into child L2 collocation acquisition.

**3.12 Limitations of the study**

The most evident limitation of this study the lack of control for frequency and word length. The measures taken to check whether this lack of control affected the significance of the congruency effect have been discussed earlier in this chapter. In this study, word length did not influence the reading times (as explained in Section 3.9.2) possibly because of the small range of word length in the stimuli. As mentioned, frequency was not taken into consideration since the collocations were extracted from textbooks and not from a corpus.

However, in future studies it is important to control for length and frequency where possible in line with other studies in this field.

Another limitation was the use of self-paced reading as the methodology which was chosen due to logistical limitations. While self-paced reading has been used in previous research on processing of lexical items (Schmitt and Underwood, 2004) and formulaic language (Wolter & Gyllstad ,2011; 2013), the most recent research has been done with eye-tracking because it provides insights that are not possible with self-paced reading. It is possible that the requirement to press the space bar to proceed to the next word or unit could mask or distort the fine time differences that are present in the processing of collocations. Additionally, self-paced reading is based on the left-to-right reading noncumulative format and does not allow participants to go back and read a portion of the text—something which occurs frequently during natural reading.

The single mode was included in this study to examine the priming effect; however it is possible that since the common consensus is that collocations are processed holistically, this may not have been the best format to present them to the learners. However, this did provide some insight into the priming effect due to congruency, but these findings are indicative and should be supplemented with eye-tracking data for a more comprehensive picture.

**3.13 Recommendations**

In order to study the correlation between vocabulary size, proficiency and reading times, it would be beneficial to either modify the tests used in this study or replace them with tests that are more appropriate for this age group and proficiency test. For the X-lex tests, it would be advisable to base locate more suitable corpora (in both Tamil and English) and base the tests on these. If this is not possible, different vocabulary size tests such as the Peabody Picture

Vocabulary tests or other alternatives can be considered because they directly test actual vocabulary knowledge. For the C-test, it is recommended that the current test is piloted with a similar sample group and if necessary, replaced or supplemented with a different proficiency test.

This study only looked at two-word collocations and did not take into account any collocation sub-categories. Other studies have found differences between the ways different collocations are acquired (Laufer and Waldman, 2011). This suggests that the different kinds of collocations could have an effect on the ways that they are processed—future studies could take this into consideration and look at how congruency affects the processing of different kinds of collocations. Future research could look at how the L1 influences acquisition and processing of different kinds of formulaic sequences such as idioms and phrasal verbs in children to build on the work already done with adults (Carrol, Conklin, & Gyllstad, 2016), (Laufer, 2000). These kinds of formulaic expressions can vary in different degrees from language to language and understanding the influence of the L1 on their acquisition and processing is important.

**3.14 Conclusion**

From the results and discussion of Study 1, it is clear that congruency does play a significant role in the processing of collocations in Tamil-English bilingual children in this context. The next chapter reports on Study 2 and will further explore how this congruency effect works and develop on the limitations of Study 1

# Chapter 4: Study 2

## 4.1 Introduction

The results of Study 1 showed that the L1 influences processing of L2 collocations for the bilingual Tamil-English children in Chennai, India who participated in the study—this was expected because the sample group was Tamil-dominant and so it was more likely for their L1 to influence their L2 processing. To further investigate the effect of congruency on the reading of English collocations, it was decided to study the effects of congruency on collocation reading on a sample of Tamil-English bilingual children in the UK for whom English was their dominant language. This group would be far more heterogenous in the exposure to Tamil and English than the sample for Study 1, and it was decided that looking at the congruency effect in relation to different levels of exposure and vocabulary scores in both languages would be beneficial for a better understanding of the congruency effect in these bilingual children Study 2 used eye-tracking rather than self-paced reading, with different stimuli than the stimuli used for Study 1 (see Section 4.6.4). The change in methodology allows the researcher to capture a more natural reading process without accounting for resources needed for key presses thus providing more detailed insight into the effects of congruency, and the change in stimuli will allow for more meaningful comparisons between the data for congruent and incongruent collocations. This chapter will begin with providing the context for this study as well as an overview of eye-tracking in general and will then report the methodology of Study 2, present the results and end with a discussion which will be expanded on in Chapter 5.

**4.2 Linguistic Context in UK Schools**

Due to the rapid rate of globalisation, the number of children with English as an Additional Language (EAL) in the UK has increased dramatically from 499,000 in 1997 to 1,306,829 in 2017, an increase of 161% (Demie, 2018) and speaking about 360 languages between them. EAL children form about 17% of the primary school population and almost 14% of the secondary school population (Department for Education, 2019). The number of and percentage of EAL pupils varies widely by area of the country, urban or semi-urban context and the type of school. The fastest growing language group in primary and secondary school is Eastern European (NALDIC, 2014) although Bangladeshi, Pakistani, Chinese, and Indian children still form the majority of the EAL group. Studies have shown that while EAL children who are the in beginning stages of developing English fluency achieve lower Key Stage 2 (KS2) test scores than their monolingual peers (Strand & Demie, 2005), EAL children who are fully fluent in English consistently score higher than their monolingual peers on all KS2 tests (Demie, Lewis, & Taplin, 2005). A report by Strand and Hessel (2018) reveals that the most important factor for EAL children's English proficiency is age: at the end of Reception, more than 55% of EAL children are still acquiring proficiency, 49% at the end of Key Stage 1, then drops to 23% at the end of Key Stage 2 and by the end of Key Stage 4, drops further to one in six children (15%). Another way of looking at it is that while only 30% of Reception EAL children are proficient in English and this increases to 85% at the end of KS4. EAL attainment is multifaceted and depends on variables such as age of arrival in the UK and home language (Hutchinson, 2018). A report by Strand, Malmberg, and Hall (2015) also found that a combination of factors including ethnicity, home language and socio-economic status (SES) interactively affect educational attainment in EAL children.

**4.3 EAL Children**

The UK Government defines an EAL pupil as anyone who has been exposed to a language other than English during childhood and continues to be exposed to this language in the home or the community. As is evident from this very broad definition, EAL children constitute a widely heterogenous group of pupils from different linguistic and ethnic backgrounds who vary significantly in their English proficiency and levels of achievement. Demie (2017) notes that the UK National Pupil Database includes in the EAL category pupils who were born in the UK and are fully fluent in English as well as pupils who have recently arrived in the UK as migrants and speak little to no English. With such a wide variance in the children who constitute the EAL group, it is misleading to treat the group homogenously and when assessing EAL children, individual differences as well as individual linguistic backgrounds should be considered.

Data from national tests of language and literacy show that EAL children are outperformed by their English-speaking monolingual peers during early years of school education in England (Bowyer-Crane, Fricke, Schaefer, Lervåg, & Hulme, 2017), which is consistent with the data from Strand and Hessel's 2018 report (see previous section). However, Babayiğit and Shapiro (2020) found in their study that EAL children's listening and reading comprehension should be considered in tandem and examined beyond the primary-school years to clarify the long-term implications of the observed EAL gap at the primary-school level. Additionally, they recommend that both vocabulary and grammar skills need to be targeted to support EAL children's listening and reading comprehension in the primary school years. These early years are crucial for the development of literacy skills and in order for appropriate language support to be extended to EAL children, a more thorough understanding of their L2 development is needed. Several studies have shown that EAL children tend to have strong reading skills but poorer reading comprehension skills (e.g.

Burgoyne, Whiteley & Hutchinson, 2011, 2013; Burgoyne, Kelly, Whitely & Spooner, 2009; Lesaux, Crosson, Kieffer, & Pierce, 2010), although there are other studies which have shown that this is not always the case (e.g. Lesaux, Rupp, & Siegel, 2007; Lesaux & Siegel, 2003). In terms of word reading skills and spelling, Smith and Murphy (2014) found that the EAL children outperformed their monolingual peers, a finding that echoes the results of a study by Burgoyne et al. (2009) which found that EAL children perform better on word reading and text accuracy tasks than their monolingual peers. Smith and Murphy (2014) observe that in the context of EAL children, research on vocabulary has mainly focused on measuring knowledge and acquisition of single words and that the field of multiword vocabulary remains relatively unexplored.  Additionally, most EAL research looks at children from different L1 backgrounds and in terms of understanding L1 influence on L2 acquisition, this is limited because the wide variance in L1s makes it very hard to determine the role of L1 in L2 acquisition. While it is good to look at children with a range of L1s since this is reflective of what we see in the classroom and is representative of this population, this approach is limited when attempting to look at the influence of specific L1s on L2 processing and having children with the same L1 is the only way to do this. As previously mentioned, EAL is a UK-specific term and although studies in other countries may be relevant, they do not really look at EAL in the same way e.g. in the United States, these students are often referred to as minority language learners.

## 4.4 EAL children in this study

Canagarajah (2008; 2013) has done extensive work investigating language attitudes and heritage language maintenance in Tamil diasporic communities in cities in the United Kingdom, Canada and the United States. Across these three communities, he found that the first-generation immigrants in these communities highly value Tamil as directly connected to their heritage and are increasingly worried and dismayed that their children seem to shift

completely to English and lose Tamil about fifteen years after migration. This rapid language attrition surpasses the typical immigrant language shift that spans three generations: the first generation migrating as monolinguals, the second generation becoming bilingual, and the third generation shifting to monolingual in the L2 (Winter & Pauwels, 2005). The older generations see this loss of Tamil among younger generations as the beginning of the demise of the Tamil community in diaspora settings, because they believe it is impossible to maintain ethnic identity without the language (Canagarajah, 2013). Interestingly, most of the teenagers and children interviewed by Canagarajah said that they did not think that their Tamil identity depended on their Tamil proficiency. In terms of language choice within the family, there were clear generational divides: 82.5% of the children said they use English to communicate with their siblings, compared to 77.6% of parents who said they use Tamil to communicate with each other. 19.9% of the children said they used a mixture of both languages, and just 2.5% reported English as their preferred language of communication with each other (Canagarajah, 2008). Among the many reasons given for this loss of Tamil as a heritage language, most parents and elders argue that the failure of the families to maintain the intergenerational transmission of Tamil is what has led to its rapid attrition in the younger generation. They view investing in Tamil cultural associations, weekend schools, media etc. as ways of engaging Tamil families and keeping the language alive in the community (Canagarajah, 2013).

The EAL children who participated in this study all attended weekend Tamil schools in the UK and they come from privileged, highly educated Tamil-speaking families which is typical of families who send their children to Tamil schools. Apart from these two factors in common, they vary widely in how long they've been in the UK and how much exposure they receive in their two languages, English and Tamil, and how often they choose to use each of the languages. Unlike the sample group in Study 1 who had similar proportions of exposure

to Tamil and English in both home and school settings, the sample for Study 2 differed quite substantially in their exposure to English and Tamil (see Section 4.8.3). Like the EAL group that they are a part of, this subset of Tamil EAL children is quite heterogenous and thus their Tamil and English vocabulary scores, as well as the scores from the language exposure questionnaires (for both Tamil and English) will be an important part of how the congruency effect is measured in this study.

This study will examine the extent to which the L1 is activated in the processing of English collocations to further our understanding of how EAL children acquire and process collocations by focusing on Tamil-speaking children in the UK. Like Study 1, the children in this study are also Tamil-English bilingual, but they vary a lot more in their proficiency in both languages as previously mentioned. The children in Study 2 also have a lot more exposure to English than the children in Study 1 because they live in an English-speaking society. As seen in the studies reviewed in Chapter 2, congruency effects tend to be smaller when the L2 proficiency is high and this is expected with the majority of the children in this study. It is also expected that exposure to Tamil and English will influence the extent of the congruency effect, with children with higher Tamil exposure showing a larger congruency effect.

**4.5 Research Questions and Hypotheses**

Based on the Literature Review in Chapter 2 and the discussion on EAL children in the previous section, the following research questions have been formulated.

1. Is there an overall congruency effect on the whole collocation i.e. do the children read the congruent collocations faster than the incongruent collocations? Is there a congruency effect on the individual words (Word 1 and Word 2) of the collocations?

2. Do vocabulary size and language exposure predict the size of the congruency effect?

Based on the above research questions, the following hypotheses have been formulated

1. It is hypothesised that the congruency effect in this study, for the whole collocations and for Word 1 and Word 2, will be smaller than the congruency effect in Study 1 due to the different contexts of both studies.

2. It is hypothesised that the children with higher Tamil vocabulary scores and higher Tamil use/exposure scores will have a larger congruency effect than those with lower Tamil vocabulary scores and lower Tamil use/exposure scores. It is also hypothesised that high English vocabulary scores will be associated with a smaller congruency effect.

**4.6 Methodology**

This section will present the methodology used in this study, including details of the participants, the experimental design, an overview of eye tracking and eye movement measures, and a description of how the stimuli were created for the eye-tracking experiment.

**4.6.1 Design**

This study had a within-participants design with one independent variable - congruency - which had two levels: congruent and incongruent: each item on the stimuli list had a congruent and incongruent version and all participants saw one version of each item, with an equal number of total items in each condition which were counterbalanced during administration. The dependent variables were four reading time measures on the collocations, as well as separate reading times for Word 1 and Word 2 of each collocation for the same four measures. The covariates were the language use/exposure scores from the questionnaire and the scores from the vocabulary tests.

### 4.6.2 Participants

The participants for this study were 80 Tamil-English bilingual children (49 girls, 31 boys), aged 8-11 years ($M = 9.6$, $SD = 1.7$), recruited from five Tamil weekend schools in southeast England, commonly called complementary schools. The children varied widely in how long they have been in the UK, ranging from 1-11 years ($M = 6$, $SD = 3$). All the schools that participated in this study followed the same curriculum and ran for one morning or afternoon a week for a duration of three hours. The primary focus of the curriculum is reading and writing, but the students also take a short oral test at the end of every term. Only children who had been attending the school regularly for at least a year at the time of testing took part in this study because they needed to be able to take the Tamil vocabulary test. All the children attended local primary state schools and so received their formal education in English.

### 4.6.3 Ethics

Prior to the data collection phase, this project was subject to ethical review by the School of Literature and Language's Ethics and Research Committee and was granted ethical approval according to the University's ethical guidelines. The researcher first obtained permission either from the headteacher of each Tamil school or the school's Board of Directors to conduct the study in the school. The researcher then met the parents of the children in each weekend school and explained the purpose of the project to them shortly before commencing the data collection. As the children who participated in this study were below the age of consent (all the children were 8– 11 years old), information sheets were given to the parents (see Appendix 10) and they were also asked to sign consent forms (see Appendix 11), prior to the commencement of the data collection (see Appendices 10 and 11). After collecting consent forms from the parents, the researcher had a session with the children to explain the elements of the study to them. All the data collection sessions were conducted at the Tamil

weekend schools during the working hours on Saturdays and Sundays. In order to fulfil requirements to work with children, the researcher obtained a DBS certificate prior to the data collection.

### 4.6.4 Preparation of stimuli

Similar to Study 1, the purpose of Study 2 was to study the processing of congruent (collocations that can be translated word-for-word from English to Tamil while still retaining the exact same meaning) and incongruent collocations (English collocations that do not have equivalents in Tamil). A new set of stimuli had to be created for two main reasons: (i) the stimuli from Study 1 had significant shortcomings (see Section 3.12), mainly because they were not controlled for length and frequency and didn't have the same sentence frame (ii) the participants for Study 2 had a very different language background from the participants of Study 1. Thus, new stimuli were created to fix the shortcomings of the first set of stimuli and also to be tailored according to the different sample group in this study.

### 4.6.4.1 Selection of collocations

Unlike Study 1 in which the collocations were taken from the participants' textbooks, the collocations for Study 2 were drawn from lists of commonly used collocations found in the *Oxford Collocations Dictionary*—since the participants of this study had far more exposure to English than the participants of Study 1 and so it was not necessary to extract the collocations from textbooks. Only two-word collocations (excluding articles and connectives) and collocations in which each word of the collocation was 3-8 letters long were extracted from the lists. These parameters were chosen to facilitate easy control of word length and collocation length. Typically, target words are at least five characters in length to minimize risk of skipping, but in this study some of the collocations had words with three or four letters

and the decision was made to retain them since there were not enough suitable collocations with both words with five characters or more. The frequencies of the individual words of each collocation were checked against the Children's Printed Word Database (CPWD). This database was first developed by researchers in the UK in 1993 as a database comprising the words that appear in the children's books during their primary school years in order to make it easier to develop stimuli for experimental work focusing on the literacy acquisition and development of young children. Although it is a small database of 12,193 types and 995,927 tokens (Masterson, Stuart, Dixon, & Lovejoy, 2010), it was chosen as the frequency determiner for this study since the children in this study all attended primary schools in the UK and this database is the closest representative database that matches the words that these children would likely encounter. This method was also the most reflective of that used in Study 1 in which collocations were extracted from the textbooks the children encountered.

### 4.6.4.2 Creation of stimuli

The researcher chose two-word collocations from this list (some containing determiners) and the length of each word was between 3 and 8 characters. The Mutual Information (MI) score for each collocation was determined using the BNC since this information is not available on CPWD. Three bands were determined based on the MI score: MI score 3-6 (low), MI score 6-9 (medium), and MI score 9-12 (high). Once the collocations were chosen and sorted according to MI score, they were divided into congruent and incongruent collocations based on whether they could be translated word-for-word into Tamil and retain their meanings. This division into congruent and incongruent collocations was done by the researcher and was cross-checked by two other Tamil-English bilingual adults. The final list of collocations was divided into congruent and incongruent and were further subdivided so that collocations in the same MI score band were grouped together. Once these divisions were finalised, the researcher created 30 sentence frames in which either a congruent or incongruent collocation

could be placed. This ensured that items across conditions were identical other than the target

collocations since the incongruent and congruent collocations for each sentence frame were

chose from the same MI band. The collocations were always placed in the middle of

sentences and so were never at the beginning or end of a sentence. These sentences were

piloted with four native English children (eight to ten years old) and they were asked to mark

which sentences they found easy and which they did not. The sentences marked as difficult

were modified to make them simpler. With the same raters, the collocations were also

checked for predictability to ensure that for each pair, one collocation was not more

predictable than the other. This was done by leaving the second word of the collocations

blank and asking the children to choose from among four options to complete them. Five

collocations were less predictable than the others and these were changed accordingly. Thus,

it was ensured that the collocations were equally predictable.

The final stimuli list consisted of thirty sentences, each with two versions: congruent or

incongruent, with a total of 60 sentences (see Appendix 15). Each participant in the eye-

tracking experiment read either of these conditions, i.e. each participant saw a total of 30

sentences, with either the congruent or incongruent version of each sentence—a total of 15

congruent and 15 incongruent collocations.

Example (collocation in bold)

The police officer asked Sam to **take a seat** immediately. (Incongruent) Tamil *oocara*--sit

The police officer asked Sam to **pay the fine** immediately. (Congruent) Tamil *kadamai--fine kettu*--pay


### 4.6.5 Research Instruments

This section will describe the research instruments used in this study: a reading efficiency test

(Test of Word Reading Efficiency; TOWRE), the vocabulary tests (X-lex test) in both

English and Tamil, the language background and use questionnaire, and the eye-tracking experiment.

### *4.6.5.1 Reading fluency test: TOWRE*

The *Test of Word Reading Efficiency-Second Edition* was designed by Torgesen, Wagner, and Rashotte (2012) to measure phonetic decoding and fluency of sight word reading for individuals aged 6 to 24. The test consists of two subtests: the Sight Word Efficiency (SWE) subtest tests the number of words an individual can read correctly in 45 seconds from a vertical list and the Phonemic Decoding Efficiency (PDE) subtest requires participants to use their decoding skills to read aloud as many nonwords as possible in 45 seconds, again from a vertical list . Although this test was normed on monolingual Americans, in this study it was used to ensure that the participants met inclusion criteria. As explained by Tarar, Meisinger, and Dickens (2015), for the regular word reading list the words are chosen on the basis of word frequency, number of syllables and complexity. The list begins with high frequency words and gradually progresses to less frequent words with more syllables. The nonwords were designed to cover a broad range of grapheme-phoneme correspondences. The nonword list begins with monosyllabic words with simple phonemes and then gradually progresses to more complex phonemes with more syllables.

The test can be used for a number of purposes including diagnosis of reading disabilities, identifying individuals who require more instruction in word reading skills and also for the purpose of the present study: to reliably assess individual word-reading skills. The TOWRE was chosen for this study because it provided a quick and efficient way to determine if each individual child possessed the word reading skills necessary to take part in the main experiment. Since this study required the children to read quite a few sentences in one sitting,

it was important that they were fairly fluent and efficient readers so they could complete the eye-tracking experiment. The score for each of the subtests is the number of words or nonwords read correctly in 45 seconds and items skipped or read erroneously are marked as incorrect. Based on the raw score, a standard score is calculated and for this study a raw score of 85 was required for the children to go on to the eye-tracking experiment. All participants achieved the standard (adjusted) score required for their age group and were allowed to participate in the reading task ($M = 66.98$; $SD = 7.83$). This indicates that this sample of children are very good readers as it would normally be expected that a few children would not achieve the required score. A full list of scores can be found in Appendix 16.

### 4.6.5.2 Vocabulary size test

Developed by Meara and Milton (2003) the X-lex 5000 test was designed for L2 learners and it measures receptive vocabulary size and it consists of both real words and pseudo-words. It focuses on vocabulary breadth and is based on the Yes/No test format wherein the participants are asked to tick the words they know. The test also includes 20 pseudo-words which were designed to look like real words but actually do not exist and they were used as foils. The X-lex tests (both English and Tamil; see Appendix 2 and 3) used in this study were the same as the ones used in Study 1. Further details of these tests have been discussed in Section 3.4.1. The only difference in the administration of these tests was that unlike Study 1, in Study 2 the researcher read the words aloud to each child for Study 2. Since other studies (e.g. Daller & Ongun, 2018) have administered the X-lex test by reading it aloud individually for each participant, the same procedure was followed for this study. This was not possible in Study 1 due to time constraints and other logistical problems.

One potential difficulty in using the X-lex tests in this study is that they were designed for L2 learners. The children in this study vary in their linguistic background: some are successive

bilinguals, and some are simultaneous bilinguals (see Section 2.1 for further details). Using these tests with children who may not be considered L2 learners could result in ceiling effects (this could be possible with the English X-lex test in this study), but the results showed that no ceiling effects were found.

### *4.6.5.3 Language Exposure and Background Questionnaire*

The questionnaire used for this study (Appendix 14) aimed to gain an overview of how often the participants used Tamil and English, in which contexts they use both languages, how long they've been in the UK and what their language preferences are. The first section of the questionnaire covered details such as age, place of birth and number of years spent in the UK. The second section asked how often the participant uses Tamil and English with different family members and with friends. The third section covered a range of common activities, such as family meals, playing with friends etc., and looked at how frequently the participant uses Tamil in these different scenarios. For the second and third sections, a 5-point Likert scale was used to assess the frequency of Tamil and English use for each question, with 0 representing "never" up to 4 representing "all the time". In this study, the questionnaire was administered orally.

### 4.6.6 Eye-tracking

This section will discuss eye-tracking and eye movements in some detail in order to provide the context for the eye-tracking experiment used in this study. According to a review by Conklin and Pellicer-Sanchez (2016), the use of eye-tracking in experimental research in L2 acquisition and applied linguistics has become increasingly popular over the last few years as it allows researchers to investigate online processing in real time instead of relying on offline measures such as judgement tasks and think aloud protocols. It is used by researchers to

record eye movements during the reading process and also to measure eye movements when looking at an image or watching a video. It is primarily used to detect and measure an eye's movements, known as saccades, and pauses, known as fixations, as well as movements back to reread a word or part of the text, known as regressions. These are basic definitions of key eye-tracking terms and will be explored in detail in the next section.

Similar to other techniques that measure online processing, there are two major assumptions, first proposed by Just and Carpenter (1980) underlying the use of eye-tracking in psychology and linguistics research: (i) the *immediacy-of-processing* assumption which stipulates that a reader process a word as soon as they encounter it without waiting for lexical and semantic ambiguities to be resolved, putting them at risk of guessing the wrong word and (ii) the *eye-mind assumption* which posits that fixating on a word or region means that that particular word or region is the one being processed during the fixation so that the amount of time spent on a region or item is reflective of the cognitive effort required to process it; thus longer fixation times indicate greater processing effort and shorter fixation times indicate less processing effort. However, these assumptions have been questioned by other researchers such as Underwood and Everatt (1992) who claim that previous knowledge and previous fixations also influence fixation times, a fact that was also acknowledged by Just and Carpenter (1983). Underwood and Everatt (1992) also present evidence that shows that words ahead of the word being fixated also influence current fixation behaviour. However, this doesn't invalidate the assumption as fixation time still reflects processing difficulty which could be influenced by the previous word or ongoing integration processes.

Conklin and Pellicer-Sanchez (2016) go on to note that eye-tracking has two significant advantages over other techniques that measure response times or reading times: (i) since eye-movements are a natural part of the reading process, eye-tracking doesn't require any external processes such as the pressing a button in self-paced reading, thus eliminating any strategic

effects and (ii) eye-tracking captures a rich, comprehensive record of eye movements and natural reading behaviour—the different movement measures allow us to measure how many times a word is encountered during reading, for how long and how many times a particular word or region is reread. Thus eye-tracking offers important advantages over self-paced reading and allows for a more in-depth investigation of the congruency effect.

### 4.6.6.1 Eye Movements and Reading

Clifton, Staub, and Rayner (2015) identify the two most robust findings that have emerged from eye movement studies on reading: (a) the fixation time on a word is shorter if the reader has a legitimate preview of the word before fixating on it and (b) when a word is easy to identify, the fixation time is shorter. This section will discuss some of the important terms related to eye movements and also the main factors that affect them.

### Key eye movement terms

Rayner (1978, 1998) was one of the first researchers to summarize much of the knowledge about typical eye movements in adult readers. In their review of the use of eye-tracking in applied linguistics research, Conklin and Pellicer-Sanchez (2016) draw on Rayner's work and provide an overview of eye movements and the associated terms relevant to this field. The basic eye movements typical of eye movement behaviour are saccades, fixations and regressions. In a discussion on the characteristics of eye movements of skilled adult readers in text reading, Keating (2014) explains that saccades are rapid eye movements that bring a new region of text into view and are usually 20 to 40 milliseconds (ms) long; the average saccade covers seven to nine characters, but it is possible for saccades to be as short as one character or as long as twenty. Leftward saccades are known as regressions and usually indicate difficulty in comprehension or occur when the eyes overshoot their intended target in the forward saccade (Rayner, 1998).  Saccades are separated by short periods where the eye

remains still, and these periods are known as fixations and these fixations are when the reader extracts meaningful information from the text. Fixations are longer than saccades, typically lasting 200 to 300 ms but can range from 50 to 500 ms depending on text complexity and reading skills of the individual (Keating, 2014). Staub and Rayner (2007) explain that it is possible to divide the text that is visible during each fixation into three regions. The "foveal" region includes text within 1° of the visual angle on either side of the fixation which is about three to four characters. Beyond this region, although the visual acuity drops considerably readers are still able to identify letters in the "parafoveal" region which extends up to 5° from the point of fixation. In the "peripheral region", readers are only aware of the general shape of the text, such as where a line ends.

Perceptual span refers to the amount of information that can be extracted from a text per fixation. It varies across different languages and is dependent on characteristics of writing systems of each language. In languages that are read from left to right, such as English, the perpetual span can extend up to 14 to 15 characters to the right of the fixation but only three to four characters to the left of the fixation (Keating, 2014). In a review of the findings of eye movement studies, Staub and Rayner (2007) conclude that the perceptual span allows readers to have access to the specific letters in a word and its phonological or sound codes, however information about its morphological composition and word meaning seem to require a direct fixation on the word.

*Eye movement measures*

Eye-tracking data is usually reported in terms of the number and length of fixations on an area of interest (AOI), otherwise known as a region of interest (ROI). Eye-tracking measures are commonly divided into early and late measures by researchers (Staub & Rayner, 2007;

Conklin, Pellicer-Sanchez & Carrol, 2018). Early measures, including first fixation duration, first pass reading time etc., are indicative of the early, more automatic stages of processing such as word recognition i.e. lexical access when a word is retrieved from the mental lexicon (Staub & Rayner, 2007). Late measures, such as total reading times, total number of fixations etc., are thought to be indicative of lexical integration and more strategic processing since they include revisits and reanalysis which reflect more conscious processing (Conklin & Pellicer-Sanchez, 2016). Total time refers to the sum of the durations of all the fixations made on a word, including any regressions and rereading time. However, it must be mentioned that these assumptions and neat divisions into early and late measures have been questioned since eye movements may not present a comprehensive overview of which stage of linguistic processing has been fully completed (Rayner & Liversedge, 2011). Using eye movements, it is perhaps easy to determine when a linguistic variable first has an effect, but it is harder to pinpoint with accuracy how long this effect plays a role in linguistic processing.

*Factors that influence eye movements*

A significant amount of eye-tracking research has focused on how words are read in isolation, but it is well-known that word processing is highly dependent on the words present in the surrounding sentence context (Libben & Titone, 2009). Thus, it is important to look at factors pertaining to individual words as well sentence-level factors while studying how words are recognised and integrated into sentences during the reading process. Rayner, Abbott and Plummer (2015) investigate these different factors in detail and a few of these relevant factors will be mentioned here: (i) word frequency is usually determined by how often a word occurs in a given corpus of texts (Juhasz & Pollastek, 2011). Word frequency influences how long readers fixate on a word, with fixation times increasing as word frequency decreases (Just & Carpenter, 1980; Rayner, 1978) In earlier studies, word frequency and word length were confounded, so Rayner and Duffy (1986) and Inhoff and

Rayner (1986) controlled for word length and manipulated for frequency and found that there was quite a large effect of frequency on first fixation and gaze duration measures. Interestingly, researchers have found that the frequency effect fades with repetition in a short text i.e. by the third occurrence of a word, there is no difference in reading times for a high or low frequency word (Rayner, Raney, & Pollastek, 1995); (ii) while word frequency is determined via corpus counts, word familiarity is determined from rating norms wherein participants rate their level of familiarity with a word. Effects of word familiarity have been found on word fixations, even when other factors have been controlled for (Juhasz & Rayner, 2003; Williams & Morris, 2004). (iii) Age-of-acquisition is an important variable that is determined from corpus counts and individual ratings. Studies have found that the effect of age-of-acquisition can be a better determinant of reading times than word frequency (Juhasz & Rayner, 2003; 2006).

### 4.6.6.2 Eye Movement Measures in This Study

For this study, the following measures were taken: first fixation duration, single fixation duration, gaze duration (early measures) and total reading time (late measure). As seen in previous eye-tracking studies, these are the most commonly used measures to study word recognition (Clifton, Staub, & Rayner, 2015) which is the main focus of this study. First fixation duration refers to the duration of the first fixation that a reader makes on a region. Gaze duration refers to the sum of durations of all the fixations a reader makes on a region, from the time they first enter it from the left to the time they exit it to the right. When a reader makes only one fixation on a region, its duration is reported as the single fixation duration. If the reader makes only one fixation on a region, the value will be the same for the single fixation duration, first fixation duration and the gaze duration. Total time refers to the sum of all the fixations on a region, including any second-pass fixations. In this study, there were two areas of interest for each item: Word 1 of the collocation and Word 2 of the collocation.

This was done so that the reading times for each word could be analysed, as well as the collocation as a whole, to study the congruency effect on the individual words as well as the whole collocation. This would also mirror the measures for Study 1. In total, there were 60 areas of interest (two for each collocation) for the whole experiment.

## 4.7 Procedure

Testing took place at the five Tamil weekend schools in a separate room with the researcher. Each child came in for a session that lasted for about 45 minutes, inclusive of a 5-minute break. Each session began with the researcher administering the language background and exposure questionnaire to the child. Next, the TOWRE test with both subtests was administered and the score was noted down. The child then completed the X-lex tests, with the researcher reading the words aloud and the child marking the words that they knew.

Finally, the child completed the main eye-tracking experiment. The child was seated at the recommended distance from the eye-tracker, placing their chin on the chin rest without any head restraints. The researcher explained the task to each child, and a 3-point calibration was used and once a successful calibration was achieved, the calibration test was repeated to validate the calibration and ensure that the tracking would be as accurate as possible. Since all the trials were only one line long, a 3-point calibration was sufficient for this experiment. Each trial of the experiment began with a display with a fixation dot in the centre of the screen to enable drift correction. This ensured that the tracker realigned with gaze location in case of any minor head movements. Once the tracker aligned with the eye, the researcher pressed a button to begin the next trial. A fixation square appeared at the point on the left of the screen where the first word of the sentence would appear. Once the tracker detected a fixation on this square, the sentence appeared, and the child read it silently. The child was allowed to complete reading it at their own pace (within a time limit of 30 seconds), after

which they pressed either the left or right button on the handheld gamepad controller to go on to the next sentence. If the child did not press one of the buttons within 30 seconds of the text appearing, the display was automatically terminated. This was the case for one participant and the data was discarded. Ten of the sentences had comprehension (yes or no) questions after them with a time limit of 15 seconds. The participants were instructed to press the left button if they thought the answer to the question was "yes" and the right button for "no". "Yes" and "no" were displayed on the screen under the questions to help the children remember which button to press for each option. All the children were given a five-minute break halfway through the experiment so that they would not get tired or restless. Another calibration was performed after the break. In total, the eye-tracking experiment lasted approximately 20 minutes per child.

## 4.8 Results

### 4.8.1 Data Cleaning

Before analysing eye-tracking data, it must be cleaned. Trials in which the tracker has lost the pupil were removed—these were trials in which no first pass reading times were recorded at all and also trials in which there were no fixations on the areas of interest. Trials in which one of the words of the target region were skipped were retained. Overall, 220 trials out of the 2400 trials (8.4 %) were deleted during the data cleaning process with a similar proportion in each condition.

### 4.8.2 Data Processing

As is usually the case with reading time data, the data were positively skewed, i.e. for all the variables $p < .05$ in the Shapiro-Wilk's test. Based on accepted practice in the field, any reading times that were 2.5 deviations above or below from the mean for each of the four conditions were excluded to eliminate outliers (Conklin & Pellicer-Sanchez, 2016) and the

data were log transformed. ==Approximately 5-7% of the data were deleted in both congruent and incongruent conditions.== This set of data were determined to be normally distributed (all $p$s > .05), for both conditions ($p$ > .0125 in the Shapiro-Wilk's test) and the analyses (parametric tests) were performed on these data.

### 4.8.3 Language Exposure/Use Questionnaire

The language use and background questionnaire was used to provide a basic overview of the language background of the children who participated in the study. It also assessed the children's own reporting of their use and exposure to both Tamil and English. The parents' responses for language use and exposure were not reported since parents were not available for questioning. The first section of the questionnaire recorded the biographical details which included the participant's place of birth. 27 of the 80 (33.75%) participants were born in the UK, 48 (60%) were born in India and the remaining five (6.25%) were born in a different country (Sri Lanka and the United Arab Emirates). The questionnaire asked which language the child considered to be his or her first language: 66 of the children (82.5%) said they considered English to be their first language and the remaining 14 children (17.5%) said they considered Tamil to be their first language. The children varied in the number of years they have lived in the UK ($M$ = 6.7, $SD$ = 3.1). None of the children reported proficient knowledge of languages apart from English and Tamil.

The second section of the questionnaire asked details about which language they used with different family members and friends. The overall trend was that the children spoke more Tamil with their mothers than with their fathers: 63.5% reported using Tamil quite often with their mother while the corresponding figure for fathers was 39.3%. Only 18.6 % of the

children reported regularly communicating with their siblings in Tamil. With regard to friends, 21% reported using Tamil frequently with their friends at Tamil school and none of them said they use Tamil at regular school. The third part of the questionnaire asked details about which contexts the children used each language in e.g. eating dinner with the family, playing with friends outside of school etc. For the questions involving activities with family (eating dinner with family, watching TV with family), 44.5% of the children reported using Tamil regularly but none of the children reported using Tamil regularly for the activities with friends (spending time with friends outside of school, talking to friends during break times at school). For both these sections, it was decided to add up the points each child gave for their Tamil use and exposure to produce an overall score that represented each child's Tamil use and exposure. The maximum score was 40 and the Tamil use/exposure scores ranged from 7 to 23 ($M$ =14, $SD$ = 4.3) (see Appendix 17 for a full list of scores).

**4.8.4 Vocabulary size tests**



*Figure 4.1* Distribution of vocabulary size scores. XLexEng = English vocabulary test;

XLexTam = Tamil vocabulary test

Figure 4.1 shows the mean raw scores of the vocabulary size test in English and Tamil. The

maximum raw score for each test is 5000. From the distribution of scores in the boxplot, it is

evident that the students have overall bigger vocabularies in English and that they vary very

widely in their Tamil vocabulary scores ($M = 2855$, $SD = 420$), but not so much for English

($M = 4275$, $SD = 545$) (see Appendices 12 and 13 for sample tests).

**4.8.5 Overall congruency effect**

As previously mentioned, the four eye-tracking measures analysed in this study were: first fixation duration (FF), single fixation duration (SF), gaze duration (Gaze), and total reading time (TT). To determine the overall congruency effect, i.e. whether the children read the congruent collocations faster than the incongruent collocations, and in line with Study 1, the entire collocation was considered as a single unit and the four eye movement measures were extracted for these units (see Appendix 18). The means and standard deviations for these reading times are listed in Table 4.1.

Table 4.1

*Mean reading times of collocations. Standard deviations in parentheses.*

|  | Congruent | Incongruent |
|---|---|---|
| First Fixation Duration (FF) | 294 (78) | 287 (56) |
| Single Fixation Duration (SF) | 291(57) | 308 (91) |
| Gaze Duration (Gaze) | 415 (124) | 411 (102) |
| Total Time (TT) | 619 (188) | 653 (153) |

Paired sample t-tests were carried out for each of the four reading measures (by-item and by-participant analyses) and a Bonferroni adjustment for multiple comparisons was applied. The differences between reading times for incongruent and congruent collocations were not statistically significant for any of the four measures, all $p$s > *.0125*.

**4.8.6 Difference in reading times for whole collocation: Individual differences**

Since there was no overall congruency effect for whole collocation, it was decided to look at whether individual differences in Tamil vocabulary and Tamil exposure influence the congruency effect on whole collocation. To do this, it was necessary to calculate the difference in reading times between the incongruent and the congruent collocations for the whole collocation. This value would show whether each individual read the congruent collocations faster or slower than they read the incongruent collocations and the size of this difference, thus making it possible to look at whether individual differences in Tamil and English vocabulary and exposure influence the congruency effect. Similar calculations have been used in the field of experimental psychology while examining the congruency effect in repetition tasks (Schmidt & Weissman, 2015) and learning and memory confounds (Weissman, Jiang, & Egner, 2014), but there is no precedent for the use of this measure in collocation processing. This difference will henceforth be referred to as Difference Congruency Effect and the new variables computed were FFDCE (Difference Congruency Effect for first fixation), SFDCE (Difference Congruency Effect for single fixation), GazeDCE (Difference Congruency Effect for gaze duration), and TTDCE (Difference Congruency Effect for total time).

A Pearson's correlation test was run to look at the correlations between the Difference Congruency Effect for the whole collocation and Tamil vocabulary scores, English vocabulary scores and Tamil use/exposure scores to determine if any of these variables had an effect on difference in reading times for the whole collocation. The correlations between the English vocabulary score and the difference in reading times for all four reading

measures—FFDCE, SFDCE, GazeDCE, TTDCE—were not significant i.e. $r(78) < .2$ and $p > .0125$ for all four measures and so linear regressions were not run.

**4.8.7 Congruency Effect on Word 1**

In their paper on the methodological aspects of using eye-tracking to study multi-word units, Carrol and Conklin (2014) conclude that perhaps the most optimal way to use eye-tracking with multiword sequences is to consider them as compositional strings as well as whole units. In light of this and since there was no overall congruency effect on the whole collocation, it was decided to analyse the data to determine whether Word 1 had a congruency effect. Although there was no overall congruency effect for the collocation, it is possible that on fixating Word 1, the children showed a small parafoveal or foveal effect. It was not possible to look at this in Study 1 since the words or chunks appeared only one-by-one in the self-paced reading experiment. In order to investigate this, reading times for Word 1 were extracted for each of the collocations for each of the four measures (see Appendix 19). Using these reading times, the mean times for Word 1 for congruent and incongruent collocations for each of the four measures were computed which are displayed in Table 4.2. Paired t-tests showed no significant differences between congruent and incongruent conditions in any of the reading time measures i.e. $p > .0125$ for all four measures. Therefore, there was no congruency effect on the Word 1 of the collocations.

Table 4.2

*Means and SDs for the reading times of Word 1*

|  | Congruent | Incongruent |
|---|---|---|
| *First Fixation (FF)* | 95 (*13*) | 93 (*16*) |
| *Single Fixation (SF)* | 101 (*22*) | 105 (*31*) |
| *Gaze Duration (Gaze)* | 150 (*34*) | 162 (*21*) |
| *Total Time (TT)* | 212 (*49*) | 206 (*69*) |

## 4.8.8 Difference in reading times for Word 1: Individual differences

Since there was no overall congruency effect for Word 1, it was decided to look at whether individual differences in Tamil vocabulary and Tamil exposure influence the congruency effect on Word 1. To do this, it was necessary to calculate the difference in reading times between Word 1 of the incongruent and the congruent collocations just as was done for the whole collocation in the previous section. This value would show whether each individual read Word 1 of the congruent collocations faster or slower than they read Word 1 of the incongruent collocations and the size of this difference, thus making it possible to look at whether individual differences in Tamil and English vocabulary and exposure influence the congruency effect on Word 1. The new variables computed were FFDCE1 (Difference Congruency Effect for first fixation), SFDCE1 (Difference Congruency Effect for single fixation), GazeDCE1 (Difference Congruency Effect for gaze duration), and TTDCE1 (Difference Congruency Effect for total time).

*Relationship between Difference Congruency Effect and Tamil vocabulary score* A
Pearson's correlation test was run to look at the correlation between the Tamil vocabulary
scores and the Difference Congruency Effect for all four measures. It was found that the
correlations between all four reading measures and Tamil vocabulary scores were positive
and significant, although the correlations were small: FFDCE1 $r(78) = .303$, $p < .01$,
SFDCE1 $r(78) = .381$, $p < .01$, GazeDCE1 $r(78) = .284$, $p < .01$ and TTDCE1 $r(78) = .239$, $p
< .01$. Based on these correlation results, simple linear regressions were run to understand the
predictive value of Tamil vocabulary scores on the congruency effect for each of the four
measures. Separate linear regressions were run instead of a multiple one due to the issue of
multicollinearity i.e. the correlations between each of the predictors were $p > .7$. All the
regressions were run with Tamil vocabulary score as the predictor and each had one of the
DCE1s as the outcome. The scatterplots for each of these linear regressions is represented in
Fig 4.2.



Fig 4.2 (a) Regression scatterplot for First Fixation Difference Congruency (Word 1) and Tamil Vocabulary

Fig 4.2 (b) Regression scatterplot for Single Fixation Difference Congruency (Word 1) and Tamil Vocabulary



Fig 4.2 (c) Regression scatterplot for Gaze Difference Congruency (Word 1) and Tamil Vocabulary

Fig 4.2 (d) Regression scatterplot for Total Time Difference Congruency (Word 1) and Tamil Vocabulary

As seen in Fig. 4.2 (a), the linear regression established that Tamil vocabulary scores

significantly predicted the Difference Congruency Effect for the first fixation duration

difference (FFDCE1) $F (1,78) = 7.87$, $p = .006$ and Tamil vocabulary scores accounted for

9.2 % of the variability of the Difference Congruency Effect in the first fixation duration.

This was also the case for single fixation difference (SFDCE1) $F (1,78) = 13.27$ $p < .001$ and

Tamil vocabulary scores accounted for 14.5 % of the variability of the Difference

Congruency Effect in the single fixation duration (Fig 4.2 (b)), as well as for gaze duration

(GazeDCE1) $F (1,78) = 6.86$, $p = .011$ and Tamil vocabulary scores accounted for 8.1 % of

the variability of the Difference Congruency Effect in gaze duration (Fig 4.2. (c)). However,

for total time, the linear regression established that Tamil vocabulary scores did not predict

the Difference Congruency Effect for the total time (TTDCE1) $F (1,78) = 4.72$, $p = .033$ as

seen in Fig 4.2 (d). This shows that Tamil vocabulary affected the congruency effect only in early processing measures.

### *Relationship between Difference Congruency Effect and English vocabulary*

A Pearson's correlation test was run to look at the correlations between the English vocabulary scores and the difference between the reading times on Word 1 for the congruent and incongruent collocations, i.e. the congruency effect, to assess whether English vocabulary had an effect on difference in reading times. The correlations between the English vocabulary score and the difference in reading times for all four reading measures— FFDCE1, SFDCE1, GazeDCE1, TTDCE1—were not significant i.e. $r(78) < .2$ and $p > .0125$ for any of the four measures and so linear regressions were not run.

### *Relationship between Difference Congruency Effect and Tamil use/exposure*

Using the Tamil use/exposure scores calculated from the questionnaire, a Pearson's correlation test was run to look at the correlation between the Tamil use/exposure scores and the Difference Congruency Effect for all four measures. It was found that although all four correlations were positive, only two were significant: small correlations for FFDCE1 ($r = .248$) and SFDCE1 ($r = .282$). The correlations for GazeDCE1 and TTDCE1 were not significant with $r = .191$ and $r = .176$ as the respective values.

Since the correlation coefficients for two of the reading measures were positive and significant, separate linear regressions were run for these two reading measures to check whether Tamil use/exposure was a predictor of the Difference Congruency Effect (DCE1). Both regressions were run with Tamil use/exposure scores as the predictor and each had a different reading measure as the outcome variable. As seen in Fig. 4.3 (b), the linear

regression established that Tamil use/exposure scores statistically significantly predicted the difference congruency effect for the single fixation difference (GazeDCE1) $F$ (1,78) = 6.73, $p$ = .011 and Tamil use/exposure scores accounted for 7.9 % of the variability of the difference congruency effect in single fixation duration. As seen in Fig. 4.3 (a), the linear regression established that Tamil use/exposure scores did not statistically significantly predict the difference congruency effect for the first fixation difference (FFDCE1) $F$ (1,78) = 5.11, $p$ = .027.



Tamil use/exposure scores

Fig 4.3 (a) Regression scatterplot for First Fixation Difference Congruency (Word 1) and Tamil Use/Exposure

Fig 4.3 (b) Regression scatterplot for Single Fixation Difference Congruency (Word 1) and Tamil Use/Exposure

## 4.8.9 Congruency Effect on Word 2

As previously mentioned in Section 4.6. as part of the discussion of different factors that affect eye movement measures, it is likely that highly predictable words have short reading times - therefore, if there is a congruency effect on Word 2 it would show that the children found Word 2 predictable. In order to assess this, reading times for Word 2 were extracted for each of the collocations for each of the four measures (see Appendix 20). Using these reading times, the mean times for Word 2 for congruent and incongruent collocations for each of the four measures was computed by participant which are displayed in Table 4.3. Paired t-tests showed no significant differences between congruent and incongruent conditions in any of the reading time measures i.e. *p > .0125* for all four measures and there was no overall congruency effect on Word 2.

Table 4.3 *Means and SDs for the reading times of Word 2*

|  | Congruent | Incongruent |
| --- | --- | --- |
| First Fixation (FF) | 287 (44) | 294 (79) |
| Single Fixation (SF) | 291 (60) | 308 (97) |
| Gaze Duration (Gaze) | 415 (116) | 411 (108) |
| Total Time (TT) | 619 (164) | 653 (191) |

**4.8.10 Difference in reading times for Word 2: Difference Congruency Effect**

Since there was no overall congruency effect, it was decided to look at whether individual differences in Tamil vocabulary and Tamil exposure influenced the congruency effect on Word 2. Similar to the procedure that was followed for DCEW1, this Difference Congruency Effect was calculated for each of the four eye movement measures for Word 1 and the new variables computed were FFDCE2 (Difference Congruency Effect for first fixation), SFDCE2 (Difference Congruency Effect for single fixation), GazeDCE2 (Difference Congruency Effect for gaze duration), and TTDCE2 (Difference Congruency Effect for total time) (see Appendix 21).

*Relationship between Difference Congruency Effect and Tamil vocabulary scores*

A Pearson's correlation test was run to look at the correlation between the Tamil vocabulary scores and the Difference Congruency Effect for all four measures. It was found that the correlations between all four reading measures and Tamil vocabulary scores were positive and significant: there was a small positive correlation for the FFDCE2 ($r = .384$) and SFDCE2 ($r = .253$) and a moderate positive correlation for GazeDCE2 ($r = .591$) and TTDCE2 ($r = .597$). Based on these correlation results, linear regressions were run to

understand the predictive value of Tamil vocabulary scores on the DCE2 for each of the four measures. Separate linear regressions were run instead of a multiple one due to the issue of multicollinearity. All the regressions were run with Tamil vocabulary score as the predictor and each had one of the DCE2s as the outcome. The scatterplots for each of these linear regressions is represented in Fig 4.4.



Fig 4.4 (a) Regression scatterplot for First Fixation Difference Congruency (Word 1) and Tamil Vocabulary

Fig 4.4 (b) Regression scatterplot for Single Fixation Difference Congruency (Word 2) and Tamil Vocabulary



Fig 4.4 (c) Regression scatterplot for Gaze Difference Congruency (Word 2) and Tamil Vocabulary

Fig 4.4 (d) Regression scatterplot for Total Time Difference Congruency (Word 2) and Tamil Vocabulary

As seen in Fig. 4.4 (a), the linear regression established that Tamil vocabulary scores statistically significantly predicted the Difference Congruency Effect for the first fixation duration difference (FFDCE2) $F (1,78) = 13.46$, $p < .001$ and Tamil vocabulary scores accounted for 14.7 % of the variability of the Difference Congruency Effect in the first fixation duration. This was also the case for the single fixation difference (SFDCE2) $F (1,78) = 5.356$, $p = 0.023$ and Tamil vocabulary scores accounted for 6.4 % of the variability of the Difference Congruency Effect in the single fixation duration (Fig. 4.4 (b)), as well as for the gaze duration (GazeDCE2) $F (1,78) = 43.186$, $p < .001$ and Tamil vocabulary scores accounted for 35.6 % of the variability of the Difference Congruency Effect in gaze duration (Fig 4.4 (c)). As seen in Fig. 4.4 (d), the linear regression established that Tamil vocabulary

scores  statistically significantly predicted the Difference Congruency Effect for the total time (TTDCE2) $F$ (1,78) = 41.952, $p < .001$ and Tamil vocabulary scores accounted for 35 % of the variability of the Difference Congruency Effect in total time.

### *Relationship between Difference Congruency Effect and English vocabulary*

A Pearson's correlation test was run to look at the correlations between the English vocabulary scores and the difference between the reading times for Word 2 for the congruent and incongruent collocations, i.e. the congruency effect, to assess whether English vocabulary had an effect on difference in reading times. The correlations between the English vocabulary score and the difference in reading times for all four reading measures—FFDCE2, SFDCE2, GazeDCE2, TTDCE2—were not significant i.e. *p > .0125* for all four measures and so linear regressions were not run.

### *Relationship between Difference Congruency Effect and English use/exposure*

Using the Tamil use/exposure scores calculated from the questionnaire, a Pearson's correlation test was run to look at the correlation between the English vocabulary scores and the Difference Congruency Effect for all four measures. It was found that the correlations between all four reading measures and English vocabulary scores were positive and significant: It was found that for three of the reading measures the correlations were positive and significant:  a moderate correlation for FFDCE2 ($r = .285$) and strong correlations for GazeDCE2 ($r = .636$) and TTDCE2 ($r = .624$). The correlation for SFDCE2 was positive but not significant ($r = .216$).
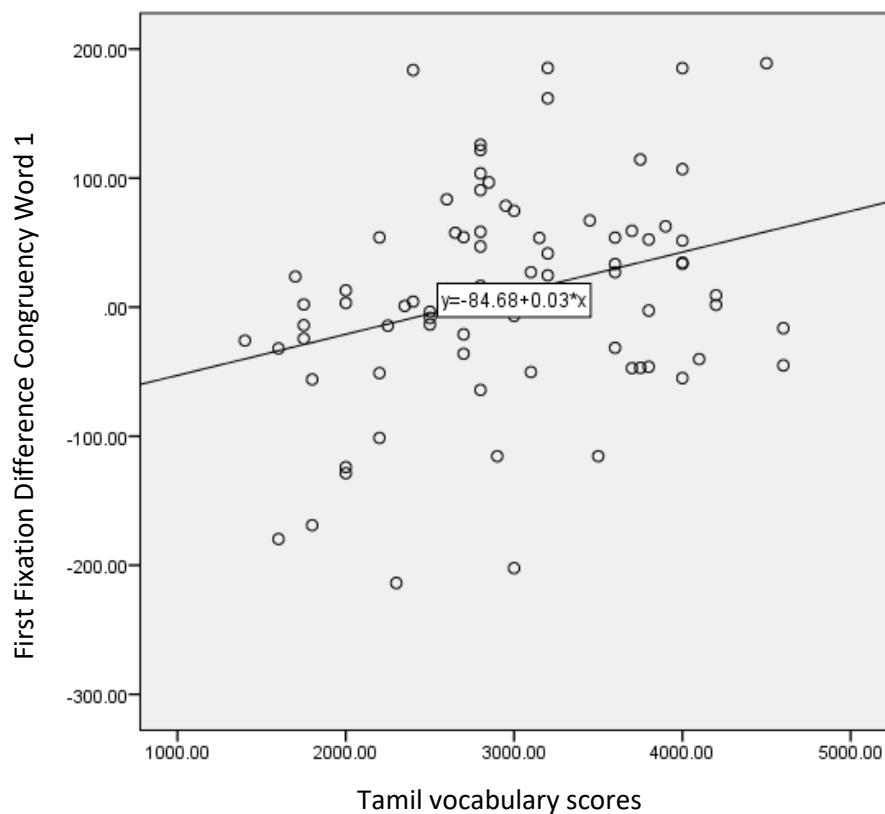
Fig 4.5 (a) Regression scatterplot for First Fixation Difference Congruency (Word 2) and English Use/Exposure



Fig 4.5 (b) Regression scatterplot for Gaze Difference Congruency (Word 2) and English Use/Exposure

Fig 4.5 (c) Regression scatterplot for Total Time Difference Congruency (Word 2) and English Use/Exposure

Since the correlation coefficients for three of the reading measures were positive and

significant, separate linear regressions were run to check whether English use/exposure is a

predictor of the Difference Congruency Effect (DCE2). The adjusted alpha level was .016

(.05/3). All the regressions were run with English use/exposure scores as the predictor and

each had a different reading measure as the outcome variable. As seen in Fig. 4.5 (a), the

linear regression established that English use/exposure scores significantly predicted the

difference congruency effect for the first fixation difference (FFDCE2) $F (1,78) = 6.914$, $p =$

.010 and English vocabulary scores accounted for 8.1 % of the variability of the difference

congruency effect in first fixation duration. As seen in Fig. 4.5 (b), the linear regression

established that English use/exposure scores could statistically significantly predict the

difference congruency effect for the gaze duration difference (GazeDCE2) $F$ (1,78) = 53.013, $p$ < .001 and English use/exposure scores accounted for 40.5 % of the variability of the difference congruency effect in gaze duration. As seen in Fig. 4.5 (c) , the linear regression established that English use/exposure scores could statistically significantly predict the difference congruency effect for the total time difference (TTDCE2) $F$ (1,78) = 53.013, $p$ < .001 and Tamil use/exposure scores accounted for 40.5 % of the variability of the difference congruency effect in gaze duration.

### 4.8.11 Vocabulary Scores Groups

Since there was no overall congruency effect but the correlations and regressions for the Difference Congruency Effect were significant for some of the reading measures, it was decided to further investigate the congruency effect. This would reveal whether there was an effect of congruency for those with high Tamil vocabulary and Tamil use/exposure scores but not for those with low Tamil vocabulary and Tamil use/exposure scores. In order to investigate the relationship between Tamil vocabulary and the extent of the congruency effect in more detail, the participants were categorised into high and low Tamil vocabulary knowledge based on their vocabulary scores. The mean and standard deviation ($M$ = 2838, $SD$ = 703) for the Tamil vocabulary scores were calculated and those participants who scored more than one standard deviation (3541) above the mean were placed in the high vocabulary group ($n$ = 16) (see Appendix 22). Those who scored more than one standard deviation (2135) below the mean were placed in the low vocabulary group ($n$ = 16) (see Appendix 23). Repeated measures ANOVA were run to determine whether there was an interaction between congruency and group. These ANOVAs were run for each of the four reading measures for the whole collocation and for Word 1 and Word 2 separately. The results of these analyses are reported below.

For the whole collocation, there was a statistically significant interaction between congruency and group only for one measure: first fixation duration, $F (1, 30) = 5.17$ $p = .008$ (see Table 4.4 for means and SDs and Fig. 4.6 for the interaction). However, post-hoc t-tests for simple effects showed that the differences between mean times for congruent and incongruent collocations in the high Tamil vocabulary group and the difference between the means for congruent and incongruent collocations in the low Tamil vocabulary group were not statistically significant, $t(15) = .685$, $p = .504$ and it was the same case for children in the high Tamil vocabulary group $t(15) = .826$, $p = .319$. The interactions between congruency and groups for the other three reading measures were not significant: single fixation $F (1, 30) = 2.49$, $p = .061$, gaze duration $F (1, 30) = .930$, $p = .901$, and total time $F (1, 30) = 3.95$, $p = .928$.

Table 4.4

*Means and SDs for the reading times of the whole collocation for both vocabulary groups*

| | High Tamil Vocab | | Low Tamil Vocab | |
|---|---|---|---|---|
| | Congruent | Incongruent | Congruent | Incongruent |
| First Fixation Duration (FF) | 224 (*76*) | 294 (*72*) | 334 (*41*) | 331 (*25*) |
| Single Fixation Duration (SF) | 275 (*105*) | 257 (*55*) | 263 (*75*) | 285 (*32*) |
| Gaze Duration (Gaze) | 453 (*245*) | 475 (*207*) | 574 (*139*) | 556 (*54*) |
| Total Time (TT) | 709 (*358*) | 823 (*485*) | 703 (*249*) | 710 (*34*) |

*Figure 4.6* Interaction between congruency and groups for first fixation for whole collocation

For Word 1, there were no statistically significant interactions between congruency and group for the four eye movement measures: for first fixation $F (1, 30) = 2.50$, $p = .124$; single fixation $F (1, 30) = 1.83$, $p = .186$; gaze duration $F (1, 30) = 1.81$, $p = .188$; and total time $F (1, 30) = 1.89$, $p = .179$, and so the tests for simple effects were not done.

For Word 2, there was a statistically significant interaction between congruency and group for two out of the four measures: first fixation and gaze duration. For first fixation, $F (1, 30) = 7.46$, $p = .01$ (see Table 4.5 for means and SDs and Fig. 4.7 for the interaction). However, post-hoc t-tests for simple effects showed that the differences between mean times for both kinds of collocations in the high Tamil vocabulary group and the difference between the

means for both kinds of collocations in the low Tamil vocabulary group were not statistically significant: for high Tamil $t(15) = 2.05$, $p = .058$ and for low Tamil $t(15) = 1.88$, $p = .079$.

Table 4.5

*Means and SDs for the reading times for Word 2 for both vocabulary groups*

|  |  | Congruent | Incongruent |
|---|---|---|---|
| **High Tamil Vocab** | First Fixation (FF) | 236 (21) | 291 (72) |
|  | Single Fixation (SF) | 260(105) | 262(35) |
|  | Gaze Duration (Gaze) | 327(108) | 374(171) |
|  | Total Time (TT) | 503(268) | 531(169) |
| **Low Tamil Vocab** | First Fixation (FF) | 294(77) | 270(53) |
|  | Single Fixation (SF) | 269(42) | 272(40) |
|  | Gaze Duration (Gaze) | 335(107) | 379(43) |
|  | Total Time (TT) | 549(163) | 561(237) |

Although the simple effects did not reach statistical significance, Table 4.5 shows that while the high Tamil vocabulary group showed longer first fixations in the incongruent condition as predicted, the low Tamil group actually showed the opposite pattern: longer first fixations on the congruent collocations. For the other three eye movement measures—single fixation, gaze duration and total time—the high Tamil vocabulary group showed longer reading times for the incongruent collocations, with the opposite effect for the low Tamil vocabulary group.

*Figure 4.7* Interaction between congruency and groups for first fixation duration for Word 2

There was also a statistically significant interaction between congruency and group for gaze duration $F(1, 30) = 5.16$, $p = .03$ (see Table 4.5 for means and SDs and Fig. 4.8 for the interaction). However, post-hoc t-tests for simple effects showed that the differences between mean times for both kinds of collocations in the high Tamil vocabulary group and the difference between the means for both kinds of collocations in the low Tamil vocabulary group were not statistically significant: for high Tamil $t(15) = 1.58$, $p = .134$ and for low Tamil $t(15) = 1.64$, $p = .122$ . Finally, the interaction between congruency and single fixation was not significant $F(1, 30) = 3.32$, $p = .078$ and the same was true for total time $F(1, 30) = 1.78$, $p = .192$.

*Figure 4.8* Interaction between congruency and groups for gaze duration for Word 2

**4.9 Discussion**

In this section, the results of the analyses from the previous section will be briefly discussed

in relation to the research questions and hypotheses presented earlier in this chapter. A more

extensive and detailed discussion of these results will be presented in Chapter 5.

**4.9.1 Hypothesis 1: Overall Congruency Effect and Congruency Effect on Word 1 and/or Word 2**

From the t-tests done with the reading times of the whole collocations, it was determined that

there was no overall congruency effect in the reading times on the whole collocations. For

this group of children, it appears that there was no influence of Tamil on the processing of

whole English collocations. Although a small (smaller than Study 1) overall congruency effect was predicted, the participants of this study were English dominant, as seen in the results of the vocabulary tests and the questionnaire. The results of the analyses also showed that there was no overall congruency effect on Word 1 or Word 2.

This absence of an overall congruency effect on the whole collocation and Word 1 and Word 2 shows that for this particular subset of EAL children, the L1 does not influence the processing of English collocations. The overall English dominance of the children and relatively low exposure to Tamil are the two factors that could have led to this result and these factors will be further explored in the next section.

### 4.9.2 Hypothesis 2: Vocabulary scores, questionnaire scores—relationship with reading times

From the results of tests to assess whether Tamil vocabulary as well as Tamil use and exposure predicted the difference in reading times between Word 2 of the congruent and incongruent collocations, it was clear that both of these factors played a significant role in the size of the congruency effect - the higher the Tamil vocabulary scores, the longer the participants took to read Word 2 of the incongruent collocations when compared to their reading times of Word 2 of the congruent collocations for all four eye-tracking measures. These findings suggest that knowledge of Tamil vocabulary does play a role in the processing of English collocations. This was further explored with the two vocabulary groups and the resulting interactions demonstrated that there may be a threshold of vocabulary knowledge above which there is a clear congruency effect. This indicated that EAL children with higher levels of vocabulary in their native language than their peers are more likely to activate their L1 when processing English collocations, although the limitations of the vocabulary tests used in this study meant that this threshold could not be explored any further. As discussed in Section 4.3, EAL children vary widely in their linguistic profiles and this is applicable to this

group of young Tamil-English bilinguals: some were born in the UK and have a fluent command of English and lower Tamil vocabulary scores, and others have been in the UK for only a few years and have a better knowledge of Tamil vocabulary. Of course, this may not be the case for all the children: it is possible that children born and raised in the UK also have high Tamil vocabulary scores due to education at Tamil school and are therefore proficient in both languages.

Similar to the results of the vocabulary analyses, the children with higher Tamil use/exposure scores took a longer time to read Word 2 of incongruent collocations when compared with Word 2 of congruent collocations. The results of the questionnaire show that there is a clear distinction between the contexts in which the children are exposed to Tamil and in which they are exposed to English: as expected, even the children who reported higher levels of Tamil use and exposure said it was only in the context of their homes and Tamil weekend schools, while they reported using and hearing English in the context of their regular schools and other forms of socialising—this is similar to the concept of domains (see Section 2.4) and how the use of different language in different domains can greatly contribute to language dominance, which influences lexical activation. From these results, it is clear that the children who use Tamil more frequently and are exposed to it more are influenced by this use and exposure while reading English collocations, and it appears that their Tamil collocational knowledge is activated, though not to the extent that it was activated for the children in Study 1. The results also show that English vocabulary scores did not predict the size of the congruency effect. The effects of language use and exposure on the processing of collocations will also be further explored in the next chapter.

### 4.9.3 Insights from results of eye-tracking measures

For the difference congruency effect on Word 1, the results showed that for the early reading measures—first fixation, single fixation, and gaze duration—Tamil vocabulary scores predicted the congruency effect, but this was not the case for total time which is a late measure. In the case of Word 2, the difference congruency effect for all four reading measures was predicted by Tamil vocabulary scores. These findings show that for Word 1, Tamil vocabulary can predict the congruency effect for the early stages of processing—this is not necessarily what we would expect since it is more probable that the congruency effect in Word 1 would be seen at the later stages of integration i.e. total time. This could be because after leaving Word 1, they most likely go to Word 2 and so by the time they go back to Word 1 they are unlikely to show any residual congruency effects. For Word 2, these findings show that the difference congruency effect is present at all stages of processing, both early and late. However, it must be noted that these eye movement measures are highly correlated and that large effects in earlier measures may be large enough to sway non-effects in later measures. With regard to Tamil exposure, for Word 1, the results showed that the Tamil use/exposure scores only predicted the difference congruency effect for one reading measure: single fixation. Since this was seen for only one out of the four measures, it can be concluded that the difference congruency effect for Word 1 is very fleeting because the children swiftly move onto Word 2. The results showed that Tamil use/exposure scores significantly predicted the size of the congruency effect for Word 2, meaning that it was present in all stages of processing and integration. The findings of this study have also contributed to our general understanding of eye movements in L2 learners and EAL children—we already know that lexical level effects such as frequency and age of exposure tend to show up quite early in the

eye movement record and these findings show that the L1 influence is present in both early and late measures at the lexical level.

This section has given an overview of what the different reading measures can tell us about the difference congruency effect. Chapter 5 will explore what these results mean for theoretical models of lexical processing and activation as well as how these results relate to other studies done in this area.

## 4.9 Limitations of the study

Due to time constraints, it was only possible to have one session with each child and this limited the type of and number of additional tests that could be done. It would have been useful to have been able to administer different kinds of vocabulary tests or proficiency tests to gain a better idea of the roles played by proficiency and vocabulary knowledge in the processing of collocations. This would have allowed for further analysis on which aspects of vocabulary and proficiency are associated most strongly with the congruency effect. Similarly, a more comprehensive and detailed assessment of individual exposure to each language could also help expand our understanding of what kind of exposure is beneficial for the acquisition of collocations.

It would have been interesting to look at the role of English in the processing of Tamil collocations with this sample group since many of them are more proficient in English than in Tamil. Due to technical constraints with the script and software, designing an experiment in Tamil wasn't possible for this study. Since the majority of the children in this study were English-dominant, it can be assumed that their knowledge of English would play a role in how they process Tamil collocations while reading.

It would have been interesting to study the processing of a particular kind of collocation such as adjective-noun collocations or adverb-adjective collocations because other studies have

shown that the L1 can influence the processing of various kinds of collocations in different ways (Yamashita & Jiang, 2010; Laufer & Waldman, 2011)—however, due to the need to control for frequency and length as well as keep the collocations simple since the sample group was primary school children, it wasn't possible to restrict the collocations to one particular kind.

**4.10 Recommendations**

As mentioned in the previous section, most of the research in collocation studies has looked at the processing of English collocations so it would be interesting for future research to look at the processing of collocations in other languages too. In particular, studies investigating cross-linguistic influence are mostly restricted to the influence of other languages on acquisition of English collocations, so it would be interesting to look at how English influences collocation acquisition in other languages.

It is also recommended that if possible, future studies incorporate a broader range of vocabulary tests, such as vocabulary depth tests, as proficiency tests to better understand how vocabulary knowledge and proficiency levels interact with cross-linguistic influence in the processing of collocations.

Finally, future research can look at cross-linguistic influence, in both directions, in other types of formulaic language such as idioms, bigrams, proverbs etc.

**4.11 Conclusion**

This study investigated the influence of Tamil on the processing of English collocations in bilingual children in the UK. It found that while there was no overall congruency effect, Tamil vocabulary scores and Tamil use and exposure significantly influenced how the children read congruent and incongruent collocations. Children with higher Tamil vocabulary

scores, as well as higher reported use of and exposure to Tamil, showed significant differences in how long they took to read the congruent and incongruent collocations—they took longer to read the incongruent collocations than the congruent ones. This suggests that at certain levels, knowledge of Tamil and exposure to it plays a role in the processing of collocations for Tamil-English bilingual children in the UK.

# Chapter 5: Discussion

Since it is well-established that language has a tendency to occur in multiword units (Schmitt, 2010), it is important to understand the influence of the L1 on L2 collocation processing and acquisition so that we can gain a better understanding of how the L1 influences L2 vocabulary acquisition. The two studies in this thesis aimed to investigate this aspect of L2 collocation processing and this chapter will discuss how the findings of this study can add to what we already know about how children who are learning a L2 process collocations. This chapter will present a detailed discussion of the results of Study 1 and Study 2 in the context of other studies done in this field and the relevant models of lexical storage, access and selection. It will also explore how the results of these studies contribute to the discussion on the bilingual mental lexicon in children. Since the methodology and stimuli for Study 1 and Study 2 were different, it is not possible to draw direct comparisons between the results of the two studies; however, general observations from the results of both studies will be discussed.

## 5.1 Congruency Effect in Study 1 and Study 2: A Brief Overview

In Study 1, the results showed there was a very large effect of congruency on the time the children took to read the whole collocations i.e. that overall, the children read the congruent collocations faster than they read incongruent collocations. Broadly speaking, this finding supports the idea of an integrated lexicon and non-selective access and retrieval (see Section 2.7) since the children appear to access their knowledge of Tamil collocations while reading and processing English collocations. Further analysis demonstrated that the children read Word 2 of the congruent collocations faster than they read Word 2 of the incongruent collocations. This shows that although Tamil and English are orthographically distant (see Section 2.14), Tamil can still influence the processing of English collocations in this group of children.

There was no overall congruency effect on the whole collocation in Study 2. This indicates that overall there was no significant activation of the L1 during the processing of English collocations for the children in Study 2. However, further analysis showed that the children with higher scores for Tamil exposure and vocabulary did show a difference in the time they took to process congruent and incongruent collocations in terms of difference congruency effect[4]: for three out of the four reading measures, the difference congruency effect was significant when considering Tamil vocabulary scores and it was significant for all four measures when considering Tamil exposure scores.

The following sections of this chapter will examine how important factors such as language dominance, language input, language exposure and vocabulary levels play a role in the extent of the congruency effect. Following this, the theoretical models relating to different features and functions of the bilingual mental lexicon will be discussed in light of the results of Study 1 and Study 2, examining how the congruency effect can be explained and understood using different aspects of these models.

## 5.2 Language Dominance

The sample groups for both Study 1 and Study 2 comprised Tamil-English bilingual children, but the children from each group came from very different contexts and so their language experiences and language backgrounds were very different. As discussed in Section 2.4, language dominance is a multidimensional construct which is measured in different ways for research purposes (Treffers-Daller, 2015; for further details see Section 2.4). Individual language dominance is one of the dimensions and refers to how individual bilinguals differ in

---

[4] Difference congruency effect: This is the value that was measured by calculating the difference in reading times between Word 1 of the incongruent and the congruent collocations. This value would show whether each individual read Word 1 of the congruent collocations faster than they read Word 2 of the incongruent collocations and the size of the difference, irrespective of whether there was a difference or not. For further details, see Section 4.8.7.

their language use and proficiency in both/all their languages (Treffers-Daller, 2015). In a discussion on individual language dominance in children, Meisel (2007) explains that key factors that affect language dominance in children are individual proficiency in each language, the status of each language in wider society, and most importantly, the levels of exposure and experience the child has with each language. Since language proficiency and exposure will be discussed in detail in other sections, this section will focus how the status of language in society influences language dominance in children and consequently, how this dominance influences the congruency effect.

## 5.2.1 Language Status in Study 1

The children in Study 1 grew up in a Tamil-speaking society even though the medium of instruction at their school is English. Due to its historical importance, regional language pride and government policy emphasis on the importance of Tamil culture and language learning, Tamil has a high status in society in the state of Tamil Nadu (Annamalai, 2005; for further details see Section 3.1). Although Chennai is a large metropolitan city and there are other languages (Malayalam, Hindi, Kannada, Telugu) spoken in very small pockets of the city, these children came from Tamil-speaking families in a Tamil-speaking part of the city so Tamil is very likely to be the only Indian language they speak and hear at home. For these children in this context, Tamil is the language of wider society, the language that is spoken at their homes and it is their main language of interaction with their friends, schoolmates and family. Data from schools and tertiary education institutes in Tamil Nadu suggest that the majority of students in English-medium institutions in Tamil Nadu lack the English proficiency needed to cope with the requirements of their education and employment opportunities (Rana, 2009; Gargesh, 2006; for further details see Section 3.1). Thus, in this context, English also has a high status but in a markedly different way: English is seen as desirable and required for upward and social mobility, but it is difficult to attain because of

lack of resources and educational infrastructure (Annamalai, 2004; for further details see Section 3.1). Although these children have incentive to learn and develop their English skills and proficiency, English is seen as a difficult language to master for children from low socioeconomic classes. Hence, for this sample group, Tamil is the dominant language, so it is not surprising that they show L1 influence in their processing of L2 collocations i.e. an overall congruency effect.

## 5.2.2 Language Status in Study 2

The children in Study 2 all currently live in the UK, although they vary in how long they have been part of this society and the amount of time they have been exposed to English. Their only exposure to their mother tongue is at home and at their weekly Tamil school sessions, as seen in the results of the questionnaire discussed in Section 4.8.3. This group of children and their families have been shown to view both English and Tamil as having a high status, although this view manifests itself in different ways. As per Canagarajah's 2008 study on the language attitudes of Tamil children and teenagers in the UK, language attrition among Tamil immigrants surpasses the typical immigrant pattern of language shift (Canagarajah, 2008; for further details see Section 4.4). Canagarajah (2006) notes a number of reasons, the most important of which is the failure of the family to pass on Tamil to the next generation—he then goes on to explain that Tamil language schools were set up in Tamil communities to combat this Tamil attrition among second and third generation immigrants (see Section 4.4). The children in this study all belonged to these Tamil language schools and likely come from families who are invested in maintaining Tamil proficiency in the younger generations, which is why they invest time and effort in sending their children to these schools. With regard to English, the children would most likely view English with a positive attitude since they are quite often surrounded by English speakers. Canagarajah (2013) also observes that first generation parents of these children view English positively for

two main reasons: (i) English is the language of the country they now live in and English proficiency is essential for progress and (ii) the Tamil communities from Tamil Nadu and Sri Lanka hold a high view of English since it is a colonial legacy that has been accepted and nativized; and this attitude is passed down through generations (for more details, see Section 4.4). These reasons could explain why English is generally held in high esteem by the Tamil diaspora in English-speaking countries such as the UK. The next section will look at the congruency effect in light of these factors.

### 5.2.3 Language dominance and congruency effects

Previous studies with children have shown that the more dominant language is more likely to have an effect on the weaker language than vice versa (Nicoladis & Gavrila, 2015; Yip & Matthews, 2000; Bernardini, 2003). In this study, it was expected that the children who had Tamil as their dominant language would show a congruency effect because there would be lexical transfer from their knowledge of Tamil collocations to their processing of English collocations. For the children in Study 2, it was expected that the congruency effect would be much smaller because they were much less likely to have Tamil as their dominant language.

The congruency effect was significant in Study 1, but there was no overall congruency effect in Study 2 for the whole collocation or for the two separate words of the collocation. Thus, for the children in Study 1 who overall had Tamil as their more dominant language, their knowledge of Tamil collocations was activated when they encountered English collocations. These children were quite familiar with Tamil collocations and so they were able to draw on their Tamil collocational knowledge when they read English collocations, which they were less familiar with. The children in Study 2, however, did not show an overall congruency effect; a majority of them (82.5%) reported their dominant language as English and this English dominance could be a contributing factor to the lack of congruency effect. However,

it must be noted that the relationship between lexical transfer and dominance is not as straightforward as it seems; Nicoladis and Gavrila (2015) note that linguistic structural differences also interact with dominance, as seen in a study with German and Italian bilinguals done by Kupisch (2007) in which the structural differences between German and Italian in German-Italian bilingual children resulted in a different pattern of use of determiners than was predicted when considering language dominance alone. For example, for the Italian-dominant children it was assumed that their knowledge of Italian determiners would influence their acquisition of German determiners, but this was not the case because of the significant structural differences between Italian and German determiners. Similarly, the structural differences between English and Tamil could play a role in the congruency effect, especially for collocations in which the order of words differs between both languages. For example, the English collocation is congruent with the Tamil equivalent "paṇattai cēmi", but in the Tamil version the noun comes first (paṇattai = money) and the verb follows it (cēmi = save). To study this, stimuli would have to be designed accordingly which is beyond the scope of this study but could be examined in future studies in this field.

## 5.3 Language Input

As discussed in Section 2.6, there is no clear consensus on exactly what kind of input best aids a child's language development, but it is agreed that the quality and quantity of language input does affect language development in bilingual children (Armon-Lotem & Meir, 2019). With regard to Tamil, the children in Study 1 received plenty of Tamil input from their parents, family members, teachers, caregivers and friends. Since most of the Tamil they hear is spoken by native speakers, they probably hear quite a few Tamil collocations used correctly and in varied contexts, especially the commonly used ones. Thus, in terms of quantity, they would be receiving adequate Tamil input containing collocations and in terms of quality they would be receiving high-quality Tamil input from native speakers (although it

is possible that native speakers have low vocabulary sizes). While the Tamil input in the context of Study 1 is relatively easy to understand even in the absence of data on relative and absolute input, the quantity and quality of English exposure these children receive is much harder to estimate. We do know that most of the children (exact details not available) in this school come from non-English speaking families, so it is safe to assume that the primary source of English input is from teachers and learning materials at school. In a study on the problems of English education in Tamil Nadu, Ponnuchamy (2012) observes that despite the state government's efforts to improve and prioritise English education at every level, low-income institutions (state-run and private) are plagued by problems such as poor school facilities, low-quality teacher training and an examination-oriented system that does not facilitate proficiency as a goal of language learning (see Section 3.1 for more details), which was also a problem in the school that participated in this study. In the context of input, Annamalai (2005) and Nehemiah (2009) single out teachers' inefficiency at communicating in English as a main reason for the low standards of ESL teaching and learning in this category of institutions in Tamil Nadu (see Section 3.1 for more details). From these studies, we can infer that the children in this context quite likely received low-quality and insufficient English language input which means that they would hear fewer collocations. It must be noted that although it is possible to draw some tentative conclusions from other studies done in the same context, a limitation of Study 1 is that this data is not available in detail.

In the context of Study 2, although detailed data on input is not available, it is possible to estimate the quantity and quality of input received from the questionnaire data on relative exposure. All the children reported hearing Tamil only at home and at Tamil school; in terms of quality, the children hear Tamil primarily from native speakers and it is likely to be high-quality as the parents and Tamil school contacts are all highly educated. In terms of input quantity, the data is harder to interpret because the questions did not specifically refer to this,

although from the results we can infer that the in general, the children had more English input than Tamil input. For English input, however, in terms of quantity it is clear that in an educational context, the children receive far more exposure to English than to Tamil and is likely to be high quality since it is in the context of UK schools. Although it is likely that some of their English input is from non-native speakers, it may be of higher quality than in India because everyone in the UK lives in an English-rich environment.

From the quantity and quality of input the children receive in Study 1, it seems that the children who received more Tamil input, i.e. more input of Tamil collocations, and less English input, i.e. less input of English collocations, tended to activate their knowledge of Tamil collocations when presented with congruent English ones. For the children in Study 2, this overall congruency effect was not present, and it could be that a relatively high level of English input is a factor that at least partially accounts for this. Wolter and Gyllstad (2013) reviewed several studies on collocational processing and concluded that while general language input is an important factor, even relatively low input of congruent collocations can provide a distinct processing advantage. The overall congruency effect seen in Study 1 supports this conclusion: while English input may have been relatively low for these children, they still showed a congruency effect. However, for a more accurate and precise understanding of the role of input in L2 collocation processing in children, more data on the quality and quantity of input is needed. This could be done with more detailed questionnaires or surveys to measure relative exposure in more detail or recording of interactions at home and at schools to measure absolute input, although this can be particularly difficult to organise (see Section 2.6).

**5.4 Language Exposure**

In vocabulary studies, it is widely acknowledged that frequency of exposure is the best predictor of effective acquisition of a lexical item and is one of the best predictors of acquisition of individual words (Nation, 2001; Schmitt, 2010, Ellis 2002a; Ellis 2002b). As discussed in Section 2.5, age of onset (AoO) is the most basic measure to determine exposure in early childhood bilingualism (for more details, see Section 2.5). In literature on reading development, age of acquisition has been found to have a significant effect on word learning i.e. the earlier a child acquires a word, the more rapidly they are likely to process it when they encounter it and the more connected it is within the semantic network (Zevin & Seidenberg, 2002). In terms of relative exposure, the children in Study 1 come under the category of successive bilinguals or ESLA (early L2 acquisition): they were exposed to Tamil from birth and were exposed to English from the time they started kindergarten at approximately four years old. In terms of relative exposure, these children would have received much more exposure to Tamil than to English and from an earlier age. This means that it would be expected that the children were more frequently exposed to Tamil collocations from an earlier age and therefore would show a congruency effect while processing English collocations in line with age of acquisition effects. This could be because results of the study showed that this overall congruency effect was present for these children, suggesting that levels of exposure in Tamil and English exposure does affect how these children process collocations. However, similar to language dominance and input, objective for language exposure is not available for Study 1 and therefore it is only possible to draw tentative conclusions.

In Study 2, some of the children were likely successive bilinguals (the ones born outside the UK) and others were simultaneous bilinguals (the ones born in the UK), However, this distinction between BFLA and ESLA for the children in Study 2 may not be as clear cut as it

seems because some of the children born in India may have had exposure to English before they moved to the UK. As noted by Armon-Lotem and Meir (2019), it can be difficult to determine age of onset for children for whom the first language is a heritage language, because it is complicated to determine. The questionnaire used in Study 2 measured the relative exposure of the children to both Tamil and English by asking them to report how much they used each language with family members, friends and in different situations. They were given a score out of 40 to indicate their level of Tamil exposure: analysis of the reading times showed that the children with higher scores for Tamil exposure showed a greater difference in some reading times for congruent and incongruent collocations i.e. Tamil exposure scores were a predictor of the difference between the reading times for both kinds of collocations.

It can be assumed that those with higher levels of Tamil exposure also heard Tamil collocations more frequently, hence acquiring them and storing them in their lexicon. This assumption is supported by studies that have looked at frequency of exposure and acquisition of collocations: Webb, Newton and Chang (2013) found that an increase of up to 15 exposures resulted in significant collocation learning gains, while Durrant and Schmitt (2010) found that in a collocation priming experiment, even one exposure had a small but significant facilitatory effect and increasing the exposures to two led to a large facilitatory effect. So, in the context of Study 2, the children who had more exposure to Tamil likely acquired more Tamil collocations and this knowledge of Tamil collocations could be what resulted in the difference congruency effect: these Tamil collocations were activated when the children read congruent English collocations. However, it must be noted that the size of the congruency effect varied as a function of children's Tamil vocabulary levels, which is why the effect was smaller than in Study 1 in which it was seen in all the children. This smaller congruency effect is likely because although for some children the degree of Tamil exposure was higher

than others, it was still relatively low—the highest score for Tamil exposure was 23 out of a maximum of 40. Additionally, Gonzalez-Fernandez and Schmitt (2015) observe that while frequency of exposure is an important factor in collocation acquisition—this is certainly the case in tightly controlled experiments and is also true in real life, we cannot measure it because generally speaking, it is impossible to ascertain how many exposures a child receives outside of a classroom environment. Further research on the different dimensions of language exposure is needed for a better understanding of how it affects the L1 influence on L2 collocation processing in children.

## 5.5 Vocabulary Scores

In Study 1, there was no correlation between the vocabulary scores and the reading times of the collocations. Reasons for this result have been discussed in detail Section 3.9.4 (Chapter 3) and one of the reasons that was taken into account for the design of Study 2 was the format of administration of the X-lex vocabulary tests. Unlike in Study 1 in which the children took the tests in groups and were asked to read the words on their own, in Study 2 the researcher individually administered the test and read out the words to each child. The Tamil vocabulary scores from Study 2 could statistically predict the size of the congruency effect for the children, but the English vocabulary scores did not have a similar effect probably because there was not enough variance in the English scores. To further investigate the relationship between the difference congruency effect and Tamil vocabulary scores, the children were split into two groups based on the mean and standard deviation of the vocabulary scores: high Tamil vocabulary and low Tamil vocabulary (for more details, see Section 4.8.10). For the whole collocation, it was found that the children in the high Tamil vocabulary group showed a larger congruency effect and this was also the case for Word 2 of the collocations— however, while the interaction between vocabulary group and congruency was significant, the pairwise comparisons were not—this could partly be due to power since the sample size

was already reduced and then split into two groups. The aim of these split analyses was to determine whether there is a particular "vocabulary threshold", above which the difference congruency effect comes into play. From the results of these analyses, it isn't possible to determine if this vocabulary threshold exists because there is no significant difference. Additionally, it is possible that just one measure of vocabulary isn't enough to capture any such threshold; it would be interesting to include other vocabulary tests that directly test vocabulary knowledge (instead of asking the child if they know the word) and perhaps tests of collocational knowledge to investigate this further.

## 5.6 Lexical Storage, Access, Activation and Selection

As discussed in the Literature Review, most models of the bilingual lexicon favour either the separate or integrated view of the bilingual lexicon but several of them have incorporated different levels of integration as well. The two main views of lexical access were also discussed in the Literature Review: the concepts of selective and non-selective access which determine if only a single language is activated during processing or if both languages are activated at the same time. Various models have differing conceptions of how exactly a lexical item is accessed, activated and selected, based on factors such as L2 proficiency, word frequencies, word exposure etc. This section will look at the results of Study 1 and Study 2 in the context of the differing views of lexical storage, access, activation and selection and which conceptual representations the results of these studies support.

### 5.6.1 Revised Hierarchical Model

In Study 1, the children showed that their L1 was activated during the processing of L2 collocations, which means that L1 collocations were activated and selected by way of lexical access when the children encountered L2 collocations. As previously mentioned, the RHM

allows for shared representations although it states that L1 and L2 lexical items are stored

separately. It is possible to assume that the shared conceptual representations could account

for the congruency effect on collocations: this would mean that for the congruent

collocations, the children had shared representations for the L1 and L2 collocations, allowing

them to read the congruent collocations faster than the incongruent collocations. However,

this doesn't account for activation at the lexical level which is what appears to be happening

in the case of these congruent collocations. Additionally, the conceptual information in the

RHM is thought to be independent of lexical information but this has plenty of experimental

contradictions (see Section 2.7.3.1). Since the RHM does not adequately explain the

congruency effect in this study, let us look at other models of that explain lexical access.


## 5.6.2 Bilingual Interactive Activation and BIA+ Models

The main obstruction for the BIA and the BIA+ models to fully explain the congruency effect

is that it assumes that orthography is shared across both the languages under consideration,

which is not the case for English and Tamil. This model posits that orthographic codes are the

first to be activated in cross-linguistic scenarios; since this is not possible across English and

Tamil, let us look at the next level of codes that are activated which are the phonological and

semantic codes. The BIA+ model suggests that these codes are activated on the basis of

subjective frequencies and other measures which mean that it is likely that L1 codes are

activated slightly before L2 codes, implying that the cross-linguistic effect would be larger

from the L1 to the L2 than for the L2 to L1. This could explain the congruency effect in

Study 1. In which there was a large congruency effect of the L1 on the L2 collocations. In the

case of Study 2, the absence of an overall congruency effect on the L2 collocations could be

accounted for by the model's stipulation that once the orthographic codes are activated, it is

possible that none of the other codes are activated. Thus, for the children in Study 2, once

their English orthographic codes are activated, all other codes are supressed which is why there is no activation of the L1 for these children. In this case, the orthographic codes are stronger than any other information or codes attached to the collocations—this is probable for these children since they are exposed to a lot more of written English than written Tamil. However, this model does not explicitly account for how proficiency plays a role in cross-linguistic influence during language processing. Although these models do not fully explain the results of the present studies, they offer an insight into a possible explanation for the role of the L1 in the processing of L2 collocations.

**5.6.3 Multilink Model**

Similar to the BIA and BIA+ models, the Multilink model also places a lot of emphasis on the orthographic aspect of word activation and goes on to explain lexical priming using word orthography. However, it also provides an explanation for semantic priming based on conceptual representations of lexical items and considers word association links to play an important role in priming. These word association links can be extended to include collocations since collocations are formed of words that are associated with each other, because they commonly occur together.

**5.6.4 Distributed Feature Model**

The DFM (Distributed Feature Model) is based on the well-supported observation that translations across languages very often capture only approximate meanings (Hofer, 2015). This is relevant to the present studies: although the English collocations chosen were verified as being either congruent or incongruent for equivalent Tamil collocations, it must be acknowledged that some of the words in the collocations are polysemous and this could have influenced how they were read and processed by the children. For example, in the collocation

*caught a cold*, the word *caught* doesn't have the same meaning as the regular use of *caught* i.e. *caught* a cold. To clarify, even though the constituent words of the chosen collocations are congruent or incongruent for the sense in which they are used in the collocations, they could have other meanings as well. Thus, for these words, the concepts attached to these individual words i.e. lexical items might only overlap partially which means that during the stage of lexical access, only the overlapping portions of the concepts would be accessed and thus activate their respective L1 and L2 lexical items. For the results of Study 1, this model would propose that the congruency effect is indicative that the congruent collocations have been stored as shared concepts, connected to both L1 and L2 lexical items and due to non-selective access, both L1 and L2 collocations are activated when the children process L2 collocations. With regard to Study 2, the absence of an overall congruency effect would suggest that the congruent collocations were not stored in the common concepts store, and instead collocations are stored in their respective language conceptual store. However, this is where this model is flawed and cannot be used to explain the further details of Study 2 results. These results show that Tamil vocabulary scores and amount of Tamil exposure affected the difference in reading time between congruent and incongruent collocations—this model does not account for factors such as these along with other crucial factors such as language proficiency, which have been shown to impact the degree and direction of cross-linguistic influence. Additionally, this model has no clear demarcations for language-specific conceptual stores which is another reason why it cannot fully explain the results of Study 2.

In the SAM (Shared Asymmetrical Model), the notion of conceptual convergence plays a key role in understanding the dynamics of L2 vocabulary acquisition. The results from Study 1 fit this model because it shows that the link between L1 and common concepts is stronger than the link between L2 and common concepts. Since this link is stronger, during the stage of lexical access, it follows that the L1 is activated when the child encounters congruent

collocations—even if they are in the L2—because of the strength of the link between the L1 and the common concepts. It is likely that since the congruent collocations are common to both the L1 and the L2, they are stored in the common conceptual store. The results of Study 2, however, which show no overall congruency effect can be explained by proposing that the children in this study did not recognise congruent collocations as being congruent i.e. even though they could be stored in the common conceptual store, these children could have stored them separately in L1 and L2 conceptual stores. If this is the case, when encountered with the collocations in English, only their English conceptual store would be activated and this is where the lexical items would be selected and activated from, thus eliminating the possibility of L1 activation. Since further analysis of the Study 2 results showed that those with higher Tamil vocabulary scores showed significantly lower reading times for congruent collocations than for incongruent collocations, for first pass reading times (first fixation, single fixation and gaze duration). This indicates that even though there was no overall congruency effect, the children in Study 2 who had high Tamil vocabulary scores were influenced by their knowledge of Tamil collocations in the initial stages of processing (as suggested by the effect on first past measures, which can be taken as indicative of early processing) of English collocations. This could mean that for these children, Tamil collocations were activated at the beginning stages of processing when they read congruent English collocations.

## 5.6.5 Modified Hierarchical Model

The MHM (Modified Hierarchical Model) aimed to address the shortcomings of the previous models mentioned here, and so let us examine whether it can adequately explain the results of the present studies. The main addition the MHM makes to existing models of the bilingual lexicon is its proposition that transitioning from explicit to implicit vocabulary knowledge is essential to complete the conceptual restructuring that it considers key to proficient L2

vocabulary acquisition. According to the MHM, this conceptual restructuring is considered to be the goal of L2 vocabulary learning in which the conceptual store is restructured to include L2-specific concepts alongside L1-specific concepts and shared concepts in the conceptual store. This conceptual store is linked to both L1 and L2 words which in turn are connected to each other. Pavlenko (2009) stresses the importance of differentiating between semantic and conceptual features and how conceptual equivalence, partial equivalence and non-equivalence determine how easily an L2 word and concept can be integrated in the mental lexicon. Although this model does not specifically mention the role of proficiency, it is understood that as a learner's L2 learning trajectory progresses they move from integrating items of conceptual equivalence to conceptual non-equivalence, thus achieving the restructuring of the conceptual store.

In terms of lexical access and activation, the results of Study 1 show that the children had integrated congruent collocations at the lexical level which is why the L1 was activated when the children processed congruent collocations, but from the design of the study it is not possible to determine whether the collocations have been integrated at the conceptual level. With regard to the results of Study 2 the results imply that even lexical integration has not occurred fully and according to the MHM model, this means that conceptual integration of collocations certainly hasn't taken place, although there is no way of determining this from the results of this study.

## 5.7 Overview

This section has brought together theoretical and experimental considerations together to examine how the results of Study 1 and Study 2 can be explained and understood by the various models of lexical representation and access, as well as dimensions of language

acquisition such as language dominance, language exposure, language input and vocabulary knowledge. The results of the study are supported by Wolter and Gyllstad's (2011) observation that when there is considerable overlap between L1 and L2, the knowledge of L1 can have a facilitative impact on L2 acquisition—the results of this study show that this can be extended to the collocational level for children too.

**5.8 Limitations and Directions for Future Research**

Since the specific limitations for Study 1 and Study 2 have been listed in their respective chapters, this section will cover the general limitations of both studies. Firstly, the use of different collocations and different sentence frames for Study 1 and Study 2 meant than the results of both studies could not be compared directly, even if adjustments were made for language exposure and vocabulary proficiency. It would be interesting to use an identical set of stimuli across both sample groups and compare the reading and processing times for congruent and incongruent collocations. Secondly, using additional assessments to measure the children's collocational knowledge in both Tamil and English could have provided more insight into the children's collocation knowledge. This wasn't possible due to the practical constraints of these studies but could be considered for future studies. Thirdly, the strength of collocations was a factor that wasn't taken into consideration for selection of collocations for the stimuli used in both experiments. This was because strength of collocations is a corpus-based measure: for Study 1, a suitable children's corpus for L2 learners isn't available and the corpus used for Study 2 (Children's Printed Word Database) did not have information on the strength of collocations. As corpora for children's language are developed, information on collocational strength is a crucial measure that should be included in order to further study how children acquire and process collocations.

In future studies, it would be interesting to study the effect of the influence of English on the processing of Tamil collocations for Tamil-English bilingual children for whom English is the dominant language. It would also be interesting to study the influence of L1 on the processing of other types of formulaic language such as pragmatisms, idioms, similes etc. Future research could also take collocation strength into account because as shown by Yamashita and Jiang (2010) this plays an important role in how learners acquire collocations and so could possibly be an important factor in understanding the cross-linguistic congruency effect in more detail.


## 5.9 Implications

To develop our understanding of how bilingual children acquire collocations, it is necessary to look at models of lexical access and processing such as the DFM and the BIA+. The results of this study indicate that these theoretical mechanisms of L1 and L2 interaction can be extended from single lexical items to collocations in these models. However, these studies have only tested lexical knowledge and not conceptual knowledge and further research is required to determine if this can be extended to the conceptual level for collocations. Conceptual knowledge is much harder to test than lexical knowledge—to test conceptual knowledge, a set of stimuli taking into consideration the conceptual meaning of each collocation would have to be designed. This would allow researchers to examine both lexical and conceptual knowledge.

In terms of teaching, it is important that child L2 learners are taught vocabulary in terms of collocations, not just in single words. A recent study by Le-Thi, Dornyei and Pellicer-Sanchez (2020) looks at how motivational methods and mental imagery can be used effectively to teach formulaic language to L2 learners and found that visualization techniques

are particularly effective. These findings can be applied to the teaching of collocations in L2 classrooms, especially for incongruent collocations that may be harder to acquire for learners who are L1-dominant. In the context of curriculum design, special focus should be given to incongruent collocations that may be more difficult to acquire. As observed by Nesselhauf (2003), the sheer number of collocations makes it impossible to include all of them in the syllabus, so emphasis should be placed on incongruent collocations in the appropriate context.

From the results of both the studies, it is clear that language exposure is an important factor in the processing of collocations. For bilingual children to attain proficiency in the L2, exposure to collocations is important, especially incongruent ones that take longer to acquire. Additionally, exposure to collocations in both languages for a bilingual child can have facilitative effects on collocation acquisition.

# Chapter 6: Conclusion

The two studies presented in this thesis examined the role of the L1 in how Tamil-English bilingual children process collocations. In this concluding chapter, a brief overview of the main differences between the methodologies, the stimuli and the sample groups will be discussed. Following this, a summary of the main findings of both studies as well as the significance of these findings will be presented to conclude this thesis.

## 6.1 Differences between Study 1 and Study 2: Methodology, Stimuli, and Sample Groups

Although both studies had the same aim, they differed with regard to the methodology, stimuli and sample group. These differences will be recapped in order to remind the reader of the differences in the findings. With regard to the methodology, both studies used different tools to measure the reading times of the collocations. In Study 1, self-paced reading measured the reading times of the collocations by recording the reading times as the children read the stimuli. While this was useful to obtain the reading times for each collocation (as a whole and as its individual constituents), it did not offer more detailed data on eye movements and also did not simulate a natural reading task. For Study 2, it was possible to use eye-tracking which offered more detailed data on the eye movements and also simulated a more natural reading task. While it would have been advantageous to use eye-tracking for Study 1, the self-paced reading task provided enough data to determine that there was an overall congruency effect for the children in Study 1. Despite an absence of an overall

congruency effect in Study 2, the use of eye-tracking enabled an analysis of which stage of processing the difference congruency effect emerges for the relevant children and it was found that it emerges at the beginning stages of processing.

The stimuli for both sample groups were designed differently for two main reasons: (i) the different nature of the sample groups and (ii) learnings from Study 1 indicated that a more thorough analysis could be conducted if the sentence frames remained the same for each pair of congruent and incongruent collocations and so this was done for Study 2. The stimuli for Study 1 were not matched for length and frequency, although this did not affect the analysis and explained in Chapter 3. However, for Study 2, the collocations were matched for length and frequency in order to avoid complications and allow for a more robust analysis.

The sample groups for both studies were both Tamil-English bilingual children but their backgrounds and language exposure were very different. The children in Study 1 had a lot more exposure to Tamil, less exposure to English and overall, they were Tamil-dominant. The children in Study 2 were largely English-dominant and had more exposure to English and less exposure to Tamil. This major difference between both sample groups is the main reason for the differences in results of both studies.

**6.2 Summary of the findings**

The most important finding of both studies is that exposure plays a crucial role in the extent of the L1 influence on the processing of congruent collocations in both sample groups. The children in Study 1 had high exposure to Tamil and this was reflected in the overall congruency effect. In Study 2, there was no overall congruency effect, perhaps because the children had lower exposure to Tamil than the children in Study 1. However, further analysis

showed that the children in Study 2 who had relatively higher levels of Tamil exposure read the congruent collocations significantly faster than they read the incongruent collocations.

Another key finding of these studies is that language dominance, which is closely linked to language exposure, is also a very important factor that affects how these bilingual children processed collocations. In Study 1, Tamil was easily activated during the reading of English collocations because the children had Tamil as their dominant language. In Study 2, 82.5 % of the children reported English as their dominant language and this is likely an important contributing factor to the absence of an overall congruency effect.

In terms of theoretical models of lexical access and processing, the data from these studies largely suggest that if the L1 is activated, it is at the beginning stages of collocational processing, which is supported by the Distributed Feature Model and the Modified Hierarchical Model: with the caveat that we know activation happens at the lexical level, but we do not have enough evidence from these studies to determine whether it happens at the conceptual level too. In terms of lexical access, the data points toward non-selective access, the extent of which is dependent on language dominance and language exposure.

**6.3 Significance of the findings**

Language exposure is a key determining factor of the congruency effect and this finding is important because it shows that just like single lexical items, increased exposure to L2 collocations is important for children to acquire them—for bilingual children, adequate exposure to both languages is important to ensure the L1 has a facilitative effect on collocation processing. Another key determining factor of the congruency effect is vocabulary size From a pedagogical perspective, the overarching finding that the L1 appears to play a significant role in the processing of L2 collocations for young learners suggests that

the L1 should not be ignored when it comes to the teaching of collocations. Teachers, as well as developers of learning materials, should give special attention to L2-only collocations. In general, collocations deserve greater emphasis in language teaching since it is clear that they play a very important role in the development of the lexical network.

More recent theoretical models of lexical access and processing posit that the L1 is activated when single lexical items in the L2 are encountered are processed, depending on factors like individual proficiency and dominance. The results of these studies show that this can be extended to the processing of L2 collocations in bilingual children. These findings are significant in the field of L2 vocabulary acquisition for children because they demonstrate that the focus should expand from single lexical items to how they are used in the context of formulaic language such as collocations. With further research into the role of the L1 in the acquisition and processing of collocations and other formulaic language, these theoretical models should be expanded to take new findings into account so that our understanding of lexical acquisition and processing is more comprehensive and not limited to single lexical items. In particular, these finding broaden our understanding of language development and lexical acquisition in bilingual children.

# References

Altarriba, J., Kroll, J. F., Sholl, A., & Rayner, K. (1996). The influence of lexical and conceptual constraints on reading mixed-language sentences: Evidence from eye fixations and naming times. Memory & Cognition, 24(4), 477-492. doi:10.3758/bf03200936

Ameel, E., Storms, G., Malt, B. C., & Sloman, S. A. (2005). How bilinguals solve the naming problem☆. Journal of Memory and Language, 53(1), 60-80. doi:10.1016/j.jml.2005.02.004

Annamalai, E. (2005). 2. Nation-building in a Globalised World: Language Choice and Education in India. Decolonisation, Globalisation, 20-37. doi:10.21832/9781853598265-004

Armon-Lotem, S., & Meir, N. (2019). The nature of exposure and input in early bilingualism. *The Cambridge Handbook of Bilingualism*, 193-212.

Armon-Lotem, S., de Jong, J., & Meir, N. (Eds.). (2015). *Assessing multilingual children: Disentangling bilingualism from language impairment*. Multilingual Matters.

Armon-Lotem, S., Joffe, S., Abutbul-Oz, H., Altman, C., & Walters, J. (2014). Language exposure, ethnolinguistic identity and attitudes in the acquisition of Hebrew as a L2 among bilingual preschool children from Russian-and English-speaking backgrounds. *Input and experience in bilingual development*, *13*, 77-98.

Arnon, I., & Snider, N. (2010). More than words: Frequency effects for multi-word phrases. *Journal of memory and language*, *62*(1), 67-82.

Ayyar, R. V. (1993). Vaidyanatha. 1993. Educational Planning and Administration in India: Retrospect and Prospect. *Journal of Educational Planning and Administration*, *7*(2), 197-214.

Babaii, E., & Ansary, H. (2001). The C-test: a valid operationalization of reduced redundancy principle? *System*, *29*(2), 209-219.

Babatsouli, E., & Nicoladis, E. (2019). The acquisition of English possessives by a bilingual child: Do input and usage frequency matter? *Journal of Child Language*, *46*(1), 170-183.

Babayiğit, S., & Shapiro, L. (2020). Component skills that underpin listening comprehension and reading comprehension in learners with English as first and additional language. *Journal of Research in Reading*, *43*(1), 78-97.

Bannard, C., & Matthews, D. (2008). Stored word sequences in language learning: The effect of familiarity on children's repetition of four-word combinations. *Psychological science*, *19*(3), 241-248.

Beauvillain, C., & Grainger, J. (1987). Accessing interlexical homographs: Some limitations of a language-selective access. *Journal of memory and language*, *26*(6), 658-672.

Bedore, L. M., Peña, E. D., Summers, C. L., Boerger, K. M., Resendiz, M. D., Greene, K.& Gillam, R. B. (2012). The measure matters: Language dominance profiles across measures in Spanish–English bilingual children. *Bilingualism (Cambridge, England)*, *15*(3), 616.

Bernardini, P. (2003). Child and adult acquisition of word order in the Italian DP. *N. Müller (ed.)*, 41-81.

Bhuvaneshwari, B., & Padakannaya, P. (2013).  Reading in Tamil: a more alphabetic and less syllabic akshara-based orthography. *South and Southeast Asian psycholinguistics*, 192.

Bialystok, E., Luk, G., Peets, K. F., & Yang, S. (2010). Receptive vocabulary differences in monolingual and bilingual children. *Bilingualism (Cambridge, England)*, *13*(4), 525.

Bod, R. (2000). Parsing with the shortest derivation. *arXiv preprint cs/0009025*.

Bod, R. (2001). What is the minimal set of fragments that achieves maximal parse accuracy? *arXiv preprint cs/0110050*.

Bogaards, P., & Laufer, B. (Eds.). (2004). *Vocabulary in a L2: Selection, acquisition, and testing* (Vol. 10). John Benjamins Publishing.

Bonk, W., & Healy, A. F. (2005). Priming effects without semantic or associative links through collocation. In *46th Annual Meeting of the Psychonomic Society, Toronto, Canada November* (pp. 10-13).

Bosch, L., & Sebastián-Gallés, N. (2001). Evidence of early language discrimination abilities in infants from bilingual environments. *Infancy*, *2*(1), 29-49.

Bosch, L., & Sebastián-Gallés, N. (2003). Simultaneous bilingualism and the perception of a language-specific vowel contrast in the first year of life. *Language and speech*, *46*(2-3), 217-243.

Bowyer-Crane, C., Fricke, S., Schaefer, B., Lervåg, A., & Hulme, C. (2017). Early literacy and comprehension skills in children learning English as an additional language and monolingual children with language weaknesses. *Reading and Writing*, *30*(4), 771-790.

Brenders, P., Van Hell, J. G., & Dijkstra, T. (2011). Word recognition in child L2 learners: Evidence from cognates and false friends. *Journal of experimental child psychology*, *109*(4), 383-396.

Brysbaert, M., & Dijkstra, T. (2006). Changing views on word recognition in bilinguals. In *Bilingualism and L2 acquisition*. Royal Academes for Science and the Arts of Belgium.

Brysbaert, M., & Duyck, W. (2010). Is it time to leave behind the Revised Hierarchical Model of bilingual language processing after fifteen years of service? *Bilingualism-Language and Cognition*, *13*(3), 359-371.

Bultena, S., Dijkstra, T., & Van Hell, J. G. (2015). Language switch costs in sentence comprehension depend on language dominance: Evidence from self-paced reading. *Bilingualism*, *18*(3), 453.

Burgoyne, K., Kelly, J. M., Whiteley, H. E., & Spooner, A. (2009). The comprehension skills of children learning English as an additional language. The British Journal of Educational Psychology, 79(4), 735–747. doi:10.1348/000709909X422530.

Burgoyne, K., Whiteley, H. E., & Hutchinson, J. M. (2011). The development of comprehension and reading-related skills in children learning English as an additional language and their monolingual, English-speaking peers. *British Journal of Educational Psychology*, *81*(2), 344-354.

Burgoyne, K., Whiteley, H. E., & Hutchinson, J. M. (2013). The role of background knowledge in text comprehension for children learning English as an additional language. *Journal of Research in Reading*, *36*(2), 132-148.

Burns, T. C., Yoshida, K. A., Hill, K., & Werker, J. F. (2007). The development of phonetic representation in bilingual and monolingual infants. *Applied Psycholinguistics*, *28*(3), 455-474.

Bybee, J. (2002). Word frequency and context of use in the lexical diffusion of phonetically conditioned sound change. *Language Variation and Change*, *14*(3), 261-290.

Bybee, J. (2006). From usage to grammar: The mind's response to repetition. *Language*, 711-733.

Byers-Heinlein, K., & Werker, J. F. (2009). Monolingual, bilingual, trilingual: infants' language experience influences the development of a word-learning heuristic. *Developmental Science*, *12*(5), 815-823.

Byers-Heinlein, K., Burns, T. C., & Werker, J. F. (2010). The roots of bilingualism in newborns. *Psychological science*, *21*(3), 343-348.

Canagarajah, A. S. (2008). Language shift and the family: Questions from the Sri Lankan Tamil diaspora 1. *Journal of Sociolinguistics*, *12*(2), 143-176.

Canagarajah, S. (2013). Reconstructing heritage language: Resolving dilemmas in language maintenance for Sri Lankan Tamil migrants. *International Journal of the Sociology of Language*, *2013*(222), 131-155.

Caramazza, A., & Brones, I. (1979). Lexical access in bilinguals. *Bulletin of the Psychonomic Society*, *13*(4), 212-214.

Carrol, G., & Conklin, K. (2020). Is all formulaic language created equal? Unpacking the processing advantage for different types of formulaic sequences. *Language and Speech*, *63*(1), 95-122.

Carrol, G., Conklin, K., & Gyllstad, H. (2016). Found in translation: The influence of the L1 on the reading of idioms in a L2. *Studies in L2 Acquisition*.

Chan, A. (2010). The Cantonese double object construction with bei2 'give'in bilingual children: The role of input. *International Journal of Bilingualism*, *14*(1), 65-85.

Chung, K. K., Liu, H., McBride, C., Wong, A. M. Y., & Lo, J. C. (2017). How socioeconomic status, executive functioning and verbal interactions contribute to early academic achievement in Chinese children. *Educational Psychology*, *37*(4), 402-420.

Chung, K. K., Liu, H., McBride, C., Wong, A. M. Y., & Lo, J. C. (2017). How socioeconomic status, executive functioning and verbal interactions contribute to early academic achievement in Chinese children. *Educational Psychology*, *37*(4), 402-420.

Cieślicka, A. B., & Heredia, R. R. (2013, May). The multiple determinants of eye movement patterns in bilingual figurative processing. In *25th APS Annual Convention, Washington, DC*.

Clifton Jr, C., Staub, A., & Rayner, K. (2007). Eye movements in reading words and sentences. In *Eye Movements* (pp. 341-371). Elsevier.

Coltheart, M., Davelaar, E., Jonasson, J. T. & Besner, D. (1977). Access to the internal lexicon. In S. Dornic (ed.), Attention and performance VI, 535-555. New York: Academic Press.

Conklin, K., & Pellicer-Sánchez, A. (2016). Using eye-tracking in applied linguistics and L2 research. *L2 Research*, *32*(3), 453-467.

Conklin, K & Carrol, G 2019, First language influence on the processing of formulaic language in a second language. in A Siyanova-Chanturia & A Pellicer-Sanchez (eds), Understanding Formulaic Language. A Second Language Acquisition Perspective. Routledge, New York, pp. 62-77

Conklin, K., & Schmitt, N. (2008). Formulaic sequences: Are they processed more quickly than nonformulaic language by native and nonnative speakers?. *Applied Linguistics*, *29*(1), 72-89.

Conklin, K., & Schmitt, N. (2012). The processing of formulaic language. *Annual Review of Applied Linguistics*, *32*.

Conklin, K., Pellicer-Sánchez, A., & Carrol, G. (2018). *Eye-tracking*. Cambridge University Press.

Carrol, G., & Conklin, K. (2014). Getting your wires crossed: Evidence for fast processing of L1 idioms in an L2. *Bilingualism: Language and Cognition*, *17*(4), 784-797.

Cook, V. (2002). Language teaching methodology and the L2 user perspective *Portraits of the L2 user*, *1*, 325.

Cook, V. J. (1992). Evidence for multicompetence. *Language Learning*, *42*(4), 557-591.

Cook, V., & Cook, V. J. (1993). *Linguistics and L2 Acquisition*. London: Macmillan.

Cristoffanini, P., Kirsner, K., & Milech, D. (1986). Bilingual lexical representation: The status of Spanish-English cognates. *The Quarterly Journal of Experimental Psychology*, *38*(3), 367-393.

Cristoffanini, P., Kirsner, K., & Milech, D. (1986). Bilingual lexical representation: The status of Spanish-English cognates. *The Quarterly Journal of Experimental Psychology*, *38*(3), 367-393.

Daller, M., & Ongun, Z. (2018). The threshold hypothesis revisited: Bilingual lexical knowledge and non-verbal IQ development. *International Journal of Bilingualism*, *22*(6), 675-694.

De Angelis, G., & Dewaele, J. M. (2009). The development of psycholinguistic research on crosslinguistic influence. *The Exploration of Multilingualism*, 63-77.

De Groot, A. M. (1992). Bilingual lexical representation: A closer look at conceptual representations. In *Advances in Psychology* (Vol. 94, pp. 389-412). North-Holland.

De Houwer, A., & Bornstein, M. (2003, April). Balancing on the tightrope: Language use patterns in bilingual families with young children. In *4th International Symposium on Bilingualism, Tempe, AZ*.

Dechert, H. W. (1983). How a story is done in a L2. *Strategies in interlanguage communication*, 175-195.

DeKeyser, R. (2008). 11 implicit and explicit learning. *The handbook of L2 acquisition*, *27*, 313.

Demie, F. (2018). English language proficiency and attainment of EAL (English as L2) pupils in England. *Journal of Multilingual and Multicultural Development*, *39*(7), 641-653.

Demie, F. (2018). English language proficiency and attainment of EAL (English as L2) pupils in England. *Journal of Multilingual and Multicultural Development*, *39*(7), 641-653.

Demie, F., Lewis, K., & Taplin, A. (2005). Pupil mobility in schools and implications for raising achievement. *Educational Studies*, *31*(2), 131-147.

DfE. 2017a. Collection of Data on Pupil Nationality, Country of Birth and Proficiency in English, Summary Report. Darlington: Department for Education. https://www.gov.uk/government/uploads/system/uploads/attachment_ data/file/665127/Data_on_pupil_nationality__country_of_birth_and_proficiency.pdf

Dijkstra, T., Timmermans, M. & Schriefers, H. (2000(b)). Cross-language effects on bilingual homo- graph recognition. *Journal of Memory and Language*, 42, 445±464.

Dijkstra, T., Van Heuven, W. J. B. & Grainger, J. (1998(a)). Simulating competitor effects with the Bilingual Interactive Activation Model. *Psychologica Belgica*, 38, 177±196

Dijkstra, T., & Rekké, S. (2010). Towards a localist-connectionist model of word translation. *The Mental Lexicon*, *5*(3), 401-420.

Dijkstra, T., & Rekké, S. (2010). Towards a localist-connectionist model of word translation. *The Mental Lexicon*, *5*(3), 401-420.

Dijkstra, T., Timmermans, M., & Schriefers, H. (2000). On being blinded by your other language: Effects of task demands on interlingual homograph recognition. *Journal of Memory and Language*, *42*(4), 445-464.

Dijkstra, T., Van Heuven, W. J., & Grainger, J. (1998). Simulating cross-language competition with the bilingual interactive activation model. *Psychologica Belgica*.

Dijkstra, T., & Heuven, W. J. (2002). The architecture of the bilingual word recognition system: From identification to decision. *Bilingualism: Language and Cognition*, 5(3), 175-197. doi:10.1017/s1366728902003012

Dong, Y., Gui, S., & MacWhinney, B. (2005). Shared and separate meanings in the bilingual mental lexicon. *Bilingualism*, *8*(3), 221.

Döpke, S. (1998). Competing language structures: The acquisition of verb placement by bilingual German-English children. *Journal of child language*, *25*(3), 555-584.

Döpke, S. (2000). Generation of and retraction from cross-linguistically motivated structures in bilingual first language acquisition. *Bilingualism Language and Cognition*, *3*(3), 209-226.

Dörnyei, Z., & Katona, L. (1992). Validation of the C-test amongst Hungarian EFL learners. *Language Testing*, *9*(2), 187-206.

Dulay, K. M., Tong, X., & McBride, C. (2017). The role of foreign domestic helpers in Hong Kong Chinese children's English and Chinese skills: A longitudinal study. *Language Learning*, *67*(2), 321-347.

Durrant, P. (2014). Corpus frequency and L2 learners' knowledge of collocations: A meta-analysis. *International Journal of Corpus Linguistics*, *19*(4), 443-477.

Durrant, P., & Schmitt, N. (2009). To what extent do native and non-native writers make use of collocations? *International Review of Applied Linguistics in Language Teaching*, *47*(2), 157-177.

Durrant, P., & Schmitt, N. (2010). Adult learners' retention of collocations from exposure. *L2 research*, *26*(2), 163-188.

Eckes, T., & Grotjahn, R. (2006). A closer look at the construct validity of C-tests. *Language Testing*, *23*(3), 290-325.

Ellis, N. C. (2002a). Frequency effects in language processing: A review with implications for theories of implicit and explicit language acquisition. *Studies in L2 Acquisition*, *24*(2), 143-188.

Ellis, N. C. (2002b). Reflections on frequency effects in language processing. *Studies in L2 Acquisition*, 297-339.

Ellis, R. (2005). Measuring implicit and explicit knowledge of a L2: A psychometric study. *Studies in L2 acquisition*, *27*(2), 141-172.

Eyckmans, J. (2004). *Measuring receptive vocabulary size. Reliability and validity of the Yes/No vocabulary test for French-speaking learners of Dutch*. Utrecht: LOT.

Fedele, E., & Kaiser, E. (2012). Comprehension of Anaphora and Cataphora in Italian: Comparing Null and Overt Pronouns. In *Poster presented at the Architectures and Mechanisms for Language Processing 2012 Conference, Riva del Garda–Italy*.

Fernández, B. G., & Schmitt, N. (2015). How much collocation knowledge do L2 learners have? The effects of frequency and amount of exposure. *ITL-international journal of applied linguistics*, *166*(1), 94-126.

Fillmore, L. W. (1991). When learning a L2 means losing the first. *Early childhood research quarterly*, *6*(3), 323-346.

Fillmore, L. W. (1991). When learning a L2 means losing the first. *Early childhood research quarterly*, *6*(3), 323-346.

Fitzpatrick, T., & Clenton, J. (2010). The challenge of validation: Assessing the performance of a test of productive vocabulary. *Language Testing*, *27*(4), 537-554.

French, R. M., & Jacquet, M. (2004). Understanding bilingual memory: Models and data. *Trends in Cognitive Sciences*, *8*(2), 87-93.

Frost, R. (2005). *Orthographic Systems and Skilled Wosrd Recognition Processes in Reading.* In M. J. Snowling & C. Hulme (Eds.), *Blackwell Handbook of Developmental Psychology. The Science of reading: A Handbook* (p. 272–295). Blackwell

Gargesh, R. (2006, June). Language issues in the context of higher education in India. In *symposium on Language Issues in English-Medium Universities across Asia, University of Hong Kong* (pp. 8-9).

Garnsey, S. M., Pearlmutter, N. J., Myers, E., & Lotocky, M. A. (1997). The contributions of verb bias and plausibility to the comprehension of temporarily ambiguous sentences. *Journal of Memory and Language*, *37*(1), 58-93.

Garnsey, S. M., Pearlmutter, N. J., Myers, E., & Lotocky, M. A. (1997). The contributions of verb bias and plausibility to the comprehension of temporarily ambiguous sentences. *Journal of Memory and Language*, *37*(1), 58-93.

Gasser, M., & McDermott, E. (1990, May). An architecture for practical delegation in a distributed system. In *Proceedings. 1990 IEEE Computer Society Symposium on Research in Security and Privacy* (pp. 20-20). IEEE Computer Society.

Gathercole, V. C. M., & Hoff, E. (2007). Input and the acquisition of language: Three questions. *Blackwell Handbook of Language Development*, 107-127.

Geeslin, K. L. (2003). A comparison of copula choice: Native Spanish speakers and advanced learners. *Language Learning*, *53*(4), 703-764.

Genesee, F. (2000). Early bilingual language development: One language or two. *The Bilingualism Reader*, 327-343.

Genesee, F. (2002). Portrait of the bilingual child. *Perspectives on the L2 user*, 170-196.

Geva, E., & Siegel, L. S. (2000). Orthographic and cognitive factors in the concurrent development of basic reading skills in two languages. *Reading and Writing*, *12*(1-2), 1-30.

Geva, E., Wade-Woolley, L., & Shany, M. (1993). The concurrent development of spelling and decoding in two different orthographies. *Journal of Reading Behavior*, *25*(4), 383-406.

Gibbs, R. W., Bogdanovich, J. M., Sykes, J. R., & Barr, D. J. (1997). Metaphor in idiom comprehension. *Journal of Memory and Language*, *37*(2), 141-154.

Gollan, T. H., Forster, K. I., & Frost, R. (1997). Translation priming with different scripts: Masked priming with cognates and noncognates in Hebrew–English bilinguals. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *23*(5), 1122.

Grodner, D., & Gibson, E. (2005). Consequences of the serial nature of linguistic input for sentenial complexity. *Cognitive Science*, *29*(2), 261-290.

Grosjean, F. (1989). Neurolinguists, beware! The bilingual is not two monolinguals in one person. *Brain and language*, *36*(1), 3-15.

Grosjean, F. (1992). Another view of bilingualism. In *Advances in Psychology* (Vol. 83, pp. 51-62). North-Holland.

Grosjean, F. (2008). *Studying Bilinguals*. Oxford University Press.

Grosjean, F. (2010). *Bilingual*. Harvard University Press.

Grosjean, F., & Li, P. (2013). *The Psycholinguistics of Bilingualism*. John Wiley & Sons.

Grotjahn, R. (2002). Konstruktion und Einsatz von C-Tests: Ein Leitfaden für die Praxis. *Der C-Test. Theoretische Grundlagen und praktische Anwendungen*, *4*(2002), 211-225.

Gutiérrez-Clellen, V. F., & Kreiter, J. (2003). Understanding child bilingual acquisition using parent and teacher reports. *Applied psycholinguistics*, *24*(2), 267.

Gyllstad, H., & Wolter, B. (2014). Processing L2 Word Combinations: What Role Does Degree of Semantic Transparency Play?. In *AAAL, 2014*.

Harsch, C., & Hartig, J. (2015). What are we aligning tests to when we report test alignment to the CEFR?. *Language Assessment Quarterly*, *12*(4), 333-362.

Hastings, A.J. 2002: Error analysis of an English C-Test: evidence for integrated processing. In Grotjahn, R., editor, Der C-Test: theoretische Grundlagen und praktische Anwendungen [The C-test: theoretical foundations and practical applications]. Volume 4. Bochum: AKS-Verlag, 53–66

He, X., & Kaiser, E. (2012). Is there a difference between 'You'and 'I'? A psycholinguistic investigation of the Chinese reflexive ziji. *University of Pennsylvania Working Papers in Linguistics*, *18*(1), 12.

Henriksen, B. (2013). Research on L2 learners' collocational competence and development–a progress report. *C. Bardel, C. Lindqvist, & B. Laufer (Eds.) L*, *2*, 29-56.

Heredia, R. R., & Cieślicka, A. B. (2014). Bilingual memory storage: Compound-coordinate and derivatives. In *Foundations of bilingual memory* (pp. 11-39). Springer, New York, NY.

Hestetræet, T. I. (2018). Vocabulary teaching for young learners. In *The Routledge handbook of teaching English to young learners* (pp. 220-233). Routledge.

Hofer, B. (2015). *On the dynamics of early multilingualism: A psycholinguistic study* (Vol. 13). Walter de Gruyter GmbH & Co KG.

Hong, Y. Y., Morris, M. W., Chiu, C. Y., & Benet-Martinez, V. (2000). Multicultural minds: A dynamic constructivist approach to culture and cognition. *American Psychologist*, *55*(7), 709.https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/808742/Attainment_of_EAL_pupils.pdf

Hulstijn, J. H. (2007). The shaky ground beneath the CEFR: Quantitative and qualitative dimensions of language proficiency. *The Modern Language Journal*, *91*(4), 663-667.

Hulstijn, J. H. (2007). The shaky ground beneath the CEFR: Quantitative and qualitative dimensions of language proficiency. *The Modern Language Journal*, *91*(4), 663-667.

Hunston, S. (2002). *Corpora in applied linguistics*. Cambridge University Press.

Hutchinson, J. (2018). *Educational outcomes of children with English as an additional language*. The Bell Foundation. https://dera.ioe.ac.uk/31500/1/EAL_Educational-Outcomes_EPI-1.pdf

Hyltenstam, K. (1977). Implicational patterns in interlanguage syntax variation. *Language Learning*, *27*(2), 383-410.

Inhoff, A. W., & Rayner, K. (1986). Parafoveal word processing during eye fixations in reading: Effects of word frequency. *Perception & Psychophysics*, *40*(6), 431-439.

Irujo, S. (1993). Stearing clear: Avoidance in the production of idioms. *IRAL: International Review of Applied Linguistics in Language Teaching*, *31*(3), 205.

Jegerski, J. (2014). Self-paced reading. In J. Jegerski & B. VanPatten (Eds.), Research methods in L2 psycholinguistics (pp. 20-49). New York: Routledge.

Jiang, N. (2000). Lexical representation and development in a L2. *Applied Linguistics*, *21*(1), 47-77.

Jiang, N., & Forster, K. I. (2001). Cross-language priming asymmetries in lexical decision and episodic recognition. *Journal of Memory and Language*, *44*(1), 32-51.

Juffs, A., & Harrington, M. (1995). Parsing Effects in L2 Sentence Processing: Subject and Object Asymmetries in wh-Extraction. *Studies in L2 Acquisition*, 483-516.

Juhasz, B. J., & Pollatsek, A. (2011). *Lexical influences on eye movements in reading.* In S. P. Liversedge, I. D. Gilchrist, & S. Everling (Eds.), *Oxford library of psychology. The Oxford handbook of eye movements* (p. 873–893). Oxford University Press.

Juhasz, B. J., & Rayner, K. (2003). Investigating the effects of a set of intercorrelated variables on eye fixation durations in reading. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *29*(6), 1312.

Juhasz, B. J., & Rayner, K. (2006). The role of age of acquisition and word frequency in reading: Evidence from eye fixation durations. *Visual Cognition*, *13*(7-8), 846-863.

Just, M. A., & Carpenter, P. A. (1980). A theory of reading: From eye fixations to comprehension. *Psychological review*, *87*(4), 329.

Kahn-Horwitz, J., Schwartz, M., & Share, D. (2011). Acquiring the complex English orthography: a triliteracy advantage?. *Journal of Research in Reading*, *34*(1), 136-156.

Keating, G. D. (2014). Eye-tracking with text. *Research methods in L2 psycholinguistics*, 69-92.

Kim, Y. S. (2009). Crosslinguistic influence on phonological awareness for Korean–English bilingual children. *Reading and Writing*, *22*(7), 843.

Klein-Braley, C. (1997). C-Tests in the context of reduced redundancy testing: An appraisal. *Language Testing*, *14*(1), 47-84.

Klein-Braley, C., & Raatz, U. (1984). A survey of research on the C-Test1. *Language Testing*, *1*(2), 134-146.

Kolers, P. A. (1963). Interlingual word associations. *Journal of verbal learning and verbal behavior*, *2*(4), 291-300.

Kroll, J. F., & De Groot, A. M. (Eds.). (2009). *Handbook of bilingualism: Psycholinguistic approaches*. Oxford University Press.

Kroll, J. F., & Stewart, E. (1994). Category interference in translation and picture naming: Evidence for asymmetric connections between bilingual memory representations. *Journal of memory and language*, *33*(2), 149.

Kroll, J. F., & Tokowicz, N. (2005). *Models of bilingual representation and processing: Looking back and to the future*. Oxford University Press.

Kroll, J. F., Van Hell, J. G., Tokowicz, N., & Green, D. W. (2010). The Revised Hierarchical Model: A critical review and assessment. *Bilingualism (Cambridge, England)*, *13*(3), 373.

Kuiper, K. (2000). On the linguistic properties of formulaic speech. *Oral Tradition,15(2),* 279-305

Kuiper, K. (2004). Formulaic performance in conventionalised varieties of speech. *Formulaic sequences: Acquisition, processing and use*, 37-54.

Kuiper, K., Columbus, G., & Schmitt, N. (2009). The acquisition of phrasal vocabulary. In *Language acquisition* (pp. 216-240). Palgrave Macmillan, London.

Kupisch, T. (2007). Determiners in bilingual German-Italian children: What they tell us about the relation between language influence and language dominance. *Bilingualism: Language and Cognition*, *10*(1), 57-78.

Laufer, B. (2000, August). Electronic dictionaries and incidental vocabulary acquisition: Does technology make a difference. In *EURALEX* (pp. 849-854). Stuttgart, Germany: Stuttgart University.

Laufer, B., & Waldman, T. (2011). Verb-noun collocations in L2 writing: A corpus analysis of learners' English. *Language learning*, *61*(2), 647-672.

Lee, J. Y. (2019). The Use of Verb-Noun Collocations in Korean EFL Writings at Different L2 Proficiency Levels. *Korean Journal of Applied Linguistics*, *35*(1), 51-78.

Leopold, W. F. (1939). *Speech development of a bilingual child: A linguist's record* (No. 11). Northwestern University.

Lesaux, N. K., & Siegel, L. S. (2003). The development of reading in children who speak English as a L2. *Developmental psychology*, *39*(6), 1005.

Lesaux, N. K., Crosson, A. C., Kieffer, M. J., & Pierce, M. (2010). Uneven profiles: Language minority learners' word reading, vocabulary, and reading comprehension skills. *Journal of applied developmental psychology*, *31*(6), 475-483.

Lesaux, N. K., Rupp, A. A., & Siegel, L. S. (2007). Growth in reading skills of children from diverse linguistic backgrounds: Findings from a 5-year longitudinal study. *Journal of Educational Psychology*, *99*(4), 821.

Lewis, M. (1993). *The lexical approach* (Vol. 1, p. 993). Hove: Language teaching publications.

Libben, M. R., & Titone, D. A. (2009). Bilingual lexical access in context: evidence from eye movements during reading. *Journal of Experimental Psychology: Learning, memory, and cognition*, *35*(2), 381.

Lin, M., & Leonard, S. (2012). *Dictionary of 1000 Chinese idioms*. Hippocrene Books, Incorporated.

Long, M. H., & Sato, C. (1984). Methodological issues in interlanguage studies: An interactionist perspective. *Interlanguage*, *253279*.

Malt, B. C., & Sloman, S. A. (2003). Linguistic diversity and object naming by non-native speakers of English. *Bilingualism: Language and cognition*, *6*(1), 47-67.

Marchman, V. A., Martínez, L. Z., Hurtado, N., Grüter, T., & Fernald, A. (2017). Caregiver talk to young Spanish-English bilinguals: comparing direct observation and parent-report measures of dual-language exposure. *Developmental science*, *20*(1), e12425.

Marian, V., & Neisser, U. (2000). Language-dependent recall of autobiographical memories. *Journal of Experimental Psychology: General*, *129*(3), 361.

Marini, A., & Fabbro, F. (2007). Psycholinguistic models of speech production in Bilingualism and Multilingualism. *Speech and language disorders in bilinguals*, 47-67.

Marsden, E., Thompson, S., & Plonsky, L. (2018). A methodological synthesis of self-paced reading in L2 research. *Applied Psycholinguistics*, *39*(5), 861-904.

Masterson, J., Stuart, M., Dixon, M., & Lovejoy, S. (2010). Children's printed word database: Continuities and changes over time in children's early reading vocabulary. *British Journal of Psychology*, *101*(2), 221-242.

McCardle, P. D., & Hoff, E. (Eds.). (2006). *Childhood bilingualism: Research on infancy through school age* (Vol. 7). Multilingual matters.

McLaughlin, B. (1995). Fostering L2 Development in Young Children: Principles and Practices. UC Berkeley: Center for Research on Education, Diversity and Excellence. Retrieved from https://escholarship.org/uc/item/23s607sr

McLaughlin, B. (Ed.). (2013). *L2 acquisition in childhood: Volume 2: School-age Children.* Psychology Press.

Meara, P. (1982). Vocabulary acquisition: A neglected aspect of language learning. *Surveys I: Eight state-of-the-art articles on key areas in language teaching*, 100-126.

Meara, P. M., & Milton, J. (2003). *X-lex: The Swansea levels test*. Express Publishing.

Meisel, J. M. (2004). The bilingual child. *The Handbook of Bilingualism*, *91*, 113.

Meisel, J. M. (2007). The weaker language in early child bilingualism: Acquiring a first language as a L2? *Applied Psycholinguistics*, *28*(3), 495.

Milton, J. (2006). X-Lex: The Swansea vocabulary levels test. In *Proceedings of the 7th and 8th Current Trends in English Language testing (CTELT) Conference* (Vol. 4, pp. 29-39). TESOL Arabia, UAE.

Milton, J. (2010). The development of vocabulary breadth across the CEFR levels. *Communicative proficiency and linguistic development: Intersections between SLA and language testing research*, 211-232.

Mitchell, D. C., & Green, D. W. (1978). The effects of context and content on immediate processing in reading. *Quarterly Journal of Experimental Psychology*, *30*(4), 609-636.

Mok, P. P. (2011). The acquisition of speech rhythm by three-year-old bilingual and monolingual children: Cantonese and English. *Bilingualism: Language and Cognition*, *14*(4), 458-472.

Montrul, S. (2013). How "native" are heritage speakers. *Heritage Language Journal*, *10*(2), 15-39.

Montrul, S. A. (2008). *Incomplete acquisition in bilingualism: Re-examining the age factor* (Vol. 39). John Benjamins Publishing.

Moon, R. (1997). Vocabulary connections: Multi-word items in English. *Vocabulary: Description, acquisition and pedagogy*, *40*, 63.

Müller, A., & Daller, M. (2019). Predicting international students' clinical and academic grades using two language tests (IELTS and C-test): A correlational research study. *Nurse education today*, *72*, 6-11.

Müller, N., & Hulk, A. (2001). Crosslinguistic influence in bilingual language acquisition: Italian and French as recipient languages. *Bilingualism: Language and cognition*, *4*(1), 1-21.

Müller, N., & Pillunat, A. (2008). Balanced bilingual children with two weak languages. *First language acquisition of morphology and syntax: Perspectives across languages and learners*, *45*, 269-294.

Müller-Lancé, J. (2003). A strategy model of multilingual learning. In *The multilingual lexicon* (pp. 117-132). Springer, Dordrecht.

Murphy, V. A. (2014). *L2 Learning in the Early School Years: Trends and Contexts-Oxford Applied Linguistics*. Oxford University Press.

Nakamura, J., & Quay, S. (2012). The impact of caregivers' interrogative styles in English and Japanese on early bilingual development. *International Journal of Bilingual Education and Bilingualism*, *15*(4), 417-434.

NALDIC. (2014). The national audit of English as an additional language training and development provision. London: National Association for Language Development in the Curriculum.

Nation, I. S. P. (2001). Learning vocabulary in another language. 2003. *Cambridge. Cambridge*.

Nattinger, J. R., & DeCarrico, J. S. (1992). *Lexical phrases and language teaching*. Oxford University Press.

Nazari, N. (2013). The effect of implicit and explicit grammar instruction on learners' achievements in receptive and productive modes. *Procedia-Social and Behavioral Sciences*, *70*, 156-162.

Nehemiah, S. (2009). Role of English as a tool for communication in Tamil society. *Language in India*, *9*(8), 10.

Nesselhauf, N. (2003). The use of collocations by advanced learners of English and some implications for teaching. *Applied linguistics*, *24*(2), 223-242.

Nesselhauf, N. (2005). *Collocations in a learner corpus* (Vol. 14). Amsterdam: John Benjamins.

Nicoladis, E. (2006). Cross-linguistic transfer in adjective-noun strings by preschool bilingual children. *Bilingualism*, *9*(1), 15.

Nicoladis, E. (2012). Cross-linguistic influence in French–English bilingual children's possessive constructions. *Bilingualism: Language and Cognition*, *15*(2), 320-328.

Nicoladis, E., & Gavrila, A. (2015). Cross-linguistic influence in Welsh–English bilingual children's adjectival constructions. *Journal of Child Language*, *42*(4), 903-916.

Nicoladis, E., & Secco, G. (2000). The role of a child's productive vocabulary in the language choice of a bilingual family. *First Language*, *20*(58), 003-28.

Northbrook, J., & Conklin, K. (2019). Is what you put in what you get out?—Textbook-derived lexical bundle processing in beginner English learners. *Applied Linguistics*, *40*(5), 816-833.

Nurmukhamedov, U. (2015). An evaluation of collocation tools for L2 writers. *Northern Arizona University*.

O'Dell, F., Read, J., & McCarthy, M. (2000). *Assessing vocabulary*. Cambridge university press.

Odlin, T. (2003). Cross-linguistic influence. *The handbook of L2 acquisition*, 436-486.

Ortega, L. (2009). Crosslinguistic influences. *L2 Acquisition*, 31-54.

Paradis, J., & Genesee, F. (1996). Syntactic acquisition in bilingual children: Autonomous or interdependent? *Studies in L2 acquisition*, 1-25.

Paradis, J., Nicoladis, E., Crago, M., & Genesee, F. (2011). Bilingual children's acquisition of the past tense: A usage-based approach. *Journal of child language*, *38*(3), 554-578.

Paradis, J., Nicoladis, E., Crago, M., & Genesee, F. (2011). Bilingual children's acquisition of the past tense: A usage-based approach. *Journal of child language*, *38*(3), 554-578.

Pavlenko, A. (2003). " I never knew I was a bilingual": Reimagining teacher identities in TESOL. *Journal of Language, Identity, and education*, *2*(4), 251-268.

Pavlenko, A. (2009). Conceptual representation in the bilingual lexicon and L2 vocabulary learning. *The bilingual mental lexicon: Interdisciplinary approaches*, 125-160.

Pavlenko, A., & Driagina, V. (2007). Russian emotion vocabulary in American learners' narratives. *The Modern Language Journal*, *91*(2), 213-234.

Pawley, A., & Syder, F. H. (1983). Two puzzles for linguistic theory: Nativelike selection and nativelike fluency. *Language and communication*, *191*, 225.

Pearson, B. Z., , S. C., Lewedeg, V., & Oller, D. K. (1997). The relation of input factors to lexical learning by bilingual infants. *Applied psycholinguistics*, *18*(1), 41-58.

Pearson, B. Z., Fernández, S. C., & Oller, D. K. (1993). Lexical development in bilingual infants and toddlers: Comparison to monolingual norms. *Language learning*, *43*(1), 93-120.

Pearson, B. Z., Fernandez, S. C., Lewedeg, V., & Oller, D. K. (1997). The relation of input factors to lexical learning by bilingual infants. *Applied psycholinguistics*, *18*(1), 41-58.

Pearson, P. D., Hiebert, E. H., & Kamil, M. L. (2007). Vocabulary assessment: What we know and what we need to learn. *Reading research quarterly*, *42*(2), 282-296.

Peirce, J. W. (2007). PsychoPy—psychophysics software in Python. *Journal of neuroscience methods*, *162*(1-2), 8-13.

Peirce, J. W. (2009). Generating stimuli for neuroscience using PsychoPy. *Frontiers in neuroinformatics*, *2*, 10.

Peña, E. D., Bedore, P., Grosjean, F., & Byers-Heinlein, K. (2018). Assessing perception and comprehension in bilingual children, without and with speech and language impairment. *The Listening Bilingual*, 220-243.

Peters, A. M. (1983). *The units of language acquisition* (Vol. 1). CUP Archive.

Pinker, S., & Bloom, P. (1990). Natural selection and natural language.

Place, S., & Hoff, E. (2011). Properties of dual language exposure that influence 2-year-olds' bilingual proficiency. *Child development*, *82*(6), 1834-1849.

Place, S., & Hoff, E. (2016). Effects and noneffects of input in bilingual environments on dual language skills in 2½-year-olds. *Bilingualism*, *19*(5), 1023.

Pliatsikas, C., & Marinis, T. (2013). Processing empty categories in a L2: When naturalistic exposure fills the (intermediate) gap. *Bilingualism: Language and Cognition*, *16*(01), 167-182.

Ponnuchamy, G. (2012). *School English-as-a-Second-Language experiences of students at tertiary institutions in Tamil Nadu, India: A phenomenological study* (Doctoral dissertation, University of Phoenix).

Poulin-Dubois, D., & Goodz, N. (2001). Language differentiation in bilingual infants: Evidence from babbling. *Trends in bilingual acquisition*, *1*, 95-106.

Prabhu, N. S. (1984). Procedural syllabuses. *Trends in language syllabus design*, 272-280.

Raatz, U., & Klein-Braley, C. (1981). The C-Test--A Modification of the Cloze Procedure.

Rana, S. (2009). Teaching language through literary texts in the ESL classroom. *Language in India*, *9*(6), 7.

Rayner, K. (1978). Eye movements in reading and information processing. *Psychological bulletin*, *85*(3), 618.

Rayner, K. (1998). Eye movements in reading and information processing: 20 years of research. *Psychological bulletin*, *124*(3), 372.

Rayner, K., & Duffy, S. A. (1986). Lexical complexity and fixation times in reading: Effects of word frequency, verb complexity, and lexical ambiguity. *Memory & cognition*, *14*(3), 191-201.

Rayner, K., & Liversedge, S. P. (2011). Linguistic and cognitive influences on eye movements during reading.

Rayner, K., Abbott, M. J., & Plummer, P. (2015). Individual differences in perceptual processing and eye movements in reading. *Handbook of individual differences in reading: Reader, text, and context*, 348-363.

Rayner, K., Raney, G. E., & Pollatsek, A. (1995). Eye movements and discourse processing.

Rebuschat, P., & Williams, J. (2009). Implicit learning of word order. In *Proceedings of the Annual Meeting of the Cognitive Science Society* (Vol. 31, No. 31).

Ross, M., Xun, W. E., & Wilson, A. E. (2002). Language and the bicultural self. *Personality and Social Psychology Bulletin*, *28*(8), 1040-1050.

Sahlin, B. H., Harding, M. G., & Seamon, J. G. (2005). When do false memories cross language boundaries in English—Spanish bilinguals?. *Memory & Cognition*, *33*(8), 1414-1421.

Salomon, R. (2000). Typological observations on the Indic script group and its relationship to other alphasyllabaries. *Studies in the Linguistic Sciences, 30*(1), 87-102. Retrieved from https://www.ideals.illinois.edu/bitstream/handle/2142/9640/SLS2000v30.1-10Salomon.pdf

Schachter, J. (1988). L2 acquisition and its relationship to Universal Grammar. *Applied linguistics*, *9*(3), 219-235.

Schmitt, N. (2008). Instructed L2 vocabulary learning. *Language teaching research*, *12*(3), 329-363.

Schmitt, N. (2010). *Researching vocabulary: A vocabulary research manual*. Springer.

Schmitt, N., & Carter, R. (2004). Formulaic sequences in action: An introduction (N. Schmitt, Ed.). In *Formulaic sequences: Acquisition, Processing and Use* (pp. 1-22). Seiten: John Benjamins Publishing.

Schmitt, N., & Underwood, G. (2004). Exploring the processing of formulaic sequences through a self-paced reading task. *Formulaic sequences: Acquisition, processing and use*, 173-189.

Schmidt, J. R., & Weissman, D. H. (2016). Congruency sequence effects and previous response times: conflict adaptation or temporal learning?. *Psychological Research*, *80*(4), 590-607.

Serratrice, L. (2013). Cross-linguistic influence in bilingual development: Determinants and mechanisms. *Linguistic Approaches to Bilingualism*, *3*(1), 3-25.

Serratrice, L. (2018). Becoming bilingual in early childhood. In A. De Houwer & L. Ortega (Eds.), *The Cambridge Handbook of Bilingualism* (pp. 15-35). Cambridge, United Kingdom: Cambridge University Press.

Shin, D., & Nation, P. (2008). Beyond single words: The most frequent collocations in spoken English. *ELT journal*, *62*(4), 339-348.

Sigott, G. (1995). The C-test: some factors of difficulty. *AAA: Arbeiten aus Anglistik und Amerikanistik*, 43-53.

Singleton, D. (2003). Critical period or general age. *Age and the acquisition of English as a foreign language. Clevedon*, 3-22.

Siyanova, A., & Schmitt, N. (2008). L2 learner production and processing of collocation: A multi-study perspective. *Canadian Modern Language Review*, *64*(3), 429-458.

Siyanova-Chanturia, A., Conklin, K., & Schmitt, N. (2011). Adding more fuel to the fire: An eye-tracking study of idiom processing by native and non-native speakers. *L2 Research*, *27*(2), 251-272.

Smith, S. A., & Murphy, V. A. (2015). Measuring productive elements of multi-word phrase vocabulary knowledge among children with English as an additional or only language. *Reading and Writing*, *28*(3), 347-369.

Sonbul, S., & Siyanova-Chanturia, A. (2019). Research on the on-line processing of collocation: Replication of Wolter and Gyllstad (2011) and Millar (2011). *Language Teaching,* 1-9. doi:10.1017/s0261444819000132

Staub, A., & Rayner, K. (2007). Eye movements and on-line comprehension processes. *The Oxford handbook of psycholinguistics*, *327*, 342.

Strand, S., & Demie, F. (2005). English language acquisition and educational attainment at the end of primary school. *Educational Studies*, *31*(3), 275-291.

Strand, S., & Hessel, A. (2018). *English as an Additional Language, proficiency in English and pupils' educational achievement: An analysis of Local Authority data.* (Rep.). Retrieved https://www.bell-foundation.org.uk/app/uploads/2018/10/EAL-PIE-and-Educational-Achievement-Report-2018-FV.pdf

Suzuki, Y., & DeKeyser, R. (2017). The interface of explicit and implicit knowledge in a L2: Insights from individual differences in cognitive aptitudes. *Language Learning*, *67*(4), 747-790.

Suzuki, Y., & DeKeyser, R. (2017). The interface of explicit and implicit knowledge in a L2: Insights from individual differences in cognitive aptitudes. *Language Learning*, *67*(4), 747-790.

Tarar, J. M., Meisinger, E. B., & Dickens, R. H. (2015). Test Review: Test of Word Reading Efficiency–Second Edition (TOWRE-2) by Torgesen, JK, Wagner, RK, & Rashotte, CA.

Thierry, G., & Wu, Y. J. (2007). Brain potentials reveal unconscious translation during foreign-language comprehension. *Proceedings of the National Academy of Sciences*, *104*(30), 12530-12535.

Thordardottir, E., Grüter, T., & Paradis, J. (2014). The typical development of simultaneous bilinguals. *Input and experience in bilingual development*, *13*, 141.

Thordardottir, E., Rothenberg, A., Rivard, M. E., & Naves, R. (2006). Bilingual assessment: Can overall proficiency be estimated from separate measurement of two languages?. *Journal of multilingual communication disorders*, *4*(1), 1-21.

Titone, D., Libben, M., Mercier, J., Whitford, V., & Pivneva, I. (2011). Bilingual lexical access during L1 sentence reading: The effects of L2 knowledge, semantic constraint, and L1–L2 intermixing. Journal of Experimental Psychology: Learning, Memory, and Cognition, 37(6), 1412–1431. doi:10.1037/a0024492

Torgesen, J. K., Wagner, R., & Rashotte, C. (2012). *Test of Word Reading Efficiency:(TOWRE-2)*. Pearson Clinical Assessment.

Traxler, M. J., & Tooley, K. M. (2008). Priming in sentence comprehension: Strategic or syntactic? *Language and Cognitive Processes*, *23*(5), 609-645.

Traxler, M. J., & Tooley, K. M. (2008). Priming in sentence comprehension: Strategic or syntactic? *Language and Cognitive Processes*, *23*(5), 609-645.

Treffers-Daller, J. Language dominance: The construct, its measurement, and operationalization. *Language Dominance in Bilinguals*, 235-265. doi:10.1017/cbo9781107375345.012

Treffers-Daller, J., & Silva-Corvalán, C. (Eds.). (2016). *Language dominance in bilinguals: Issues of measurement and operationalization*. Cambridge University Press.

Tremblay, A., & Baayen, R. H. (2010). Holistic processing of regular four-word sequences: A behavioral and ERP study of the effects of structure, frequency, and probability on immediate free recall. *Perspectives on formulaic language: Acquisition and communication*, *151*, 173.

Tuller, L. (2015). 11 clinical use of parental questionnaires in multilingual contexts. *Assessing multilingual children: Disentangling bilingualism from language impairment*, *13*, 301-330.

—

Underwood, G., & Everatt, J. (1992). The role of eye movements in reading: some limitations of the eye-mind assumption. In *Advances in psychology* (Vol. 88, pp. 111-169). North-Holland.

Underwood, G., Schmitt, N., & Galpin, A. (2004). The eyes have it. *Formulaic sequences: Acquisition, Processing, and Use*, *9*, 153.

Valdés, G. (1999). Heritage language students: Profiles and possibilities. *Heritage languages in America: Preserving a national resource*, 37-80.

Van Hell, J. G., & De Groot, A. M. (1998). Conceptual representation in bilingual memory: Effects of concreteness and cognate status in word association. *Bilingualism: Language and cognition*, *1*(3), 193-211.

Van Heuven, W. J. B., Dijkstra, T., & Grainger, J. (1998). Orthographic neighborhood effects in bilingual word recognition. Journal of Memory and Language, 39, 458±483.

Van Heuven, W. J., Schriefers, H., Dijkstra, T., & Hagoort, P. (2008). Language conflict in the bilingual brain. *Cerebral cortex*, *18*(11), 2706-2716.

Vihman, M. M. (1985). Language differentiation by the bilingual infant. *Journal of child language*, *12*(2), 297-324.

Vijayalakshmi, M., & Babu, M. S. (2014). A brief history of English language teaching in India. *International Journal of scientific and research Publications*, *4*(5), 1-4.

Volterra, V., & Taeschner, T. (1978). The acquisition and development of language by bilingual children. *Journal of child language*, *5*(2), 311-326.

Volterra, V., & Taeschner, T. (1978). The acquisition and development of language by bilingual children. *Journal of child language*, *5*(2), 311-326.

Warren, H. (Ed.). (1994). *Oxford learner's dictionary of English idioms*. Oxford University Press.

Webb, S., Newton, J., & Chang, A. (2013). Incidental learning of collocation. *Language Learning*, *63*(1), 91-120.

Wei, L. (Ed.). (2020). *The Bilingualism Reader*. Routledge.

Weinreich, U. (2010). *Languages in contact: Findings and problems* (No. 1). Walter de Gruyter.

Weissman, D. H., Jiang, J., & Egner, T. (2014). Determinants of congruency sequence effects without learning and memory confounds. *Journal of Experimental Psychology: Human Perception and Performance*, *40*(5), 2022.

Werker, J. F., & Byers-Heinlein, K. (2008). Bilingualism in infancy: First steps in perception and comprehension. *Trends in cognitive sciences*, *12*(4), 144-151.

White, L. (1989). *Universal grammar and L2 acquisition* (Vol. 1). John Benjamins Publishing.

White, L., & White, L. (2003). *L2 acquisition and universal grammar*. Cambridge University Press.

Williams, R., & Morris, R. (2004). Eye movements, word familiarity, and vocabulary acquisition. *European Journal of Cognitive Psychology*, *16*(1-2), 312-339.

Winter, J., & Pauwels, A. (2005). Gender in the construction and transmission of ethnolinguistic identities and language maintenance in immigrant Australia. *Australian Journal of Linguistics*, *25*(1), 153-168.

Wolter, B. (2001). Comparing the L1 and L2 mental lexicon: A depth of individual word knowledge model. *Studies in L2 acquisition*, 41-69.

Wolter, B. (2002). Assessing proficiency through word associations: is there still hope?. *System*, *30*(3), 315-329.

Wolter, B., & Gyllstad, H. (2011). Collocational links in the L2 mental lexicon and the influence of L1 intralexical knowledge. *Applied Linguistics*, *32*(4), 430-449.

Wolter, B., & Gyllstad, H. (2013). Frequency of input and L2 collocational processing: A comparison of congruent and incongruent collocations. *Studies in L2 Acquisition*, *35*(3), 451-482.

Wolter, B., & Yamashita, J. (2018). Word frequency, collocational frequency, L1 congruency, and proficiency in L2 collocational processing: What accounts for L2 performance?. *Studies in L2 Acquisition*, *40*(2), 395-416.

Wray, A. (2000). Formulaic sequences in L2 teaching: Principle and practice. *Applied Linguistics*, *21*(4), 463-489.

Wray, A. (2005). *Formulaic language and the lexicon*. Cambridge University Press.

Yamashita, J., & Jiang, N. A. N. (2010). L1 influence on the acquisition of L2 collocations: Japanese ESL users and EFL learners acquiring English collocations. *Tesol Quarterly*, *44*(4), 647-668.

Yip, V., & Matthews, S. (2000). Syntactic transfer in a Cantonese–English bilingual child. *Bilingualism: Language and cognition*, *3*(3), 193-208.

Yip, V., & Matthews, S. (2007). *The bilingual child*. Ernst Klett Sprachen.

Zevin, J. D., & Seidenberg, M. S. (2002). Age of acquisition effects in word reading and other tasks. *Journal of Memory and Language*, *47*(1), 1-29.

Zobl, H. (1980). The formal and developmental selectivity of LI influence on L2 acquisition. *Language Learning*, *30*(1), 43-57.

# Appendices: Part I

## Appendix 1: Information Sheets and Consent Forms

*Appendix 1 (a): Study 1 Information Sheet (Tamil)*

**University of Reading**

School of Literature and Languages

Department of Applied Linguistics

Department of English Language and Applied Linguistics

HumSS Building
The University of Reading
Whiteknights, PO Box 218
Reading RG6 6AA

*Phone*


*Email appling@reading.ac.uk*
          p.a.thompson@reading.ac.uk

**Researcher**:
Roopa Leonard
*Email:* r.k.leonard@pgr.reading.ac.uk

**Supervisor**:
Dr Michael Daller
*Phone*:
*Email*: m.daller@reading.ac.uk

## தகவல தாள

*ஆங்கில மொழி சொற்களஞ்சியம் எவ்வாறு பெறுகிறது என்பதை மாணவர்களின் சொந்த மொழி (தமிழ்) எவ்வாறு பாதிக்கிறது என்பதை இந்தத் திட்டம் ஆராய்கிறது. தமிழ் பேசும் இளம் மாணவர்களிடையே இந்த ஆய்வு கவனம் செலுத்துவதால் இந்த பங்கேற்பாளர்கள் தேர்வு செய்யப்பட்டுள்ளனர். மாணவர்கள் படிப்பிற்கான இரண்டு பகுதிகளிலும் பங்கு பெற வேண்டுமென கோரியுள்ளனர்: ஒரு சொல்லகராதி சோதனை மற்றும் ஒரு சுய-வாசிப்பு வாசிப்பு பரிசோதனை. இந்த சோதனைகளின் மதிப்பெண்கள் பாதுகாப்பான கணினியில் கடவுச்சொல் பாதுகாக்கப்பட்ட இயக்ககத்தில் சேகரிக்கப்பட்டு சேமிக்கப்படும். ஆராய்ச்சியாளர் (நானே) மற்றும் என் மேற்பார்வையாளர்கள் மட்டுமே தரவு அணுக முடியும். தரவு பாதுகாப்பு சட்டத்தின் விதிமுறைகளால் வரையறுக்கப்பட்டு, கல்விக் கல்வியின் நோக்கத்திற்காக மட்டுமே பயன்படுத்தப்படும். பங்கேற்பாளர்களின் அடையாளங்கள் திட்டத்தின் எந்த நேரத்திலும் வெளிப்படுத்தப்படாது மற்றும் கடுமையான இரகசியத்தன்மை பராமரிக்கப்படும். அவர்கள் விரும்பினால், பங்கேற்பாளர்கள் ஆய்வின் எந்த கட்டத்திலும் திரும்ப பெற சுதந்திரமாக உள்ளனர்.*

*இந்த திட்டம் பள்ளி ஒழுக்கவியல் குழு மூலம் நெறிமுறை மறு ஆய்வுக்கு உட்பட்டது, மற்றும் ஆராய்ச்சிக் நெறிமுறைகளின் வழிகாட்டலுக்கான பல்கலைக்கழகங்களுக்கான குறிப்புகள் பத்தி 6-ல் குறிப்பிட்டுள்ள விதிவிலக்கு வழிமுறைகளின் கீழ் தொடர அனுமதிக்கப்பட்டுள்ளது.*


*நீங்கள் எந்த வினாக்களும் இருந்தால் அல்லது படிப்பிற்கான எதையும் தெளிவுபடுத்த விரும்பினால், தயவுசெய்து என் மேற்பார்வையாளரை தொடர்பு கொள்ளவும். மேலே உள்ள முகவரி அல்லது மின்னஞ்சல் மூலம்* m.daller@reading.ac.uk.


*கையெழுத்திட்ட*

*Appendix 1 (b): Study 1 Information Sheet (English)*

**University of Reading**

Department of English Language and Applied Linguistics

HumSS Building
The University of Reading
Whiteknights, PO Box 218
Reading RG6 6AA

School of Literature and Languages

Department of Applied Linguistics

*Phone*

**Researcher**:
Roopa Leonard
*Email:* r.k.leonard@pgr.reading.ac.uk

*Email appling@reading.ac.uk*
p.a.thompson@reading.ac.uk

**Supervisor**:
Dr Michael Daller
*Phone*:
*Email*: m.daller@reading.ac.uk

## Information Sheet

This project explores how the native language of the students (Tamil) influences how they acquire English vocabulary. These participants have been selected because the study focuses on young learners who are Tamil speakers. The students are requested to take part in the two parts of the study: a vocabulary test and a self-paced reading experiment. The scores of these tests will be collected and stored in a password-protected drive on a secure computer. Only the researcher (myself) and my supervisors will have access to the data. The data will only be used for the purposes of academic study, restricted by terms of the Data Protection Act. The identities of the participants will not be revealed at any point of the project and strict confidentiality will be maintained. If they wish, the participants are at liberty to withdraw at any stage of the study.

This project has been subject to ethical review by the School Ethics Committee, and has been allowed to proceed under the exceptions procedure as outlined in paragraph 6 of the University's *Notes for Guidance* on research ethics.

If you have any queries or wish to clarify anything about the study, please feel free to contact my supervisor at the address above or by email at m.daller@reading.ac.uk.

Signed

*Appendix 1 (c): Study 1 Consent Form (Tamil)*

## University of READING

**School of Literature and Languages**
**Department of English Language and Applied Linguistics**

நெறிமுறை குழு

ஒப்புதல் படிவம்

திட்டத்தின் தலைப்பு: இளம் கற்கைகளுக்கான L2 இடமாற்றத்தின் மீதான L1 இன் செல்வாக்கு

இந்த ஆராய்ச்சியின் நோக்கம் எனக்கு புரிகிறது, எனக்கு என்ன தேவை என்பதை புரிந்துகொள்கிறேன்; ரூப லியோனாரால் எனக்கு விளக்கப்பட்டுள்ளது இந்த திட்டம் தொடர்பான தகவல் தாள், படித்து புரிந்து. என் பங்கிற்கு தொடர்புபடுத்தியதில் இதுவரை தகவல் தாள் விவரித்துள்ள ஏற்பாடுகளை நான் ஒப்புக்கொள்கிறேன்.

என் பிள்ளையின் பங்களிப்பு முற்றிலும் தானாகவே உள்ளது மற்றும் அவர் எந்த நேரத்திலும் திட்டத்திலிருந்து விலக்குவதற்கான உரிமையுடையவர் என்று எனக்கு புரிகிறது.

இந்த ஒப்புதலுக்கான படிவம் மற்றும் அதனுடன் தொடர்புடைய தகவல் தாள் ஆகியவற்றை நான் பெற்றுள்ளேன்.

பெயர்:

ஒப்பந்தம்:

நாள்:

*Appendix 1 (d): Study 1 Consent Form (English)*

**School of Literature and Languages**
**Department of English Language and Applied Linguistics**

University of
**Reading**

ETHICS COMMITTEE

*Consent Form*

Project title**: The Influence of the L1 on the L2 Collocation Acquisition of Young ESL**

**Learners**

I understand the purpose of this research and understand what is required of me; I have read and understood the Information Sheet relating to this project, which has been explained to me by Roopa Leonard. I agree to the arrangements described in the Information Sheet in so far as they relate to my participation.

I understand that my child's participation is entirely voluntary and that he/she has the right to withdraw from the project at any time.

I have received a copy of this Consent Form and of the accompanying Information Sheet.

Name:

Signed:

Date:

**Appendix 2: X-lex test (English)**

**X-Lex Vocabulary Size Test**

Please look at these words. Some of these words are real English words and some are invented but are made to look like real words. Please tick the words that you know or can use. Here is an example.

**dog ✓**

Thank you for your help.

| | | | | | |
|---|---|---|---|---|---|
| that | both | sumption | sandy | lessen | independent |
| with | century | stream | fishlock | oak | woman |
| cantileen | cup | normal | impress | antique | kennard |
| person | discuss | everywhere | staircase | horobin | dish |
| feel | gillen | deny | daily | limp | military |
| round | path | shot | essential | permission | before |
| early | tower | refer | hyslop | headlong | park |
| table | weather | darrock | conduct | violent | humble |
| question | wheel | feeling | relative | frequid | fade |
| effect | alden | bullet | upward | rake | sorrow |
| market | perform | juice | publish | trunk | provide |
| waygood | pity | nod | insult | mercy | arrive |
| stand | probable | gentle | cardboard | anxious | gumis |
| believe | signal | slip | pardoe | pedestrian | whole |
| fine | gazard | diamond | contract | arrow | treadaway |
| instead | earn | press | mount | feeble | cliff |
| produce | sweat | candlin | tube | hobrow | horozone |
| group | trick | drum | moreover | brighten | associate |
| litholect | manage | reasonable | crisis | dam | manomize |
| difficult | mud | boil | Jug | outlet | antique |

*Appendix 2 (a): X-lex English Samples (Study 1)*

# Sample 1



**Score: 4200/5000**

**Sample 2**



**Score: 3150/5000**

## Sample 3



CW: 69   IW: 4

2450

### X-Lex Vocabulary Size Test

Please look at these words. Some of these words are real English words and some are invented but are made to look like real words. Please tick the words that you know or can use. Here is an example.

**dog**

Thank you for your help.

| | | | | | |
|---|---|---|---|---|---|
| that ✓ | both ✓ | sumption ✗ | sandy ✓ | lesser ✓ | independent ✓ |
| with ✓ | century | stream ✓ | fishlock ✗ | oak ✓ | woman |
| cantileen | cup ✓ | normal | impress ✓ | antique ✓ | kennard |
| person | discuss ✓ | everywheré | staircase | horobin ✓ | dish ✓ |
| feel ✓ | gillen | deny ✓ | daily ✓ | limp ✓ | military |
| round ✓ | path ✓ | shot ✓ | essential | permission | before ✓ |
| early | tower | eter ✓ | hyslop ✗ | headlong | park ✓ |
| table ✓ | weather ✓ | darrook ✓ | conduct | violent ✓ | humble |
| question | wheel ✓ | feeling ✓ | relative | frequid | fade ✓ |
| effect | alden ✓ | bullet ✓ | upward ✓ | rake ✓ | sorrow ✓ |
| market ✓ | perform ✓ | juice ✓ | publish | trunk ✓ | provide ✓ |
| waygood ✗ | pity ✓ | nod ✓ | insult ✓ | mercy ✓ | arrive ✓ |
| stand ✓ | probable | gentle | cardboard | anxious | gumis |
| believe ✓ | signal | slip ✓ | pardoe ✗ | pedestrian | whole ✓ |
| fine ✓ | gazard | diamond | contract ✓ | arrow ✓ | treadaway |
| instead ✓ | earn ✓ | press ✓ | mount ✓ | feeble ✓ | cliff ✓ |
| produce | swear ✓ | candlin | tube ✓ | hobrow | horozone ✓ |
| group | trick ✓ | drum ✓ | moreover ✓ | brighten | associate ✓ |
| litholect | manage | reasonable | crisis | dam ✓ | manomize |
| difficult ✓ | mud ✓ | boil ✓ | jug ✓ | outlet ✓ | antique ✓ |

**Score: 2450/5000**

**Appendix 3: X-lex Tamil test**

| | | | | | |
|---|---|---|---|---|---|
| நான் | முடியாத | நாலயம் | அடிக்கடி | இறுதி | மாநாடு |
| மக்கள் | தடப்பன் | பணம் | படாசம் | மாலை | கிடைத்த |
| சூழல் | தெரியும் | கூடிய | கலை | பெறுந்து | வெளி |
| நாடகால | உலகம் | இல்லாத | சாதிதிரும் | வாங்கி | மார் |
| அரசியல் | பகுதி | முன் | அறிவு | அஷோ | கண் |
| மேலும் | கோடி | எழுதி | காலை | தர | கால |
| நாள் | ஆலவரம் | நோய் | மாதி | தலை | பின்பு |
| மிகவும் | கண்டு | தாடம் | பால | பன்னி | கத்தறி |
| ஆண்டு | மணம் | திரும் | கெருப்பு | பட்ட | பேசி |
| பெயர் | சரிப்பு | சாற்று | தாய் | பொங்கன் | பிரிவு |
| பத்தொன்ம் | சிறிய | வீடு | மகன் | மருந்து | சீனா |
| நேரம் | நீர் | பற்றியும் | நம்மை | பிடித்து | ரொபாச்கில் |
| காலம் | இங்கே | படி | கஜிதை | பயண | கட்டி |
| வெற்றிகி | வெற்றி | சாங்காடு | கை | வரி | பலரும் |
| போய்ச் | பெற்ற | ஏற்பட்ட | அம்மா | அரசு | வாழ்விலில் |
| இடிவு | அளவு | உயர் | பல்லிடம் | கடல் | பிரிவு |
| நாடு | உடல் | இரல் | தாம் | சிறியண் | கயிரு |
| கல்வி | பார்வை | காதல் | கொலை | நடிகர் | தலை |
| நன்றி | இடிவு | பயணம் | உணர்வு | ரொம்ப | நாயகம் |
| சிங்கள | மொழி | மாநில | நன்றாக | வாநிலை | தோடர் |

*Appendix 3 (a): X-lex Tamil Samples*

**Sample 1**



**Score: 1050/5000**

**Sample 2**



நான் முடியாத தாலயம் அடிங்கடி இறுதி மாநாடு
மக்கள் நடப்பன் பணம் லட்சம் மாலை கிளைத்த
மீளம் தெரியும் கூடிய கலை பேருந்து வெளி
நாட்கால உலகம் இல்லாத சாத்திரம் வாங்கி மனர்
அரசியல் பதிகி மீன் அறிவு கீழா கண
மேலும் கோடி எழுதி காலை நூறு கோல்
தான் ஆலகரம் நோய் மாழி கண்ல பின்பு
மிகவும் கண்கு தாலயம் பால் பன்னி கத்தரி
ஆண்கு மனம் திருவ கெருப்பு பட்ட பேனி
பெயர் சாய்ப்பு சற்று தால் பொங்கன் பிரிவு
பத்தஎணம் சிறிய வீடு மகன் ழிருந்து சீனா
நேரம் நீர் பற்றியும் நம்மை பிடித்து நீர் பாய்ச
காயம் இங்கே பழ் கவிதை பயன் கூட்டி
வெற்றி உண்மை சாங்காடு நக வரி பழரும்
போய்ச் பெற்ற ஏற்பாட்ட அம்மா அரம வாழ்வில்
இடிவு அன்பு உயர் பலயிடம் கடல் பிரிவு
நாகு உடல் குரல் தாம் சூரியன் குலிகு
கல்லி பாணச காதல் கொலை நடிகர் தலை
நண்றி முடிவு பயணம் உணர்வு தெரம்பு நாமகு
திங்கள் மொழி மாநில நன்றாக வாழிலை தோடர்

CW:83  IW:4  3150

**Score: 3150/5000**

# Sample 3

| நான் ✓ | முடியாத ✓ | நாலயம் | அடிக்கடி ✓ | இறுதி ✓ | மாநாடு ✓ |
| மக்கள் | தடப்பன் | பணம் ✓ | வட ஈழ | மாலை | கிடைத்த |
| மீசம் | தெரியும் | கூடிய ✓ | கலை | பெறுந்து | வெளி |
| நாட்கால் ✓ | உலகம் | இல்லாத | சாத்திரம் | வாங்கி | மன் |
| அரசியல் ✓ | பதிகுதி | மின் | அறிவு | கிழா | கண் |
| மேலும் ✓ | கோடி ✓ | எழுதி | நாலை | நர | கால் |
| தான் ✓ | ஆலகரம் | நோய் ✓ | மாறி | தலை | பின்பு |
| மிகவும் | கண் ✓ | நாலயம் ✓ | பால் ✓ | பன்னி | கதறறி |
| ஆண்டு ✓ | மணம் ✓ | திரும் ✓ | கெருப்பு | பட்ட | பேசி |
| பெயர் | சாரிப்பு ✓ | சற்று ✓ | தாய் ✓ | பொங்கன் | மிறிவு |
| பத்தனம் | சிறிய ✓ | வீடு ✓ | மகன் | பிரிந்து ✓ | சீனா ✓ |
| நேரம் ✓ | நீர் ✓ | பற்றியும் ✓ | நம்மை | அடித்து ✓ | நீர் பாங்க |
| காலம் ✓ | இங்கே | படி | கவிதை | பயன் ✓ | கட்டி ✓ |
| வெற்றி ✓ | உண்மை | சாங்காடு ✓ | கை | வரி | மண்ரு ✓ |
| போய்ச் | பெற்ற | ஏற்பட்ட | அம்மா | அரசு | வாழ்வி |
| முடிஉ ✓ | அனுவு | உயர் ✓ | பல்லிடம் | கடல் | அறிவு |
| நாடு ✓ | உடல் ✓ | குரல் ✓ | தாம் | சூரியன் | கயிரு |
| கல்வி ✓ | பார்வை | காதல் ✓ | கொலை ✓ | நடிகள் | தலை |
| நன்றி ✓ | முடிவு ✓ | பயணம் | உணர்வு | ரொம்ப | நாயக |
| திங்கள | மொழி | மாநில ✓ | நன்றாக ✓ | வாநிலை | தொடர |

CW:99  IW:10    2450

**Score: 2450/5000**

**Appendix 4: C-test**

Please look at the example and fill in the blanks for the other three paragraphs.

**Example**

Plants such as trees and grass are at the bottom of the food chain. Plants g_____ th____ energy fr____ the s____. Animals su____ as de____ and rab_____ get th____ energy b_ eating pla_____. They a____ called herbi_____, which me_____ 'plant eat_____'. There a____ many mo____ herbivores o_ our pla_____ than carni_____, which a____ animals th____eat me____. Predators such as wolves and lions are at the top of the food chain.

**At the beach**

One of the best ways to escape the heat of the summer months is to head down to the beach. To co____ off, ma____ people g_ swimming i_ the oce____. When th____ are do__ in th_ water, th____ use a tow____ to dr_ off an_ then bui__ a sandc_____ in th_ sand. At t____ beach, it's f____ to sea_____ for thi_____ that wa____ up o_ the sh_____. There are also many curly shells, which children like to collect.

**Back to school**

To get ready for the new school year, students will need to buy some school supplies. They wi__ need pe___ to wri___ with, noteb_____ to wri__ in, cray___ to col____ with, scis____ to cu_ with, eras____ to era___ with, and gu_ to pas__ things toge____. Finally, they'll need a bag to carry everything in.

**Evaporation**

On a warm, sunny day, water in a glass of water seems to slowly disappear. The wat___ disappears beca___ the ene____ from th_ sun i_ heat____ the wat___ up an_ turn___ the liq____ water in__ water vap____. This proc___ is cal___ evaporation. When the water evaporates, it becomes an invisible gas in the atmosphere.

*Appendix 4 (a): C-test Samples*

# Sample 1

Please look at the example and fill jn the blanks for the other three paragraphs.

**Example**

Plants such as trees and grass are at the bottom of the food chain. Plants get their energy from the sun. Animals such as deer and rabbits get their energy by eating plants. They are called herbivores, which means 'plant eaters'. There are many more herbivores on our planet than carnivores, which are animals that eat meat. Predators such as wolves and lions are at the top of the food chain.

**At the beach**

One of the best ways to escape the heat of the summer months is to head down to the beach. To cool off, many people go swimming in the ocean. When they are done in the water, they use a towel to dry off and then build a sandcastle in the sand. At the beach, it's fun to search for things that wash up on the shore. There are also many curly shells, which children like to collect.

**Back to school**

To get ready for the new school year, students will need to buy some school supplies. They will need pencil to write with, notebook to write in, crayons to color with, scissors to cut with, erasers to erase with, and gum to paste things together. Finally, they'll need a bag to carry everything in.

14

**Evaporation**

On a warm, sunny day, water in a glass of water seems to slowly disappear. The water disappears because the energy from the sun it heat up the water up and turn to the liquid water in the water vapour. This process is called evaporation. When the water evaporates, it becomes an invisible gas in the atmosphere.

10

**Score: 40/5**

# Sample 2

Please look at the example and fill in the blanks for the other three paragraphs.

**Example**

Plants such as trees and grass are at the bottom of the food chain. Plants _get_ th_eir_ energy fr_om_ the s_un_. Animals su_ch_ as de_er_ and rabb_its_ get th_eir_ energy b_y_ eating pla_nts_. They a_re_ called herbi_vores_, which me_ans_ 'plant eat_ers_'. There a_re_ many mo_re_ herbivores o_n_ our pla_net_ than carnivor_es_, which a_re_ animals tha_t_ eat me_at_. Predators such as wolves and lions are at the top of the food chain.

**At the beach**

One of the best ways to escape the heat of the summer months is to head down to the beach. To co_ol_ off, ma_ny_ people g_o_ swimming i_n_ the oce_an_. When th_ey_ are do_ne_ in th_e_ water, th_ey_ use a tow_el_ to dr_y_ off and then bui_ld_ a sandc_astle_ in th_e_ sand. At th_e_ beach, it's f_un_ to sea_rch_ for thi_ngs_ that wa_sh_ up o_n_ the sh_ore_. There are also many curly shells, which children like to collect.

**Back to school**

To get ready for the new school year, students will need to buy some school supplies. They will need pe pay to write with, notebook to write in, cray  to collect with, scisptist to cut with, eraserssto era  with, and guf to past things together. Finally, they'll need a bag to carry everything in.

**Evaporation**

On a warm, sunny day, water in a glass of water seems to slowly disappear. The water disappears because the energt from the sun in heatsun the water up and turn into the liquid water into water vapour . This process is called evaporation. When the water evaporates, it becomes an invisible gas in the atmosphere.

**Score 31/50**

# Sample 3

Please look at the example and fill jn the blanks for the other three

paragraphs.

**Example**

(25)

Plants such as trees and grass are at the bottom of the food chain.

Plants g et  th eir  energy fr om  the s un . Animals su ch  as de er

and rabb its  get their  energy by eating pla nts  . They a re  called

herbi vores  , which me ans  'plant eat ers '. There a re  many mo re

herbivores on our pla net  than carniv ores , which a re  animals

tha t eat me at . Predators such as wolves and lions are at the top of

the food chain.

**At the beach**

One of the best ways to escape the heat of the summer months is to

head down to the beach. To c__ off, ma__ people __ swimming t__

the oce__. When th__ are do__ in th__ water, th__ use a tow__ to

dr__ off and then bu__ a sandc__ in th__ sand. At th__ beach, it's

f__ to sea __ for thi__ that wa__ up __ he sh__. There are

also many curly shells, which children like to collect.

**Back to school**

To get ready for the new school year, students will need to buy some school supplies. They will need pen to write with, notebook to write in, crayon to color with, scis to cut with, eras and to erase with, and gu to paste things together. Finally, they'll need a bag to carry everything in.

**Evaporation**

On a warm, sunny day, water in a glass of water seems to slowly disappear. The water disappears because the enemy from the sun is heat and the water up and turn out the liquid water into water vapour. This process is called evaporation. When the water evaporates, it becomes an invisible gas in the atmosphere.

**Score 25/50**

**Appendix 5: List of Collocations**

| Collocation Number | Collocation | Congruency | MI Score |
|---|---|---|---|
| S1.1 | heavy traffic | Incongruent | 7.21 |
| S1.2 | woken up | Incongruent | 7.03 |
| S1.3 | get ready | Incongruent | 4.54 |
| S1.4 | busy road | Congruent | 5.48 |
| S1.5 | paying attention | Incongruent | 6.29 |
| S2.1 | gives advice | Congruent | 5.39 |
| S2.2 | waste time | Congruent | 3.89 |
| S2.3 | hard work | Congruent | 7.77 |
| S2.4 | save money | Congruent | 11.38 |
| S2.5 | play games | Congruent | 6.36 |
| S3.1 | takes notes | Congruent | 3.58 |
| S3.2 | correct answer | Congruent | 6.27 |
| S3.3 | hard work | Congruent | 7.77 |
| S3.4 | made friends | Incongruent | 3.37 |
| S3.5 | have lunch | Incongruent | 5.85 |
| S4.1 | rainy season | Congruent | 11.13 |
| S4.2 | first time | Congruent | 5.45 |
| S4.3 | strong wind | Congruent | 6.17 |
| S4.4 | broad daylight | Incongruent | 10.48 |
| S4.5 | fly away | Incongruent | 5 |
| S5.1 | late evening | Incongruent | 4.62 |
| S5.2 | long way | Congruent | 5.45 |
| S5.3 | got lost | Incongruent | 3.65 |
| S5.4 | straight ahead | Incongruent | 7.81 |
| S6.1 | book tickets | Incongruent | 6.42 |
| S6.2 | short trip | Incongruent | 5.12 |
| S6.3 | bad luck | Congruent | 9.16 |
| S6.4 | sudden change | Congruent | 6.47 |
| S6.5 | next time | Congruent | 4.13 |
| S7.1 | very well | Congruent | 6.48 |
| S7.2 | best friends | Incongruent | 6.11 |
| S7.3 | quite often | Incongruent | 5.09 |
| S7.4 | get along | Incongruent | 6.78 |
| S8.1 | work together | Congruent | 7.233 |
| S8.2 | bright idea | Incongruent | 5.28 |
| S8.3 | first step | Incongruent | 6.39 |
| S8.4 | made progress | Incongruent | 4.25 |

| S9.1 | waste time | Congruent | 3.89 |
|------|-----------|-----------|------|
| S9.2 | keep quiet | Incongruent | 6.47 |
| S9.3 | caused trouble | Incongruent | 6.14 |
| S9.4 | read aloud | Incongruent | 9.87 |
| S10.1 | good news | Congruent | 6.67 |
| S10.2 | warm welcome | Incongruent | 8.63 |
| S10.3 | serve dinner | Incongruent | 4.99 |
| S11.1 | lose weight | Incongruent | 8.46 |
| S11.2 | take action | Congruent | 4.79 |
| S11.3 | keep fit | Congruent | 7.83 |
| S11.4 | free time | Congruent | 4.51 |

**Appendix 6: Ministories**

**Single Mode**

1) It was Monday morning so there was **heavy traffic** in the city. Sammy was on his way to school on the school bus. He was sleepy because he had woken up early to **get ready** for school. Suddenly, a **young child** ran out on to the **busy road** right in front of the bus. Thankfully, the bus driver had been **paying attention** so he managed to stop the bus before it hit the child. (5)

2) My grandmother is very wise and always **gives advice** to her grandchildren. She tells us that we must not **waste time** and that **hard work** is very important. She advises us to **save money** even though we are still young. She says we must enjoy life and encourages us to **play games** every day. I love my grandmother and I try to follow her advice. (5)

3) Mary is a good student. She **takes notes** in class and always knows the **right answer** to the teacher's questions. The teachers praise her for her **hard work** and sincere attitude. K was a new student who needed help. Mary **made friends** with her and they began to **have lunch** together every day. (5)

4) It was the **rainy season** and Timmy couldn't go out to play. It was the **first time** this year that he had seen such a **strong wind** blowing outside. In **broad daylight**, he saw all the birds **fly away** to take shelter in the trees. (5)

5) The sun had set and it was **late evening** by the time David finished his work. He had missed the last bus. He had to walk a **long way** home and he **got lost** because he didn't know the way back. David was scared, but then he saw his favourite hotel **straight ahead** and he knew where he was. Je was relieved and he reached home safely. (4)

**Chunk Mode**

6) Anita and Naveen wanted to go on a holiday. They decided to **book tickets** to go on a **short trip** the following month. Unfortunately, they had a stroke of **bad luck** and there was a **sudden change** in their plans. They had to cancel the tickets. They hoped that **next time**, their plan would work out. (5)

7) Monica and Payal know each other **very well** because they have been **best friends** since nursery school. Monica likes Payal because Payal always encourages her. Payal likes Monica

because she **quite often** helps her. They **get along** with each other and spend a lot of time with each other. (4)

8) Susan and James had to **work together** for their Science class project. As they were thinking of what they could do, Susan had the **bright idea** of working on the different plants on the school campus. For the **first step**, they would have to go around campus and take pictures of the different kinds of plants. Susan and James started working on their project. Soon they had **made progress** and their teacher was very happy with them. (4)

9) The teacher did not want to **waste time** because it was almost time for exams. She told the class to **keep quiet** so she could give them instructions. She warned the class that anyone who **caused trouble** would be sent to the Principal's office. She asked the class leader to **read aloud** the important questions for the exam. (4)

10) My mother told us the **good news** that she had received a promotion and she wanted to celebrate. She invited her friends home for dinner. I helped my mother to give them a **warm welcome** and to **serve dinner** to them. Everyone was very happy and they enjoyed the meal. (3)

**Appendix 7: Mean reading times for congruent and incongruent, and single and chunked collocations for each participant (in milliseconds)**

| Participant Number | Mean for Reading Time in Single mode (Congruent) | Mean for Reading Time in Single mode (Incongruent) | Mean for Reading Time in Chunk mode (Congruent) | Mean for Reading Time in Chunk mode (Incongruent) |
|---|---|---|---|---|
| 1 | 1676.53 | 2285.91 | 1511.72 | 2165.55 |
| 2 | 1625.76 | 1936.96 | 1426.61 | 1661.84 |
| 3 | 1483.13 | 1747.88 | 1249.31 | 1894.32 |
| 4 | 3090.16 | 3475.86 | 2081.35 | 2660.59 |
| 5 | 1625.96 | 2354.12 | 1653.52 | 1603.57 |
| 6 | 3840.01 | 4429.07 | 2524.70 | 2486.82 |
| 7 | 2987.41 | 3053.72 | 1923.08 | 2396.16 |
| 8 | 1537.02 | 1897.20 | 1695.17 | 1826.54 |
| 9 | 2772.31 | 3476.01 | 2236.67 | 2748.53 |
| 10 | 2202.36 | 2514.11 | 1790.24 | 2073.54 |
| 11 | 2133.09 | 2141.04 | 1832.88 | 1703.52 |
| 12 | 1970.35 | 2236.92 | 1282.91 | 1821.45 |
| 13 | 2041.70 | 2050.52 | 1535.08 | 1430.90 |
| 14 | 1223.27 | 1513.78 | 1183.48 | 1779.12 |
| 15 | 2120.15 | 2230.38 | 1447.54 | 1626.64 |
| 16 | 1913.45 | 2616.78 | 1765.91 | 1733.82 |
| 17 | 1995.99 | 2868.89 | 1444.80 | 1377.75 |
| 18 | 2817.75 | 4178.49 | 3191.15 | 2779.32 |
| 19 | 1350.18 | 1381.25 | 1434.71 | 1473.15 |
| 20 | 1438.41 | 1743.95 | 1248.08 | 1423.50 |
| 21 | 2145.30 | 2317.98 | 933.86 | 1091.48 |
| 22 | 6130.22 | 9776.19 | 0.00 | 0.00 |
| 23 | 1685.40 | 1893.74 | 1187.39 | 1149.09 |
| 24 | 3034.71 | 3866.13 | 1858.77 | 2045.01 |
| 25 | 1440.67 | 1636.70 | 1255.44 | 1378.63 |
| 26 | 2536.04 | 2964.26 | 1896.82 | 2268.57 |
| 27 | 3466.84 | 3377.39 | 2033.86 | 2654.11 |
| 28 | 2183.23 | 2425.85 | 1377.76 | 1858.93 |
| 29 | 2649.39 | 2814.14 | 1568.91 | 1821.96 |
| 30 | 2230.72 | 2401.87 | 1449.03 | 1448.82 |
| 31 | 3894.78 | 3779.90 | 2620.90 | 2468.15 |
| 32 | 2176.17 | 2373.04 | 1373.08 | 1496.87 |
| 33 | 3255.79 | 3213.27 | 2119.45 | 2458.70 |
| 34 | 1514.44 | 1795.84 | 1258.99 | 1248.18 |

| | | | | |
|---|---|---|---|---|
| 35 | 1297.90 | 1981.99 | 979.75 | 1458.87 |
| 36 | 2678.79 | 2685.44 | 1622.02 | 2198.87 |
| 37 | 1972.99 | 2139.67 | 1254.37 | 1683.26 |
| 38 | 2745.97 | 2767.83 | 1919.90 | 2122.09 |
| 39 | 2138.36 | 2330.40 | 1124.08 | 1567.77 |
| 40 | 2539.53 | 2278.58 | 1342.00 | 1631.64 |
| 41 | 1889.03 | 1711.29 | 793.91 | 906.98 |
| 42 | 2007.73 | 2124.36 | 1201.15 | 1278.73 |
| 43 | 2569.80 | 2359.35 | 1431.44 | 1299.97 |
| 44 | 2553.43 | 2582.23 | 1573.54 | 2072.07 |
| 45 | 2246.64 | 2439.85 | 1273.73 | 1405.91 |
| 46 | 1682.05 | 1762.92 | 1277.68 | 1158.65 |
| 47 | 2032.21 | 1920.96 | 1250.19 | 1316.41 |
| 48 | 3472.62 | 4431.29 | 1097.04 | 1049.73 |
| 49 | 2134.92 | 1832.07 | 940.94 | 1195.74 |
| 50 | 1764.40 | 2132.90 | 1480.11 | 1640.83 |
| 51 | 2344.42 | 2408.48 | 1486.11 | 2094.24 |
| 52 | 2527.50 | 2325.38 | 1338.67 | 1534.59 |
| 53 | 1584.57 | 1433.62 | 874.84 | 917.15 |
| 54 | 2398.43 | 2643.15 | 1315.72 | 1647.59 |
| 55 | 2107.98 | 3805.79 | 2002.55 | 2140.98 |
| 56 | 3447.32 | 3885.12 | 1909.99 | 2006.63 |
| 57 | 2204.86 | 2148.13 | 1224.14 | 1345.46 |
| 58 | 2181.80 | 1814.71 | 1395.71 | 1569.27 |
| Mean | 2322.62 | 2633.01 | 1508.67 | 1712.04 |

**Appendix 8: Mean Reading Times per Word in the Single Condition for Each Participant (in milliseconds)**

| Participant No. | Word 1 (Congruent) | Word 2 (Congruent) | Word 1 (Incongruent) | Word 2 (Incongruent) |
|---|---|---|---|---|
| 1 | 844.63 | 831.90 | 966.91 | 1319.00 |
| 2 | 838.19 | 787.57 | 915.99 | 1020.97 |
| 3 | 726.93 | 756.20 | 801.31 | 946.57 |
| 4 | 1569.35 | 1520.81 | 1471.29 | 2004.57 |
| 5 | 744.03 | 881.93 | 1155.19 | 1198.93 |
| 6 | 1939.81 | 1900.20 | 1732.75 | 2696.32 |
| 7 | 1554.82 | 1432.59 | 1549.29 | 1504.42 |
| 8 | 822.61 | 714.42 | 766.75 | 1130.46 |
| 9 | 1587.17 | 1185.14 | 2004.45 | 1471.56 |
| 10 | 1228.92 | 973.43 | 1291.02 | 1223.10 |
| 11 | 1113.85 | 1019.24 | 1054.70 | 1086.34 |
| 12 | 972.86 | 997.49 | 1007.50 | 1229.42 |
| 13 | 1027.69 | 1014.01 | 929.50 | 1121.01 |
| 14 | 646.02 | 577.26 | 796.38 | 717.40 |
| 15 | 1106.71 | 1013.44 | 1125.98 | 1104.40 |
| 16 | 857.75 | 1055.70 | 1327.67 | 1289.11 |
| 17 | 968.58 | 1027.41 | 1368.22 | 1500.67 |
| 18 | 1319.42 | 1498.33 | 1901.86 | 2276.63 |
| 19 | 675.97 | 674.22 | 718.50 | 662.75 |
| 20 | 724.75 | 713.66 | 743.16 | 1000.79 |
| 21 | 1084.97 | 1060.34 | 1125.77 | 1192.20 |
| 22 | 2579.73 | 3550.50 | 5549.03 | 4227.17 |
| 23 | 887.28 | 798.12 | 932.78 | 960.96 |
| 24 | 1492.12 | 1542.59 | 2029.18 | 1836.95 |
| 25 | 740.96 | 699.71 | 813.67 | 823.03 |
| 26 | 1406.44 | 1129.60 | 1439.25 | 1525.00 |
| 27 | 1863.06 | 1603.78 | 1611.90 | 1765.49 |
| 28 | 1092.00 | 1091.24 | 1225.13 | 1200.72 |
| 29 | 1341.56 | 1307.83 | 1462.70 | 1351.44 |
| 30 | 1114.58 | 1116.14 | 1213.86 | 1188.02 |
| 31 | 1846.28 | 2048.50 | 1760.38 | 2019.52 |
| 32 | 1150.46 | 1025.71 | 1185.78 | 1187.26 |
| 33 | 1640.62 | 1615.17 | 1628.76 | 1584.51 |
| 34 | 763.93 | 750.51 | 857.67 | 938.17 |
| 35 | 684.26 | 613.64 | 1166.27 | 815.73 |
| 36 | 1367.05 | 1311.74 | 1384.74 | 1300.69 |
| 37 | 1025.51 | 947.48 | 1056.00 | 1083.67 |
| 38 | 1377.37 | 1368.60 | 1290.74 | 1477.10 |
| 39 | 985.38 | 1152.98 | 1225.60 | 1104.80 |

| 40 | 1504.96 | 1034.56 | 1212.66 | 1065.92 |
|---|---|---|---|---|
| 41 | 793.90 | 1095.13 | 915.16 | 796.13 |
| 42 | 1036.67 | 971.06 | 1078.34 | 1046.02 |
| 43 | 1321.66 | 1248.14 | 1160.97 | 1198.37 |
| 44 | 1356.64 | 1196.79 | 1336.62 | 1245.61 |
| 45 | 1150.73 | 1095.91 | 1202.08 | 1237.77 |
| 46 | 770.71 | 911.34 | 947.93 | 814.99 |
| 47 | 986.03 | 1046.18 | 956.94 | 964.02 |
| 48 | 1636.79 | 1835.84 | 2051.44 | 2379.85 |
| 49 | 1262.30 | 872.62 | 941.16 | 890.92 |
| 50 | 1019.03 | 745.37 | 1046.62 | 1086.28 |
| 51 | 1270.45 | 1073.97 | 1145.56 | 1262.91 |
| 52 | 1469.92 | 1057.58 | 1348.79 | 976.60 |
| 53 | 839.45 | 745.12 | 691.97 | 741.65 |
| 54 | 1351.20 | 1047.23 | 1196.48 | 1446.67 |
| 55 | 1088.08 | 1019.90 | 1980.49 | 1825.29 |
| 56 | 1739.12 | 1708.20 | 1672.03 | 2213.09 |
| 57 | 1078.81 | 1126.05 | 893.47 | 1254.65 |
| 58 | 891.90 | 1289.90 | 816.75 | 997.96 |
| Mean | 1177.28 | 1145.35 | 1296.26 | 1336.75 |

**Appendix 9: Vocabulary and Proficiency Test Scores for Each Participant (Study 1)**

| Participant No. | X-lex English | X-lex Tamil | C-test |
|---|---|---|---|
| 1 | 2850 | 2600 | 25 |
| 2 | 1450 | 0 | 29 |
| 3 | 3750 | 1950 | 39 |
| 4 | 1500 | 2700 | 15 |
| 5 | 3100 | 2000 | 27 |
| 6 | 2850 | 1950 | 21 |
| 7 | 1950 | 1800 | 18 |
| 8 | 1300 | 1000 | 34 |
| 9 | 2900 | 3250 | 28 |
| 10 | 2500 | 1800 | 40 |
| 11 | 3500 | 1900 | 36 |
| 12 | 2950 | 2450 | 36 |
| 13 | 2700 | 2700 | 30 |
| 14 | 2700 | 2350 | 28 |
| 15 | 3000 | 2350 | 30 |
| 16 | 3400 | 2700 | 30 |
| 17 | 2600 | 3150 | 39 |
| 18 | 2000 | 1850 | 39 |
| 19 | 3000 | 2650 | 33 |
| 20 | 2250 | 2700 | 38 |
| 21 | 4050 | 1950 | 39 |
| 22 | 2100 | 2150 | 32 |
| 23 | 550 | 500 | 29 |
| 24 | 3750 | 2050 | 43 |
| 25 | 3200 | 3100 | 31 |
| 26 | 4200 | 3350 | 42 |
| 27 | 3200 | 1600 | 36 |
| 28 | 1450 | 1200 | 27 |
| 29 | 3950 | 1800 | 37 |
| 30 | 3000 | 1900 | 24 |
| 31 | 3150 | 3350 | 31 |
| 32 | 1300 | 1550 | 30 |
| 33 | 1050 | 2100 | 19 |
| 34 | 2800 | 2950 | 20 |
| 35 | 2450 | 350 | 30 |
| 36 | 3600 | 2150 | 40 |
| 37 | 2450 | 2100 | 27 |
| 38 | 3050 | 2200 | 38 |
| 39 | 3700 | 3000 | 40 |
| 40 | 4000 | 3950 | 40 |
| 41 | 4000 | 2600 | 30 |

| 42 | 1800 | 700 | 28 |
| 43 | 2500 | 3400 | 39 |
| 44 | 1800 | 2450 | 41 |
| 45 | 4350 | 3350 | 40 |
| 46 | 3650 | 2450 | 43 |
| 47 | 2150 | 1450 | 40 |
| 48 | 3500 | 3550 | 37 |
| 49 | 900 | 2400 | 12 |
| 50 | 2200 | 1450 | 31 |
| 51 | 2650 | 1950 | 26 |
| 52 | 3300 | 3300 | 36 |
| 53 | 3300 | 1050 | 27 |
| 54 | 1500 | 1450 | 30 |
| 55 | 850 | 1000 | 22 |
| 56 | 4150 | 2100 | 41 |
| 57 | 3550 | 750 | 36 |
| 58 | 3700 | 2800 | 44 |
| **Mean** | **2742.93** | **2161.20** | **32.39** |

**University of Reading**

# Appendices Part II

## Appendix 10: Study 2 Information Sheet

**Parent/carer information sheet**

**Research Project: The Influence of the L1 on the L2 Collocation Acquisition of Young Tamil-English Bilingual Children**

**Researcher: Roopa Leonard (r.k.leonard@pgr.reading.ac.uk)**

**Supervisors: Dr Michael Daller (m.daller@reading.ac.uk) and Dr Holly Joseph (h.joseph@reading.ac.uk)**

We would like to invite your child to take part in a research study about the vocabulary skills of Tamil-English bilingual children who live in the UK.

**What is the study?**
The aim of our research is to explore how Tamil-English bilingual children acquire vocabulary, specifically investigating the influence of Tamil knowledge of English vocabulary skills. We explore these questions by monitoring children's eye movements as they read sentences containing specially designed word clusters. We hope that our research will help to shed light on how children use their knowledge of vocabulary in two languages as they read. We also hope that our findings will contribute to the literature on bilingual vocabulary acquisition and possibly be used to further understand the benefits of bilingualism. Since there is very limited work in vocabulary acquisition with Indian children, this study aims to address that gap.

**Why has my child been chosen to take part?**
Your child has been invited to take part in the project because he/she is between 8-11 years of age and studies Tamil at a Tamil school in the UK.

**Does my child have to take part?**
It is up to you and your child to decide whether or not to take part. Giving permission is entirely voluntary. If you agree to allow your child to take part, please sign the attached consent form. If you decide to take part, you and your child are still free to withdraw at any time and without giving a reason.

**What will happen if my child takes part?**
Your child will be seen by a researcher for a session that lasts approx. 30-40 minutes. All sessions will be conducted by the researcher (Roopa Leonard) who has full DBS clearance which means that she has no relevant convictions and is safe to work with children. First your child will be asked to complete some short vocabulary knowledge tasks in Tamil and English. Followed by this they will be asked to complete a reading fluency task as well as a short interview with the researcher about their language use. Finally, your child will be asked to read simple sentences in English on a computer monitor. Children's eye movements will be recorded during this task using an eye tracker. The eye tracker is a small camera which sits below the computer monitor and videos children's patterns of eye-movements while they read. During an eye tracking session, children would be asked to sit at a desk and lean against a head and chin rest. They would not be restrained in any way. A picture of a child seated at an eye tracker is shown below.

**What are the risks and benefits of taking part?**
The information we collect from your child will remain confidential and will only be seen by the Researcher and Supervisors involved in the project. Neither you, your child or the school will be identifiable in any published report resulting from the study. Taking part will in no way influence the grades your child receives at school. Information about individuals will not be shared with the school.

**What will happen to the data?**
Any data collected will be held in strict confidence and no real names will be used in this study or in any subsequent publications. The records of this study will be kept private. No identifiers linking you, your child or the school to the study will be included in any sort of report that might be published. We anonymise all recorded data before analysing the results. Children will be assigned a number and will be referred to by that number on tasks. Research records will be stored securely in a locked filing cabinet and on a password-protected computer and only the research team will have access to the records. The data will be destroyed securely once the findings of the study are written up, after five years. The results of the study will be presented at national and international conferences, and in written reports and articles.

**Who has reviewed the study?**
This project has been reviewed following the procedures of the University Research Ethics Committee and has been given a favourable ethical opinion for conduct. The University has the appropriate insurances in place. Full details are available on request.

**What happens if I/ my child change our mind?**
You/your child can change your mind at any time without any repercussions. During the research, your child can stop completing the activities at any time by telling the researcher or telling his/her teacher. If you change your mind after data collection has ended, contact Dr Daller or Dr Joseph with the contact details provided above.

**What happens if something goes wrong?**
In the unlikely case of concern or complaint, you can contact Dr Daller or Dr Joseph with the contact details provided above.

Thank you for taking the time to read this letter.

Yours sincerely,

Roopa Leonard

**Appendix 11: Study 2 Consent Form**

**School of Literature and Languages**
**Department of English Language and Applied Linguistics**

# Ethics Committee

**Consent Form**

Project title**: The Influence of the L1 on the L2 Collocation Acquisition of Young Tamil-**

**English Bilingual Children**

I understand the purpose of this research and understand what is required of my child and I; I have read and understood the Information Sheet relating to this project, which has been explained to me by Roopa Leonard. I agree to the arrangements described in the Information Sheet in so far as they relate to my own participation and my child's participation.

I understand that my participation and my child's participation is entirely voluntary and that he/she has the right to withdraw from the project at any time.

I have received a copy of this Consent Form and of the accompanying Information Sheet.

Name:

Signed:

Date:

# Appendix 12: X-lex English Samples
## Sample 1

| | | | | | |
|---|---|---|---|---|---|
| that | both | cliff | horobin | lessen | discuss |
| with | century | stream | military | horozone | bullet |
| before | cup | treadaway | impress | antique | kennard |
| person | darrock | everywhere | staircase | chart | market |
| feel | park | deny | gillen | limp | trick |
| pardoe | path | shot | essential | permission | humble |
| early | tower | refer | associate | headlong | daily |
| table | weather | independent | conduct | candlin | round |
| question | wheel | feeling | relative | fade | diamond |
| effect | whole | waygood | upward | rake | hobrow |
| gazard | perform | juice | publish | trunk | violent |
| woman | alden | nod | insult | mercy | sorrow |
| stand | probable | gentle | cardboard | anxious | fine |
| believe | signal | slip | cantileen | pedestrian | pity |
| gumm | dish | frequid | contract | arrow | normal |
| instead | earn | press | mount | feeble | tube |
| produce | sweat | provide | sumption | litholect | oak |
| group | fishlock | drum | moreover | brighten | hyslop |
| arrive | manage | reasonable | crisis | dam | difficult |
| manomize | mud | bolt | jug | outlet | sandy |

4200

**Score: 4200/5000**

**Sample 2**

| | | | | | |
|---|---|---|---|---|---|
| that | both | cliff | horobin | lessen | discuss |
| with | century | stream | military | horozone | bullet |
| before | cup | treadaway | impress | antique | kennard |
| person | darrock | everywhere | staircase | chart | market |
| feel | park | deny | gillen | limp | trick |
| pardoe | path | shot | essential | permission | humble |
| early | tower | refer | associate | headlong | daily |
| table | weather | independent | conduct | candlin | round |
| question | wheel | feeling | relative | fade | diamond |
| effect | whole | waygood | upward | rake | hobrow |
| gazard | perform | juice | publish | trunk | violent |
| woman | alden | nod | insult | mercy | sorrow |
| stand | probable | gentle | cardboard | anxious | fine |
| believe | signal | slip | cantileen | pedestrian | pity |
| gumm | dish | frequid | contract | arrow | normal |
| instead | earn | press | mount | feeble | tube |
| produce | sweat | provide | sumption | litholect | oak |
| group | fishlock | drum | moreover | brighten | hyslop |
| arrive | manage | reasonable | crisis | dam | difficult |
| manomize | mud | boil | jug | outlet | sandy |

4400

**Score: 4400/5000**

# Sample 3

| | | | | | |
|---|---|---|---|---|---|
| that | both | cliff | horobin | lessen | discuss |
| with | century | stream | military | horozone | bullet |
| before | cup | treadaway | impress | antique | kennard |
| person | darrock | everywhere | staircase | chart | market |
| feel | park | deny | gillen | limp | trick |
| pardoe | path | shot | essential | permission | humble |
| early | tower | refer | associate | headlong | daily |
| table | weather | independent | conduct | candlin | round |
| question | wheel | feeling | relative | fade | diamond |
| effect | whole | waygood | upward | rake | hobrow |
| gazard | perform | juice | publish | trunk | violent |
| woman | alden | nod | insult | mercy | sorrow |
| stand | probable | gentle | cardboard | anxious | fine |
| believe | signal | slip | cantileen | pedestrian | pity |
| gumm | dish | frequid | contract | arrow | normal |
| instead | earn | press | mount | feeble | tube |
| produce | sweat | provide | sumption | litholect | oak |
| group | fishlock | drum | moreover | brighten | hyslop |
| arrive | manage | reasonable | crisis | dam | difficult |
| manomize | mud | boil | jug | outlet | sandy |

B900

**Score 3900/5000**

# Appendix 13: X-lex Tamil Samples
## Sample 1

**Sample 2**

| | | | | | |
|---|---|---|---|---|---|
| நான் | முடியாத | நாலயம் | அடிக்கடி | இறுதி | மாநாடு |
| மக்கள் | தடப்பன் | பணம் | லட்சம் | மாலை | கிடைத்த |
| மூலம் | தெரியும் | கூடிய | கலை | பேருந்து | வெளி |
| நாட்கால | உலகம் | இல்லாத | சாத்திரம் | வாங்கி | மணர் |
| அரசியல | பகுதி | முன் | அறிவு | விழா | கண் |
| மேலும் | கோடி | எழுதி | காலை | தர | கால் |
| தான் | ஆலகரம் | நோய் | மாறி | தலை | பின்பு |
| மிகவும் | கண்டு | நாலயம் | பால் | பள்ளி | கத்தறி |
| ஆண்டு | மணம் | தரும் | நெருப்பு | பட்ட | பேசி |
| பெயர் | சாய்ப்பு | சற்று | தாய் | பொங்கள் | பிரிவு |
| புத்தளம் | சிறிய | செடு | மகன் | மருந்து | சீனா |
| நேரம் | நீர் | பற்றியும் | நம்மை | பண | நீர்பாச |
| காலம் | இங்கே | படி | கவிதை | வரி | கட்டி |
| வெற்றி | பெற்ற | சாக்காடு | கை | அரச | பலரும் |
| போய்ச் | அளவு | ஏற்பட்ட | அம்மா | கடல் | வாழ்வில் |
| முடிவு | உடல் | உயர் | பல்லிடம் | சூரியன் | பிரிவு |
| நாடு | பார்வை | குரல் | கொலை | நடிகர் | தலை |
| கல்வி | முடிவு | காதல் | உணர்வு | ரொம்ப | கயிரு |
| நன்றி | உண்மை | பயணம் | தாம் | வானிலை | நாயகம் |
| திங்கள் | மொழி | மாநில | நன்றாக | பிடித்து | தொடர் |

2400

**Score: 2400/5000**

254

**Sample 3**

| | | | | | |
|---|---|---|---|---|---|
| நான் | முடியாத | நாலயம் | அடிக்கடி | இறுதி | மாநாடு |
| மூங்கள் | தடப்பன் | பணம் | லட்சம் | மாலை | கிணைத்த |
| மூலம் | தெரியும் | கூடிய | கலை | பேருந்து | வெளி |
| நாட்கால | உலகம் | இல்லாத | சாத்திரம் | வாங்கி | மளார் |
| அரசியல | பகுதி | முன் | அறிவு | விடா | கண் |
| மேலும் | கோடி | எழுதி | காலை | தர | கால் |
| தான் | ஆலகரம் | நோய் | மாறி | தலை | பின்பு |
| மிகவும் | கண்டு | நாலயம் | பால் | பள்ளி | கத்தறி |
| ஆண்டு | மணம் | தரும் | நெருப்பு | பட்ட | பேசி |
| பெயர் | சாரிப்பு | சற்று | தாய் | பொங்கள் | பிரிவு |
| பத்தளாம் | சிறிய | வீடு | மகன் | மருந்து | சீனா |
| நேரம் | நீர் | பற்றியும் | நம்மை | பயண | நீர்பாச |
| காலம் | இங்கே | படி | கவிதை | வரி | கட்டி |
| வெற்றி | பெற்ற | சாக்காடு | கை | அரசு | பலரும் |
| போய்ச் | அளவு | ஏற்பட்ட | அம்மா | கடல் | வாழ்வில் |
| டிடிவு | உடல் | உயர | பல்லிடம் | சூரியன் | பிரிவு |
| ாடு | பார்வை | குரல் | கொலை | நடிகர் | தலை |
| ல்வி | முடிவு | காதல் | உணர்வு | ரொம்ப | கயிறு |
| ன்றி | உண்மை | பயணம் | தாம் | வானிலை | நாயகம் |
| ிங்கள | மொழி | மாநில | நன்றாக | பிடித்து | தொடர் |

3600

**Score: 3600/5000**

# Appendix 14: Language Use and Exposure Questionnaire

| Name: | | Initials | |
|---|---|---|---|

| Date of Birth | (date/ month/ year):        /      / |
|---|---|
| Place of birth: | City/ town:<br><br>Country: |
| Number of years in UK | |
| Number of years learning English | |

1. How many languages do you speak?...........................

2. Which language did you learn to speak first? ...............................................................

3. What do you consider to be your first language? .....................................................

4. How many languages do you speak at home? ...........................

5. What language(s) do you use when speaking to the following people (including English)?

**0**= Never  **1**=Rarely (1-30 % of the time)   **2**=Some (30-60 % of the time)   **3**=Most of the time (61- 90% of the time)    **4**=All of the time

| | Language 1 | | Language 2 | | Language 3 | | Language 4 | |
|---|---|---|---|---|---|---|---|---|
| Mother | | | | | | | | |
| Father | | | | | | | | |
| Siblings | | | | | | | | |

| School friends | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Other friends | | | | | | | | |

6. How often do you **speak English** in the following situations?

**0**= Never  **1**=Rarely (1-30 % of the time)   **2**=Some (30-60 % of the time)   **3**=Most of the time (61- 90% of the time)     **4**=All of the time

| Whilst having dinner with your immediate family (parents, brother, sister) | |
|---|---|
| When watching TV with your immediate family (parents, brother, sister) | |
| When talking to your friends at school during break | |
| When talking with your friends outside of school | |
| When you socialise with your friends outside of school | |

6. How often do you **speak Tamil** in the following situations?

**0**= Never  **1**=Rarely (1-30 % of the time)   **2**=Some (30-60 % of the time)   **3**=Most of the time (61- 90% of the time)     **4**=All of the time

| Whilst having dinner with your immediate family (parents, brother, sister) | |
|---|---|
| When watching TV with your immediate family (parents, brother, sister) | |
| When talking to your friends at school during break | |
| When talking with your friends outside of school | |
| When you socialise with your friends outside of school | |

**Appendix 15: Stimuli**

*Sentence 1 in each pair contains the congruent collocation and Sentence 2 contains the*

*incongruent collocation.*

1. John had to **keep in mind** that the weather might not be good for the football match.
   John had to **take a chance** that the weather might not be good for the football match.

2. Anna put **a pile of books** on the table very carefully when nobody was looking.
   Anna put **a bottle of water** on the table very carefully when nobody was looking.

3. Andrew didn't want to play football because it was **freezing cold** outside.
   Andrew didn't want to play football because it was **boiling hot** outside.

4. Annie wanted to **take a photo** at the beach before she went home.
   Annie wanted to **take a look** at the beach before she went home.

5. Early in the morning, my teacher asked me to **come prepared** for the meeting.
   Early in the morning, my teacher asked me to **get ready** for the meeting.

6. Bobby **made an effort** to finish his homework even though he was tired.
   Bobby **made a decision** to finish his homework even though he was tired.

7. My sister **broke her leg** last week so she can't come swimming today.
    My sister **caught a cold** last week so she can't come swimming today.

8. My mother wanted to **do the cooking** before we went shopping for clothes.
   My mother wanted to **have a haircut** before we went shopping for clothes.

9. Last week, the boys **broke the window** when they were playing football.
   Last week, the boys **broke the rules** when they were playing football.

10. Matthew **caught the ball** and won the game for his team.
    Matthew **broke a record** and won the game for his team.

11. Anna didn't **have a bath** before she went to bed last night.
    Anna didn't **brush her teeth** before she went to bed last night.

12. At the police station, the police officer asked Sam to **take a seat** immediately.
    At the police station, the police officer asked Sam to **pay the fine** immediately.

13. Mike decided to **go fishing** with his friend Ross on Saturday.
    Mike decided to **have dinner** with his friend Ross on Saturday.

14. Sally knew **make a mistake** if her brother didn't help her with her project.
    Sally knew she would **make a mess** if her brother didn't help her with her project.

15. When she visited her friend, Sophie had to **catch a bus** to get to the zoo.
    When she visited her friend, Sophie had to **take a taxi** to get to the zoo.

16. Alex was excited to **go abroad** during the summer holidays.
    Alex was excited to **go sailing** during the summer holidays.

17. Ben knew that his brother would **make trouble** if he didn't share his chocolate.
    Ben knew that his brother would **get angry** if he didn't share his chocolate.

18. She knew she could **solve the problem** if she really tried.
    She knew she could **make a difference** if she really tried.

19. My mother said I have to **do my homework** before I went to play.
    My mother said I have to **do the shopping** before I went to play.

20. I knew that I could not **be late** on my first day at work.
    I knew that I could not **get lost** on my first day at work.

21. My brother had to **pass an exam** to get the job at the new firm.
    My brother had to **take an exam** to get the job at the new firm.

22. I had to **make a list** before I went shopping at the supermarket.
    I had to **do the dishes** before I went shopping at the supermarket.

23. My friend wanted to buy a **pack of cards** from the stationery store.
    My friend wanted to buy a **pad of paper** from the stationery store.

24. My mother warned me not to **make a noise** in case the baby woke up.
    My mother warned me not to **make the call** in case the baby woke up.

25. Jeff had to be careful with his money so he could **pay the bills** at the end of the month.
    Jeff had to be careful with his money so he could **throw a party** at the end of the month.

26. My mother told me to always **lock the door** before I leave my house.
    My mother told me to always **make my bed** before I leave my house.

27. Laura worked hard to get her **happy ending** despite the challenges.
    Laura worked hard to get her **dream job** despite the challenges.

28. My aunt is going to **have a baby** in June so she applied for leave from work.

My aunt is going to **take a holiday** in June so she applied for leave from work.

29. I knew that I could not **waste time** because I had to finish my project.
    I knew that I could not **get upset** because I had to finish my project.

30. My father **spent time** to make sure that the washing machine would work properly.
    My father **took action** to make sure that the washing machine would work properly.

**Appendix 16: Age, Gender and TOWRE Scores of Participants**

| Participant | Age | Gender | TOWRE Reg | TOWRE Irreg |
|---:|---:|:---|---:|---:|
| 1 | 8.2 | F | 64 | 33 |
| 2 | 9.3 | M | 66 | 43 |
| 3 | 11.2 | M | 74 | 49 |
| 4 | 8.5 | F | 60 | 37 |
| 5 | 8.6 | F | 62 | 44 |
| 6 | 10.4 | F | 59 | 44 |
| 7 | 8.7 | F | 55 | 43 |
| 8 | 8.3 | F | 72 | 55 |
| 9 | 9.1 | F | 63 | 49 |
| 10 | 8.8 | F | 73 | 44 |
| 11 | 9.8 | F | 65 | 38 |
| 12 | 9.4 | F | 88 | 60 |
| 13 | 9.7 | F | 62 | 42 |
| 14 | 10.3 | F | 83 | 59 |
| 15 | 8.5 | M | 59 | 31 |
| 16 | 8.9 | F | 66 | 45 |
| 17 | 9.3 | M | 72 | 56 |
| 18 | 8.6 | M | 65 | 52 |
| 19 | 10 | M | 82 | 56 |
| 20 | 8.1 | F | 62 | 30 |
| 21 | 10.1 | M | 81 | 54 |
| 22 | 8.6 | M | 67 | 48 |
| 23 | 8.2 | M | 58 | 37 |
| 24 | 8.2 | M | 59 | 35 |
| 25 | 8.4 | M | 70 | 41 |
| 26 | 11.4 | F | 75 | 52 |
| 27 | 8.3 | F | 70 | 47 |
| 28 | 8.7 | F | 82 | 61 |
| 29 | 8.6 | M | 65 | 47 |
| 30 | 8.8 | F | 72 | 49 |
| 31 | 11.5 | M | 82 | 55 |
| 32 | 8.6 | F | 59 | 31 |
| 33 | 11.3 | M | 75 | 50 |
| 34 | 10.7 | F | 61 | 36 |
| 35 | 10.3 | F | 65 | 42 |
| 36 | 10 | F | 81 | 57 |
| 37 | 10.5 | F | 79 | 48 |
| 38 | 10.2 | M | 66 | 39 |
| 39 | 9.3 | F | 64 | 42 |
| 40 | 10.2 | M | 71 | 55 |
| 41 | 9 | M | 73 | 50 |
| 42 | 9.1 | M | 63 | 39 |

| 43 | 8.1 | F | 60 | 41 |
|---|---|---|---|---|
| 44 | 9.3 | F | 61 | 38 |
| 45 | 10.4 | F | 70 | 51 |
| 46 | 9.6 | F | 67 | 50 |
| 47 | 10.1 | F | 69 | 45 |
| 48 | 10.5 | F | 60 | 42 |
| 49 | 8.3 | F | 64 | 44 |
| 50 | 8.6 | F | 59 | 39 |
| 51 | 9.7 | F | 60 | 42 |
| 52 | 11.5 | M | 79 | 58 |
| 53 | 8.5 | M | 68 | 48 |
| 54 | 9.3 | M | 65 | 44 |
| 55 | 9.8 | M | 54 | 35 |
| 56 | 8.9 | M | 63 | 40 |
| 57 | 9.6 | F | 70 | 43 |
| 58 | 8.7 | F | 58 | 36 |
| 59 | 9.5 | F | 69 | 38 |
| 60 | 8.8 | F | 69 | 47 |
| 61 | 11.3 | F | 72 | 48 |
| 62 | 11.3 | M | 60 | 41 |
| 63 | 8.6 | F | 61 | 46 |
| 64 | 10.7 | M | 73 | 55 |
| 65 | 11.3 | F | 76 | 54 |
| 66 | 10.8 | F | 72 | 56 |
| 67 | 9.4 | F | 67 | 41 |
| 68 | 9.9 | F | 60 | 41 |
| 69 | 9.8 | M | 63 | 40 |
| 70 | 9.7 | M | 54 | 37 |
| 71 | 10.7 | M | 67 | 43 |
| 72 | 10.8 | M | 62 | 40 |
| 73 | 10.9 | M | 66 | 43 |
| 74 | 11.3 | F | 64 | 42 |
| 75 | 11.7 | F | 69 | 44 |
| 76 | 11.5 | F | 57 | 36 |
| 77 | 10.5 | M | 59 | 38 |
| 78 | 10.7 | F | 64 | 40 |
| 79 | 11.3 | F | 67 | 45 |
| 80 | 10.8 | F | 58 | 47 |
| Mean | 9.69 | | 66.82 | 44.78 |

**Appendix 17: Vocabulary Scores and Language Use/Exposure Scores**

| Participant | Tamil Vocab | English Vocab | Tamil use/exposure |
|---|---|---|---|
| 1 | 3600 | 4300 | 18 |
| 2 | 2800 | 3800 | 14 |
| 3 | 2850 | 4950 | 11 |
| 4 | 1700 | 4250 | 16 |
| 5 | 3900 | 4400 | 7 |
| 6 | 1750 | 3550 | 9 |
| 7 | 1800 | 3650 | 17 |
| 8 | 3750 | 4050 | 9 |
| 9 | 2000 | 3950 | 16 |
| 10 | 2500 | 4500 | 12 |
| 11 | 1750 | 3900 | 20 |
| 12 | 4000 | 4850 | 19 |
| 13 | 3000 | 4150 | 13 |
| 14 | 4600 | 4950 | 21 |
| 15 | 2800 | 3900 | 14 |
| 16 | 2550 | 3750 | 14 |
| 17 | 3600 | 4450 | 10 |
| 18 | 2950 | 3900 | 11 |
| 19 | 3300 | 4500 | 9 |
| 20 | 1600 | 3650 | 14 |
| 21 | 2200 | 4300 | 19 |
| 22 | 2400 | 4450 | 17 |
| 23 | 1600 | 3900 | 19 |
| 24 | 1400 | 4100 | 10 |
| 25 | 3700 | 4150 | 7 |
| 26 | 2200 | 4000 | 12 |
| 27 | 2000 | 4400 | 13 |
| 28 | 3300 | 4600 | 8 |
| 29 | 4600 | 3850 | 11 |
| 30 | 3500 | 4500 | 11 |
| 31 | 2300 | 4750 | 13 |
| 32 | 2600 | 4200 | 13 |
| 33 | 1750 | 4400 | 21 |
| 34 | 2800 | 3800 | 17 |
| 35 | 3000 | 4200 | 10 |
| 36 | 3600 | 4550 | 7 |
| 37 | 3600 | 4400 | 19 |
| 38 | 2000 | 3900 | 22 |
| 39 | 2550 | 4000 | 10 |

| 40 | 2250 | 4300 | 17 |
|---|---|---|---|
| 41 | 3000 | 4400 | 12 |
| 42 | 2900 | 3750 | 17 |
| 43 | 2400 | 3950 | 20 |
| 44 | 2650 | 3700 | 18 |
| 45 | 3150 | 4750 | 16 |
| 46 | 4100 | 4600 | 15 |
| 47 | 3450 | 4250 | 11 |
| 48 | 2200 | 4000 | 10 |
| 49 | 2800 | 4350 | 18 |
| 50 | 2400 | 3900 | 12 |
| 51 | 2450 | 3800 | 16 |
| 52 | 2700 | 4550 | 17 |
| 53 | 4150 | 4600 | 20 |
| 54 | 3000 | 4400 | 19 |
| 55 | 3750 | 3900 | 9 |
| 56 | 3600 | 4250 | 10 |
| 57 | 3200 | 4600 | 7 |
| 58 | 2800 | 3900 | 19 |
| 59 | 2900 | 4200 | 12 |
| 60 | 3100 | 4250 | 19 |
| 61 | 2700 | 4300 | 12 |
| 62 | 2350 | 4000 | 12 |
| 63 | 1600 | 4100 | 21 |
| 64 | 3000 | 4300 | 12 |
| 65 | 1800 | 4550 | 12 |
| 66 | 3600 | 4400 | 16 |
| 67 | 3300 | 4200 | 19 |
| 68 | 3150 | 3900 | 8 |
| 69 | 3000 | 4300 | 20 |
| 70 | 2800 | 3800 | 9 |
| 71 | 2850 | 4300 | 11 |
| 72 | 2700 | 4000 | 16 |
| 73 | 2900 | 4100 | 13 |
| 74 | 3300 | 4100 | 7 |
| 75 | 3200 | 4500 | 21 |
| 76 | 2400 | 3800 | 11 |
| 77 | 3100 | 3900 | 11 |
| 78 | 3000 | 4000 | 9 |
| 79 | 3200 | 4100 | 22 |
| 80 | 2900 | 3800 | 13 |

**Appendix 18: Reading times for whole collocations for all four measures (in milliseconds)**

*FFcong: first fixation congruent, FFincong; first fixation incongruent; SFcong: single fixation congruent, SFincong: single fixation incongruent; Gazecong: gaze congruent, Gazeincong: gaze incongruent; TTcong: total time congruent, TTincong: total time incongruent*

| Participant | FFcong | FFincong | SFcong | SFcong | Gazecong | Gazeincong | TTcong | TTIncong |
|---|---|---|---|---|---|---|---|---|
| 1 | 203.00 | 196.90 | 161.25 | 188.33 | 203.00 | 233.10 | 259.50 | 339.00 |
| 2 | 262.17 | 239.00 | 190.57 | 238.22 | 379.67 | 323.17 | 472.31 | 594.83 |
| 3 | 208.50 | 209.60 | 178.00 | 300.67 | 208.50 | 223.80 | 544.86 | 468.46 |
| 4 | 218.00 | 204.33 | 191.00 | 185.60 | 305.64 | 234.17 | 584.83 | 526.92 |
| 5 | 246.33 | 260.83 | 207.33 | 319.33 | 355.33 | 260.83 | 688.43 | 652.31 |
| 6 | 352.73 | 246.91 | 300.67 | 276.29 | 396.91 | 352.55 | 903.38 | 1175.69 |
| 7 | 511.07 | 420.27 | 431.64 | 266.67 | 643.64 | 601.73 | 675.07 | 585.83 |
| 8 | 230.18 | 322.14 | 282.40 | 426.80 | 279.82 | 350.86 | 379.27 | 539.71 |
| 9 | 272.91 | 239.09 | 284.67 | 271.00 | 376.36 | 282.55 | 637.85 | 559.54 |
| 10 | 310.00 | 229.33 | 313.60 | 310.00 | 324.33 | 243.00 | 499.54 | 485.23 |
| 11 | 358.00 | 291.86 | 273.75 | 257.25 | 420.73 | 410.71 | 686.31 | 618.71 |
| 12 | 256.50 | 275.27 | 247.80 | 275.11 | 371.00 | 428.36 | 465.33 | 543.33 |
| 13 | 405.00 | 530.00 | 622.67 | 387.43 | 594.86 | 754.00 | 929.29 | 1030.15 |
| 14 | 206.67 | 245.56 | 299.75 | 225.71 | 256.67 | 272.44 | 442.31 | 387.60 |
| 15 | 478.30 | 456.14 | 344.40 | 441.43 | 703.60 | 478.36 | 733.30 | 648.36 |
| 16 | 228.73 | 335.11 | 191.00 | 391.71 | 450.91 | 624.89 | 663.60 | 784.14 |
| 17 | 369.43 | 387.71 | 342.00 | 330.00 | 397.71 | 625.71 | 483.71 | 596.00 |
| 18 | 403.85 | 391.11 | 388.20 | 573.57 | 518.38 | 571.89 | 627.13 | 649.11 |
| 19 | 248.62 | 267.38 | 228.00 | 318.67 | 329.08 | 300.92 | 500.00 | 364.62 |
| 20 | 327.43 | 270.73 | 310.00 | 302.67 | 524.86 | 433.20 | 1128.43 | 846.55 |
| 21 | 203.07 | 324.15 | 228.33 | 209.71 | 321.07 | 546.62 | 382.67 | 918.71 |
| 22 | 297.20 | 279.33 | 409.00 | 302.67 | 313.20 | 416.67 | 463.27 | 994.00 |
| 23 | 272.86 | 245.82 | 479.00 | 227.50 | 502.57 | 307.09 | 631.75 | 565.64 |
| 24 | 203.71 | 196.00 | 299.00 | 258.67 | 262.86 | 196.00 | 1057.56 | 684.73 |
| 25 | 190.00 | 187.00 | 147.71 | 265.50 | 260.00 | 239.40 | 314.77 | 425.27 |
| 26 | 232.77 | 260.91 | 248.36 | 316.44 | 343.85 | 336.18 | 495.69 | 429.83 |
| 27 | 341.60 | 293.86 | 271.43 | 248.40 | 601.40 | 415.43 | 1071.00 | 856.14 |
| 28 | 407.00 | 371.50 | 342.00 | 330.00 | 431.75 | 579.75 | 507.00 | 558.44 |
| 29 | 406.46 | 488.71 | 376.75 | 358.29 | 520.46 | 773.23 | 1009.29 | 1029.86 |
| 30 | 199.14 | 188.80 | 279.00 | 163.50 | 199.14 | 188.80 | 316.25 | 375.00 |
| 31 | 203.33 | 206.57 | 211.20 | 202.00 | 271.00 | 298.57 | 270.50 | 414.00 |
| 32 | 336.00 | 332.36 | 354.44 | 369.67 | 454.29 | 483.45 | 683.43 | 656.17 |
| 33 | 262.40 | 235.43 | 250.20 | 265.60 | 358.67 | 392.00 | 844.57 | 737.71 |
| 34 | 257.83 | 326.31 | 272.29 | 432.91 | 459.33 | 551.08 | 1272.62 | 1087.60 |
| 35 | 254.00 | 277.00 | 266.00 | 239.71 | 375.09 | 323.60 | 979.08 | 961.45 |

transcription>

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 36 | 239.00 | 222.00 | 196.67 | 264.25 | 379.71 | 349.09 | 494.57 | 649.27 |
| 37 | 188.83 | 188.60 | 149.67 | 265.50 | 264.67 | 260.80 | 314.67 | 469.82 |
| 38 | 303.64 | 298.00 | 405.71 | 289.40 | 571.64 | 352.78 | 724.57 | 666.69 |
| 39 | 257.92 | 175.13 | 182.50 | 180.67 | 266.77 | 189.13 | 459.54 | 582.75 |
| 40 | 243.31 | 246.30 | 228.86 | 223.75 | 408.85 | 370.10 | 656.69 | 443.92 |
| 41 | 225.00 | 244.29 | 248.25 | 301.60 | 278.17 | 244.29 | 355.14 | 413.50 |
| 42 | 261.78 | 265.57 | 241.00 | 276.71 | 429.22 | 382.57 | 448.22 | 616.23 |
| 43 | 386.82 | 327.90 | 447.57 | 299.50 | 472.36 | 613.80 | 606.23 | 780.45 |
| 44 | 290.67 | 296.70 | 347.67 | 532.75 | 491.07 | 551.90 | 830.80 | 650.50 |
| 45 | 295.77 | 214.71 | 269.00 | 283.50 | 437.69 | 399.86 | 549.85 | 574.57 |
| 46 | 210.93 | 267.93 | 231.63 | 182.86 | 391.07 | 301.00 | 604.87 | 619.43 |
| 47 | 260.70 | 263.46 | 194.71 | 261.83 | 380.60 | 379.00 | 476.08 | 631.14 |
| 48 | 309.17 | 289.50 | 398.43 | 354.25 | 389.58 | 353.75 | 756.21 | 511.43 |
| 49 | 364.00 | 242.00 | 280.00 | 395.00 | 607.00 | 431.22 | 645.30 | 709.40 |
| 50 | 423.54 | 360.31 | 478.00 | 325.67 | 975.23 | 910.31 | 1170.79 | 1119.38 |
| 51 | 295.33 | 336.00 | 246.00 | 328.00 | 669.25 | 713.11 | 765.31 | 788.80 |
| 52 | 306.00 | 233.91 | 221.20 | 247.29 | 519.10 | 301.73 | 774.23 | 562.53 |
| 53 | 199.14 | 188.80 | 279.00 | 163.50 | 199.14 | 188.80 | 316.25 | 375.00 |
| 54 | 219.46 | 228.54 | 202.18 | 208.00 | 324.85 | 308.15 | 384.77 | 356.85 |
| 55 | 289.00 | 241.67 | 331.57 | 339.60 | 526.73 | 621.22 | 786.17 | 1170.80 |
| 56 | 248.20 | 292.30 | 232.11 | 256.00 | 330.67 | 386.80 | 528.87 | 526.91 |
| 57 | 223.82 | 231.55 | 291.17 | 249.50 | 380.73 | 518.09 | 500.79 | 618.67 |
| 58 | 243.18 | 261.82 | 311.17 | 660.50 | 412.45 | 318.27 | 851.45 | 793.17 |
| 59 | 259.14 | 223.43 | 284.23 | 309.70 | 407.14 | 421.93 | 577.00 | 805.64 |
| 60 | 203.00 | 196.90 | 161.25 | 188.33 | 203.00 | 233.10 | 259.50 | 339.00 |
| 61 | 273.33 | 284.45 | 288.50 | 252.30 | 327.50 | 333.91 | 523.42 | 454.33 |
| 62 | 200.88 | 263.50 | 183.50 | 181.20 | 302.13 | 303.63 | 498.30 | 543.17 |
| 63 | 308.77 | 277.88 | 291.88 | 328.00 | 454.31 | 615.75 | 498.40 | 869.22 |
| 64 | 203.69 | 229.18 | 153.50 | 235.67 | 214.77 | 342.09 | 304.77 | 471.18 |
| 65 | 234.58 | 230.44 | 338.00 | 223.00 | 379.08 | 243.89 | 608.57 | 817.25 |
| 66 | 182.67 | 213.85 | 360.00 | 265.29 | 182.67 | 240.23 | 492.14 | 340.15 |
| 67 | 275.55 | 239.62 | 236.90 | 234.30 | 358.18 | 455.77 | 429.38 | 694.86 |
| 68 | 336.29 | 269.80 | 301.33 | 242.50 | 480.14 | 407.13 | 560.00 | 535.93 |
| 69 | 273.85 | 256.00 | 273.00 | 255.00 | 314.85 | 380.00 | 453.77 | 467.57 |
| 70 | 383.62 | 402.67 | 329.60 | 393.60 | 739.54 | 668.64 | 946.92 | 1006.53 |
| 71 | 231.54 | 265.20 | 268.50 | 264.00 | 371.77 | 335.80 | 569.21 | 435.92 |
| 72 | 268.73 | 257.56 | 247.17 | 201.20 | 357.09 | 351.56 | 896.57 | 1015.38 |
| 73 | 283.29 | 291.58 | 275.50 | 289.80 | 352.43 | 378.33 | 454.29 | 512.50 |
| 74 | 418.43 | 278.33 | 462.88 | 235.57 | 626.86 | 278.33 | 804.55 | 703.42 |
| 75 | 269.86 | 297.25 | 262.50 | 270.25 | 269.86 | 297.25 | 516.91 | 572.09 |
| 76 | 480.55 | 544.83 | 299.86 | 555.29 | 617.55 | 616.67 | 852.31 | 865.62 |
| 77 | 587.29 | 432.91 | 608.29 | 625.60 | 691.14 | 543.27 | 856.30 | 709.82 |
| 78 | 366.30 | 443.07 | 316.50 | 535.29 | 503.80 | 562.73 | 562.90 | 858.80 |
| 79 | 445.83 | 381.36 | 333.38 | 696.50 | 540.67 | 605.27 | 593.54 | 618.00 |
| 80 | 607.17 | 562.89 | 425.83 | 502.56 | 698.83 | 800.56 | 685.08 | 938.45 |
| | 294.30 | 287.41 | 291.74 | 308.15 | 415.19 | 411.45 | 619.36 | 653.71 |

**Appendix 19: Reading times for Word 1 for all four measures (in milliseconds)**

| Participant | FFWord1Cong | FFWord1Incong | SFWord1Cong | SFWord1Incong | GazeWord1Cong | GazeWord1Incong | TTWord1Cong | TTWord1Incong |
|---|---|---|---|---|---|---|---|---|
| 1 | 161.25 | 188.33 | 203 | 176.29 | 161.25 | 188.33 | 216 | 266.5 |
| 2 | 188 | 235 | 316.29 | 185.33 | 218 | 256.6 | 364.4 | 342.33 |
| 3 | 204 | 300.67 | 208.55 | 212.57 | 204 | 300.67 | 656 | 312 |
| 4 | 188 | 211.71 | 206.29 | 206.4 | 204.89 | 262 | 384.67 | 487.45 |
| 5 | 212 | 274.6 | 294.67 | 255.33 | 212 | 292.4 | 446 | 862.5 |
| 6 | 300.67 | 276.29 | 396.5 | 262.86 | 300.67 | 276.29 | 692.77 | 845.64 |
| 7 | 431.64 | 262.64 | 591.29 | 461.38 | 431.64 | 261.3 | 442.58 | 387.58 |
| 8 | 282.4 | 396.86 | 224.33 | 340.5 | 282.4 | 430.57 | 292.67 | 455.09 |
| 9 | 267.71 | 271 | 297.2 | 293.2 | 267.71 | 271 | 441.45 | 307.33 |
| 10 | 307.64 | 294.2 | 335.4 | 230.22 | 307.64 | 311 | 409.83 | 364.18 |
| 11 | 260.44 | 262.55 | 383.43 | 362 | 260.44 | 262.55 | 424.62 | 436.62 |
| 12 | 249.33 | 283 | 261.75 | 283.6 | 284 | 301.4 | 284 | 301.4 |
| 13 | 608.77 | 406.5 | 498 | 527.2 | 669.23 | 406.5 | 798.71 | 525.2 |
| 14 | 270.6 | 225.5 | 216.25 | 256.67 | 286.6 | 287.75 | 387.85 | 309.08 |
| 15 | 344.4 | 402.75 | 582.17 | 463.6 | 344.4 | 402.75 | 344.4 | 426.8 |
| 16 | 197 | 386 | 279.33 | 321.33 | 254.67 | 436.44 | 543.17 | 610.44 |
| 17 | 342 | 351.2 | 374 | 280.67 | 342 | 401.2 | 342 | 443.67 |
| 18 | 388.2 | 573.57 | 464.11 | 405.67 | 388.2 | 573.57 | 454.25 | 582.13 |
| 19 | 228 | 318.67 | 262.75 | 281 | 228 | 318.67 | 285.67 | 380.92 |
| 20 | 306 | 274 | 527.5 | 259.33 | 507.56 | 298.75 | 647.71 | 531.17 |
| 21 | 217.83 | 219.6 | 186 | 418.5 | 300.5 | 275 | 462.33 | 472.46 |
| 22 | 311 | 264.8 | 344.67 | 342 | 449 | 320 | 385.33 | 568.67 |
| 23 | 387.67 | 208 | 367 | 183.2 | 482 | 270.57 | 612.75 | 378.4 |
| 24 | 284.57 | 258.67 | 198.4 | 206 | 309.71 | 258.67 | 788.5 | 684.67 |
| 25 | 158.22 | 217.33 | 190.33 | 198.57 | 158.22 | 250.67 | 307.33 | 364.44 |
| 26 | 248.36 | 302.36 | 245.2 | 243.14 | 248.36 | 327.64 | 304.17 | 495.27 |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 27 | 253.45 | 266.4 | 359 | 336.25 | 375.82 | 417.6 | 564.91 | 655 |
| 28 | 308.57 | 319.67 | 416.29 | 275 | 308.57 | 377.33 | 347.43 | 417.14 |
| 29 | 340.83 | 324.55 | 426.67 | 470.86 | 457.83 | 432 | 592.62 | 513.08 |
| 30 | 279 | 163.5 | 199.14 | 187.75 | 279 | 163.5 | 400 | 363.4 |
| 31 | 252.33 | 202 | 189 | 245.5 | 252.33 | 202 | 397 | 262.33 |
| 32 | 326.36 | 340.6 | 362 | 454.5 | 381.09 | 454.4 | 610 | 746.4 |
| 33 | 252.31 | 238.31 | 284 | 239 | 285.69 | 250.77 | 659.2 | 461.38 |
| 34 | 266.83 | 388.53 | 439.33 | 354.8 | 422.83 | 435.07 | 1077.69 | 858.8 |
| 35 | 234.67 | 318.22 | 367.5 | 313 | 260.83 | 334.67 | 838.15 | 683.82 |
| 36 | 208.5 | 262.44 | 243.67 | 222.67 | 208.5 | 262.44 | 437.8 | 407.23 |
| 37 | 161 | 239.6 | 187.6 | 202 | 161 | 239.6 | 281.09 | 343.33 |
| 38 | 425.67 | 297 | 340 | 356.6 | 471.11 | 455.11 | 438.77 | 509.23 |
| 39 | 182.5 | 180.67 | 259 | 181.86 | 182.5 | 180.67 | 278.29 | 518.33 |
| 40 | 250.6 | 236.22 | 237.5 | 197.33 | 250.6 | 268.56 | 371.15 | 314.91 |
| 41 | 227 | 301.6 | 230.22 | 287.6 | 271.2 | 301.6 | 344 | 305.43 |
| 42 | 235.09 | 276.71 | 291.33 | 270.67 | 258.36 | 276.71 | 291.45 | 482.23 |
| 43 | 402.27 | 355 | 505 | 372.5 | 468.09 | 601.25 | 661 | 771.3 |
| 44 | 266.58 | 324.2 | 406.33 | 214 | 446 | 444.9 | 616.79 | 678 |
| 45 | 223.14 | 276.75 | 304.2 | 298.33 | 295 | 336.42 | 378.22 | 549.14 |
| 46 | 235.43 | 195.14 | 219.17 | 275.7 | 305.57 | 292.64 | 571.6 | 466.93 |
| 47 | 194.71 | 261.83 | 249 | 232 | 194.71 | 261.83 | 290.38 | 389.7 |
| 48 | 364.75 | 313.67 | 285.38 | 272.17 | 387.38 | 377.83 | 519.22 | 382.33 |
| 49 | 239.11 | 342.78 | 339 | 255.33 | 335.22 | 388.22 | 450.4 | 596.5 |
| 50 | 332 | 336.3 | 421 | 179 | 465.77 | 388.9 | 764.85 | 666.5 |
| 51 | 238.2 | 271.6 | 274.25 | 490.33 | 370.4 | 398.4 | 537 | 584.8 |
| 52 | 194.78 | 249 | 382.75 | 246.4 | 194.78 | 273.63 | 569.45 | 442.5 |
| 53 | 279 | 163.5 | 199.14 | 187.75 | 279 | 163.5 | 400 | 363.4 |
| 54 | 202.18 | 198.73 | 222.89 | 232 | 202.18 | 198.73 | 275.73 | 337.55 |
| 55 | 347.18 | 300.14 | 290.6 | 329 | 463.36 | 367.86 | 543.09 | 656.89 |
| 56 | 246.45 | 215 | 291.17 | 328.17 | 255.27 | 292.75 | 362.15 | 475.67 |
| 57 | 289.71 | 234.75 | 223.8 | 334 | 320.29 | 260.25 | 397.64 | 408.9 |

| | | | | | | | |
|------|--------|--------|--------|--------|--------|--------|--------|--------|
| 58 | 270.11 | 395.88 | 233 | 395.33 | 335.56 | 578.5 | 885.7 | 728.31 |
| 59 | 275.21 | 309.7 | 307 | 318 | 275.21 | 309.7 | 336.2 | 488.21 |
| 60 | 161.25 | 188.33 | 203 | 176.29 | 161.25 | 188.33 | 216 | 266.5 |
| 61 | 273.4 | 252.3 | 275.71 | 311.14 | 286.8 | 252.3 | 317.38 | 412.1 |
| 62 | 183.5 | 184.5 | 174 | 280.14 | 183.5 | 184.5 | 322.5 | 399.83 |
| 63 | 236 | 288.43 | 290.6 | 361 | 285.43 | 378.43 | 353.2 | 538.88 |
| 64 | 175.33 | 226.86 | 209.1 | 246.25 | 175.33 | 226.86 | 304.18 | 347.33 |
| 65 | 285.88 | 229.86 | 244.4 | 227 | 332.13 | 229.86 | 554.67 | 534.82 |
| 66 | 360 | 236 | 222.5 | 223.64 | 360 | 236 | 236.67 | 307.38 |
| 67 | 236.9 | 234.3 | 254.71 | 258 | 236.9 | 234.3 | 292.82 | 436.08 |
| 68 | 328.73 | 227.4 | 350.5 | 318.2 | 328.73 | 227.4 | 636.55 | 352.33 |
| 69 | 270.09 | 263.27 | 268 | 333.4 | 431.36 | 312.45 | 562.69 | 477.38 |
| 70 | 336.82 | 353.33 | 292 | 452 | 336.82 | 410.08 | 581.62 | 480.86 |
| 71 | 262.6 | 254.25 | 257.2 | 285.63 | 285.6 | 308.13 | 265 | 527 |
| 72 | 237.29 | 201.2 | 252 | 269.57 | 383.43 | 201.2 | 441.08 | 302.8 |
| 73 | 276.85 | 289.8 | 297.63 | 339.5 | 299.23 | 289.8 | 324.47 | 337.38 |
| 74 | 439.78 | 226 | 559 | 248.17 | 439.78 | 226 | 578.5 | 591.64 |
| 75 | 262.5 | 287.2 | 263.83 | 271.14 | 262.5 | 425.4 | 407.08 | 578.36 |
| 76 | 302.1 | 485.8 | 451.2 | 621.63 | 302.1 | 485.8 | 558.7 | 695.46 |
| 77 | 608.29 | 544.13 | 676.2 | 477.83 | 608.29 | 589.5 | 697.33 | 646.09 |
| 78 | 350.11 | 535.29 | 395.43 | 662.17 | 591.78 | 535.29 | 748.67 | 500.73 |
| 79 | 313.6 | 475.4 | 499.11 | 338.8 | 365.7 | 701.4 | 442.23 | 622.5 |
| 80 | 378.5 | 485.4 | 678.3 | 733.25 | 378.5 | 485.4 | 451.63 | 588.18 |
| Mean | 280.38 | 290.44 | 319.43 | 308.08 | 319.49 | 329.48 | 468.76 | 486.43 |

**Appendix 20: Reading times for Word 2 for all four measures (in milliseconds)**

| Participant | FFWord2 Cong | FFWord2 Incong | SFWord2 Cong | SFWord2 Incong | GazeWord2 Cong | GazeWord2 Incong | TTWord2 Cong | TTWord2 Incong |
|---|---|---|---|---|---|---|---|---|
| 1 | 203 | 196.9 | 161.25 | 188.33 | 203 | 233.1 | 259.5 | 339 |
| 2 | 310 | 229.33 | 313.6 | 310 | 324.33 | 243 | 499.54 | 485.23 |
| 3 | 358 | 291.86 | 273.75 | 257.25 | 420.73 | 410.71 | 686.31 | 618.71 |
| 4 | 256.5 | 275.27 | 247.8 | 275.11 | 371 | 428.36 | 465.33 | 543.33 |
| 5 | 405 | 530 | 622.67 | 387.43 | 594.86 | 754 | 929.29 | 1030.15 |
| 6 | 206.67 | 245.56 | 299.75 | 225.71 | 256.67 | 272.44 | 442.31 | 387.6 |
| 7 | 478.3 | 456.14 | 344.4 | 441.43 | 703.6 | 478.36 | 733.3 | 648.36 |
| 8 | 228.73 | 335.11 | 191 | 391.71 | 450.91 | 624.89 | 663.6 | 784.14 |
| 9 | 369.43 | 387.71 | 342 | 330 | 397.71 | 625.71 | 483.71 | 596 |
| 10 | 403.85 | 391.11 | 388.2 | 573.57 | 518.38 | 571.89 | 627.13 | 649.11 |
| 11 | 248.62 | 267.38 | 228 | 318.67 | 329.08 | 300.92 | 500 | 364.62 |
| 12 | 262.17 | 239 | 190.57 | 238.22 | 379.67 | 323.17 | 472.31 | 594.83 |
| 13 | 327.43 | 270.73 | 310 | 302.67 | 524.86 | 433.2 | 1128.43 | 846.55 |
| 14 | 203.07 | 324.15 | 228.33 | 209.71 | 321.07 | 546.62 | 382.67 | 918.71 |
| 15 | 297.2 | 279.33 | 409 | 302.67 | 313.2 | 416.67 | 463.27 | 994 |
| 16 | 272.86 | 245.82 | 479 | 227.5 | 502.57 | 307.09 | 631.75 | 565.64 |
| 17 | 203.71 | 196 | 299 | 258.67 | 262.86 | 196 | 1057.56 | 684.73 |
| 18 | 190 | 187 | 147.71 | 265.5 | 260 | 239.4 | 314.77 | 425.27 |
| 19 | 232.77 | 260.91 | 248.36 | 316.44 | 343.85 | 336.18 | 495.69 | 429.83 |
| 20 | 341.6 | 293.86 | 271.43 | 248.4 | 601.4 | 415.43 | 1071 | 856.14 |
| 21 | 407 | 371.5 | 342 | 330 | 431.75 | 579.75 | 507 | 558.44 |
| 22 | 406.46 | 488.71 | 376.75 | 358.29 | 520.46 | 773.23 | 1009.29 | 1029.86 |
| 23 | 208.5 | 209.6 | 178 | 300.67 | 208.5 | 223.8 | 544.86 | 468.46 |
| 24 | 199.14 | 188.8 | 279 | 163.5 | 199.14 | 188.8 | 316.25 | 375 |
| 25 | 203.33 | 206.57 | 211.2 | 202 | 271 | 298.57 | 270.5 | 414 |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 26 | 336 | 332.36 | 354.44 | 369.67 | 454.29 | 483.45 | 683.43 | 656.17 |
| 27 | 262.4 | 235.43 | 250.2 | 265.6 | 358.67 | 392 | 844.57 | 737.71 |
| 28 | 257.83 | 326.31 | 272.29 | 432.91 | 459.33 | 551.08 | 1272.62 | 1087.6 |
| 29 | 254 | 277 | 266 | 239.71 | 375.09 | 323.6 | 979.08 | 961.45 |
| 30 | 239 | 222 | 196.67 | 264.25 | 379.71 | 349.09 | 494.57 | 649.27 |
| 31 | 188.83 | 188.6 | 149.67 | 265.5 | 264.67 | 260.8 | 314.67 | 469.82 |
| 32 | 303.64 | 298 | 405.71 | 289.4 | 571.64 | 352.78 | 724.57 | 666.69 |
| 33 | 257.92 | 175.13 | 182.5 | 180.67 | 266.77 | 189.13 | 459.54 | 582.75 |
| 34 | 218 | 204.33 | 191 | 185.6 | 305.64 | 234.17 | 584.83 | 526.92 |
| 35 | 243.31 | 246.3 | 228.86 | 223.75 | 408.85 | 370.1 | 656.69 | 443.92 |
| 36 | 225 | 244.29 | 248.25 | 301.6 | 278.17 | 244.29 | 355.14 | 413.5 |
| 37 | 261.78 | 265.57 | 241 | 276.71 | 429.22 | 382.57 | 448.22 | 616.23 |
| 38 | 386.82 | 327.9 | 447.57 | 299.5 | 472.36 | 613.8 | 606.23 | 780.45 |
| 39 | 290.67 | 296.7 | 347.67 | 532.75 | 491.07 | 551.9 | 830.8 | 650.5 |
| 40 | 295.77 | 214.71 | 269 | 283.5 | 437.69 | 399.86 | 549.85 | 574.57 |
| 41 | 210.93 | 267.93 | 231.63 | 182.86 | 391.07 | 301 | 604.87 | 619.43 |
| 42 | 260.7 | 263.46 | 194.71 | 261.83 | 380.6 | 379 | 476.08 | 631.14 |
| 43 | 309.17 | 289.5 | 398.43 | 354.25 | 389.58 | 353.75 | 756.21 | 511.43 |
| 44 | 364 | 242 | 280 | 395 | 607 | 431.22 | 645.3 | 709.4 |
| 45 | 246.33 | 260.83 | 207.33 | 319.33 | 355.33 | 260.83 | 688.43 | 652.31 |
| 46 | 423.54 | 360.31 | 478 | 325.67 | 975.23 | 910.31 | 1170.79 | 1119.38 |
| 47 | 295.33 | 336 | 246 | 328 | 669.25 | 713.11 | 765.31 | 788.8 |
| 48 | 306 | 233.91 | 221.2 | 247.29 | 519.1 | 301.73 | 774.23 | 562.53 |
| 49 | 199.14 | 188.8 | 279 | 163.5 | 199.14 | 188.8 | 316.25 | 375 |
| 50 | 219.46 | 228.54 | 202.18 | 208 | 324.85 | 308.15 | 384.77 | 356.85 |
| 51 | 289 | 241.67 | 331.57 | 339.6 | 526.73 | 621.22 | 786.17 | 1170.8 |
| 52 | 248.2 | 292.3 | 232.11 | 256 | 330.67 | 386.8 | 528.87 | 526.91 |
| 53 | 223.82 | 231.55 | 291.17 | 249.5 | 380.73 | 518.09 | 500.79 | 618.67 |
| 54 | 243.18 | 261.82 | 311.17 | 660.5 | 412.45 | 318.27 | 851.45 | 793.17 |
| 55 | 259.14 | 223.43 | 284.23 | 309.7 | 407.14 | 421.93 | 577 | 805.64 |
| 56 | 352.73 | 246.91 | 300.67 | 276.29 | 396.91 | 352.55 | 903.38 | 1175.69 |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **57** | 203 | 196.9 | 161.25 | 188.33 | 203 | 233.1 | 259.5 | 339 |
| **58** | 273.33 | 284.45 | 288.5 | 252.3 | 327.5 | 333.91 | 523.42 | 454.33 |
| **59** | 200.88 | 263.5 | 183.5 | 181.2 | 302.13 | 303.63 | 498.3 | 543.17 |
| **60** | 308.77 | 277.88 | 291.88 | 328 | 454.31 | 615.75 | 498.4 | 869.22 |
| **61** | 203.69 | 229.18 | 153.5 | 235.67 | 214.77 | 342.09 | 304.77 | 471.18 |
| **62** | 234.58 | 230.44 | 338 | 223 | 379.08 | 243.89 | 608.57 | 817.25 |
| **63** | 182.67 | 213.85 | 360 | 265.29 | 182.67 | 240.23 | 492.14 | 340.15 |
| **64** | 275.55 | 239.62 | 236.9 | 234.3 | 358.18 | 455.77 | 429.38 | 694.86 |
| **65** | 336.29 | 269.8 | 301.33 | 242.5 | 480.14 | 407.13 | 560 | 535.93 |
| **66** | 273.85 | 256 | 273 | 255 | 314.85 | 380 | 453.77 | 467.57 |
| **67** | 511.07 | 420.27 | 431.64 | 266.67 | 643.64 | 601.73 | 675.07 | 585.83 |
| **68** | 383.62 | 402.67 | 329.6 | 393.6 | 739.54 | 668.64 | 946.92 | 1006.53 |
| **69** | 231.54 | 265.2 | 268.5 | 264 | 371.77 | 335.8 | 569.21 | 435.92 |
| **70** | 268.73 | 257.56 | 247.17 | 201.2 | 357.09 | 351.56 | 896.57 | 1015.38 |
| **71** | 283.29 | 291.58 | 275.5 | 289.8 | 352.43 | 378.33 | 454.29 | 512.5 |
| **72** | 418.43 | 278.33 | 462.88 | 235.57 | 626.86 | 278.33 | 804.55 | 703.42 |
| **73** | 269.86 | 297.25 | 262.5 | 270.25 | 269.86 | 297.25 | 516.91 | 572.09 |
| **74** | 480.55 | 544.83 | 299.86 | 555.29 | 617.55 | 616.67 | 852.31 | 865.62 |
| **75** | 587.29 | 432.91 | 608.29 | 625.6 | 691.14 | 543.27 | 856.3 | 709.82 |
| **76** | 366.3 | 443.07 | 316.5 | 535.29 | 503.8 | 562.73 | 562.9 | 858.8 |
| **77** | 445.83 | 381.36 | 333.38 | 696.5 | 540.67 | 605.27 | 593.54 | 618 |
| **78** | 230.18 | 322.14 | 282.4 | 426.8 | 279.82 | 350.86 | 379.27 | 539.71 |
| **79** | 607.17 | 562.89 | 425.83 | 502.56 | 698.83 | 800.56 | 685.08 | 938.45 |
| **80** | 272.91 | 239.09 | 284.67 | 271 | 376.36 | 282.55 | 637.85 | 559.54 |
| **Mean** | 287.3045 | 294.408875 | 291.74475 | 308.145125 | 415.19425 | 411.4465 | 619.3638 | 653.7085 |

**Appendix 21: Difference Congruency Effect (DCE) for Word 1 and Word 2**

| FFWord1DCE | SFWord1DCE | GazeWord1DCE | TTWord1DCE | FFWord2DCE | SFWord2DCE | GazeWord2DCE | TTWord2DCE |
|---|---|---|---|---|---|---|---|
| 27.08 | -26.71 | -53.79 | -27.08 | 26.71 | 53.79 | 27.08 | -26.71 |
| 47 | -130.96 | -177.96 | -47 | 130.96 | 177.96 | 47 | -130.96 |
| 96.67 | 4.02 | -92.65 | -96.67 | -4.02 | 92.65 | 96.67 | 4.02 |
| 23.71 | 0.11 | -23.6 | -23.71 | -0.11 | 23.6 | 23.71 | 0.11 |
| 62.6 | -39.34 | -101.94 | -62.6 | 39.34 | 101.94 | 62.6 | -39.34 |
| -24.38 | -133.64 | -109.26 | 24.38 | 133.64 | 109.26 | -24.38 | -133.64 |
| -169 | -129.91 | 39.09 | 169 | 129.91 | -39.09 | -169 | -129.91 |
| 114.46 | 116.17 | 1.71 | -114.46 | -116.17 | -1.71 | 114.46 | 116.17 |
| 3.29 | -4 | -7.29 | -3.29 | 4 | 7.29 | 3.29 | -4 |
| -13.44 | -105.18 | -91.74 | 13.44 | 105.18 | 91.74 | -13.44 | -105.18 |
| 2.11 | -21.43 | -23.54 | -2.11 | 21.43 | 23.54 | 2.11 | -21.43 |
| 33.67 | 21.85 | -11.82 | -33.67 | -21.85 | 11.82 | 33.67 | 21.85 |
| -202.27 | 29.2 | 231.47 | 202.27 | -29.2 | -231.47 | -202.27 | 29.2 |
| -45.1 | 40.42 | 85.52 | 45.1 | -40.42 | -85.52 | -45.1 | 40.42 |
| 58.35 | -118.57 | -176.92 | -58.35 | 118.57 | 176.92 | 58.35 | -118.57 |
| 189 | 42 | -147 | -189 | -42 | 147 | 189 | 42 |
| 9.2 | -93.33 | -102.53 | -9.2 | 93.33 | 102.53 | 9.2 | -93.33 |
| 185.37 | -58.44 | -243.81 | -185.37 | 58.44 | 243.81 | 185.37 | -58.44 |
| 90.67 | 18.25 | -72.42 | -90.67 | -18.25 | 72.42 | 90.67 | 18.25 |
| -32 | -268.17 | -236.17 | 32 | 268.17 | 236.17 | -32 | -268.17 |
| 1.77 | 232.5 | 230.73 | -1.77 | -232.5 | -230.73 | 1.77 | 232.5 |
| -46.2 | -2.67 | 43.53 | 46.2 | 2.67 | -43.53 | -46.2 | -2.67 |
| -179.67 | -183.8 | -4.13 | 179.67 | 183.8 | 4.13 | -179.67 | -183.8 |
| -25.9 | 7.6 | 33.5 | 25.9 | -7.6 | -33.5 | -25.9 | 7.6 |
| 59.11 | 8.24 | -50.87 | -59.11 | -8.24 | 50.87 | 59.11 | 8.24 |
| 54 | -2.06 | -56.06 | -54 | 2.06 | 56.06 | 54 | -2.06 |
| 12.95 | -22.75 | -35.7 | -12.95 | 22.75 | 35.7 | 12.95 | -22.75 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 11.1 | -141.29 | -152.39 | -11.1 | 141.29 | 152.39 | 11.1 | -141.29 |
| -16.28 | 44.19 | 60.47 | 16.28 | -44.19 | -60.47 | -16.28 | 44.19 |
| -115.5 | -11.39 | 104.11 | 115.5 | 11.39 | -104.11 | -115.5 | -11.39 |
| -50.33 | 56.5 | 106.83 | 50.33 | -56.5 | -106.83 | -50.33 | 56.5 |
| 14.24 | 92.5 | 78.26 | -14.24 | -92.5 | -78.26 | 14.24 | 92.5 |
| -14 | -45 | -31 | 14 | 45 | 31 | -14 | -45 |
| 121.7 | -84.53 | -206.23 | -121.7 | 84.53 | 206.23 | 121.7 | -84.53 |
| 83.55 | -54.5 | -138.05 | -83.55 | 54.5 | 138.05 | 83.55 | -54.5 |
| 53.94 | -21 | -74.94 | -53.94 | 21 | 74.94 | 53.94 | -21 |
| 78.6 | 14.4 | -64.2 | -78.6 | -14.4 | 64.2 | 78.6 | 14.4 |
| -128.67 | 16.6 | 145.27 | 128.67 | -16.6 | -145.27 | -128.67 | 16.6 |
| -1.83 | -77.14 | -75.31 | 1.83 | 77.14 | 75.31 | -1.83 | -77.14 |
| -14.38 | -40.17 | -25.79 | 14.38 | 40.17 | 25.79 | -14.38 | -40.17 |
| 74.6 | 57.38 | -17.22 | -74.6 | -57.38 | 17.22 | 74.6 | 57.38 |
| 41.62 | -20.66 | -62.28 | -41.62 | 20.66 | 62.28 | 41.62 | -20.66 |
| -47.27 | -132.5 | -85.23 | 47.27 | 132.5 | 85.23 | -47.27 | -132.5 |
| 57.62 | -192.33 | -249.95 | -57.62 | 192.33 | 249.95 | 57.62 | -192.33 |
| 53.61 | -5.87 | -59.48 | -53.61 | 5.87 | 59.48 | 53.61 | -5.87 |
| -40.29 | 56.53 | 96.82 | 40.29 | -56.53 | -96.82 | -40.29 | 56.53 |
| 67.12 | -17 | -84.12 | -67.12 | 17 | 84.12 | 67.12 | -17 |
| -51.08 | -13.21 | 37.87 | 51.08 | 13.21 | -37.87 | -51.08 | -13.21 |
| 103.67 | -83.67 | -187.34 | -103.67 | 83.67 | 187.34 | 103.67 | -83.67 |
| 4.3 | -242 | -246.3 | -4.3 | 242 | 246.3 | 4.3 | -242 |
| 33.4 | 216.08 | 182.68 | -33.4 | -216.08 | -182.68 | 33.4 | 216.08 |
| 54.22 | -136.35 | -190.57 | -54.22 | 136.35 | 190.57 | 54.22 | -136.35 |
| -115.5 | -11.39 | 104.11 | 115.5 | 11.39 | -104.11 | -115.5 | -11.39 |
| -3.45 | 9.11 | 12.56 | 3.45 | -9.11 | -12.56 | -3.45 | 9.11 |
| -47.04 | 38.4 | 85.44 | 47.04 | -38.4 | -85.44 | -47.04 | 38.4 |
| -31.45 | 37 | 68.45 | 31.45 | -37 | -68.45 | -31.45 | 37 |
| -54.96 | 110.2 | 165.16 | 54.96 | -110.2 | -165.16 | -54.96 | 110.2 |
| 125.77 | 162.33 | 36.56 | -125.77 | -162.33 | -36.56 | 125.77 | 162.33 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| **34.49** | 11 | -23.49 | -34.49 | -11 | 23.49 | 34.49 | 11 |
| **27.08** | -26.71 | -53.79 | -27.08 | 26.71 | 53.79 | 27.08 | -26.71 |
| **-21.1** | 35.43 | 56.53 | 21.1 | -35.43 | -56.53 | -21.1 | 35.43 |
| **1** | 106.14 | 105.14 | -1 | -106.14 | -105.14 | 1 | 106.14 |
| **52.43** | 70.4 | 17.97 | -52.43 | -70.4 | -17.97 | 52.43 | 70.4 |
| **51.53** | 37.15 | -14.38 | -51.53 | -37.15 | 14.38 | 51.53 | 37.15 |
| **-56.02** | -17.4 | 38.62 | 56.02 | 17.4 | -38.62 | -56.02 | -17.4 |
| **-124** | 1.14 | 125.14 | 124 | -1.14 | -125.14 | -124 | 1.14 |
| **-2.6** | 3.29 | 5.89 | 2.6 | -3.29 | -5.89 | -2.6 | 3.29 |
| **-101.33** | -32.3 | 69.03 | 101.33 | 32.3 | -69.03 | -101.33 | -32.3 |
| **-6.82** | 65.4 | 72.22 | 6.82 | -65.4 | -72.22 | -6.82 | 65.4 |
| **16.51** | 160 | 143.49 | -16.51 | -160 | -143.49 | 16.51 | 160 |
| **-8.35** | 28.43 | 36.78 | 8.35 | -28.43 | -36.78 | -8.35 | 28.43 |
| **-36.09** | 17.57 | 53.66 | 36.09 | -17.57 | -53.66 | -36.09 | 17.57 |
| **12.95** | 41.87 | 28.92 | -12.95 | -41.87 | -28.92 | 12.95 | 41.87 |
| **-213.78** | -310.83 | -97.05 | 213.78 | 310.83 | 97.05 | -213.78 | -310.83 |
| **24.7** | 7.31 | -17.39 | -24.7 | -7.31 | 17.39 | 24.7 | 7.31 |
| **183.7** | 170.43 | -13.27 | -183.7 | -170.43 | 13.27 | 183.7 | 170.43 |
| **-64.16** | -198.37 | -134.21 | 64.16 | 198.37 | 134.21 | -64.16 | -198.37 |
| **185.18** | 266.74 | 81.56 | -185.18 | -266.74 | -81.56 | 185.18 | 266.74 |
| **161.8** | -160.31 | -322.11 | -161.8 | 160.31 | 322.11 | 161.8 | -160.31 |
| **106.9** | 54.95 | -51.95 | -106.9 | -54.95 | 51.95 | 106.9 | 54.95 |

**Appendix 22: High Tamil Vocabulary Group**

| Participant | FFWord1 Cong | FFWord1 Incong | SFWord1 Cong | SFWord1 Incong | GazeWord1 Cong | GazeWord1 Incong | TTWord1 Cong | TTWord1 Incong | FFWord2 Cong | FFWord2 Incong | SFWord2 Cong | SFWord2 Incong | GazeWord2 Cong |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 47 | 194.71 | 261.83 | 249 | 232 | 194.71 | 261.83 | 290.38 | 389.7 | 260.7 | 263.46 | 194.71 | 261.83 | 380.6 |
| 30 | 279 | 163.5 | 199.14 | 187.75 | 279 | 163.5 | 400 | 363.4 | 199.14 | 188.8 | 279 | 163.5 | 199.14 |
| 17 | 342 | 351.2 | 374 | 280.67 | 342 | 401.2 | 342 | 443.67 | 369.43 | 387.71 | 342 | 330 | 397.71 |
| 36 | 208.5 | 262.44 | 243.67 | 222.67 | 208.5 | 262.44 | 437.8 | 407.23 | 239 | 222 | 196.67 | 264.25 | 379.71 |
| 56 | 246.45 | 215 | 291.17 | 328.17 | 255.27 | 292.75 | 362.15 | 475.67 | 248.2 | 292.3 | 232.11 | 256 | 330.67 |
| 66 | 360 | 236 | 222.5 | 223.64 | 360 | 236 | 236.67 | 307.38 | 182.67 | 213.85 | 360 | 265.29 | 182.67 |
| 25 | 158.22 | 217.33 | 190.33 | 198.57 | 158.22 | 250.67 | 307.33 | 364.44 | 190 | 187 | 147.71 | 265.5 | 260 |
| 8 | 282.4 | 396.86 | 224.33 | 340.5 | 282.4 | 430.57 | 292.67 | 455.09 | 230.18 | 322.14 | 282.4 | 426.8 | 279.82 |
| 55 | 347.18 | 300.14 | 290.6 | 329 | 463.36 | 367.86 | 543.09 | 656.89 | 289 | 241.67 | 331.57 | 339.6 | 526.73 |
| 5 | 212 | 274.6 | 294.67 | 255.33 | 212 | 292.4 | 446 | 862.5 | 246.33 | 260.83 | 207.33 | 319.33 | 355.33 |
| 12 | 249.33 | 283 | 261.75 | 283.6 | 284 | 301.4 | 284 | 301.4 | 256.5 | 275.27 | 247.8 | 275.11 | 371 |
| 46 | 235.43 | 195.14 | 219.17 | 275.7 | 305.57 | 292.64 | 571.6 | 466.93 | 210.93 | 267.93 | 231.63 | 182.86 | 391.07 |
| 53 | 279 | 163.5 | 199.14 | 187.75 | 279 | 163.5 | 400 | 363.4 | 199.14 | 188.8 | 279 | 163.5 | 199.14 |
| 14 | 270.6 | 225.5 | 216.25 | 256.67 | 286.6 | 287.75 | 387.85 | 309.08 | 206.67 | 245.56 | 299.75 | 225.71 | 256.67 |
| 29 | 340.83 | 324.55 | 426.67 | 470.86 | 457.83 | 432 | 592.62 | 513.08 | 406.46 | 488.71 | 376.75 | 358.29 | 520.46 |
| 1 | 161.25 | 188.33 | 203 | 176.29 | 161.25 | 188.33 | 216 | 266.5 | 203 | 196.9 | 161.25 | 188.33 | 203 |

## Appendix 23: Low Tamil Vocabulary Group

| Participant | FFWord1 Cong | FFWord1 Incong | SFWord1 Cong | SFWord1 Incong | GazeWord1 Cong | GazeWord1 Incong | TTWord1 Cong | TTWord1 Incong | FFWord2 Cong | FFWord2 Incong | SFWord2 Cong | SFWord2 Incong | GazeWord2 Cong |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 24 | 284.57 | 258.67 | 198.4 | 206 | 309.71 | 258.67 | 788.5 | 684.67 | 203.71 | 196 | 299 | 258.67 | 262.86 |
| 20 | 306 | 274 | 527.5 | 259.33 | 507.56 | 298.75 | 647.71 | 531.17 | 327.43 | 270.73 | 310 | 302.67 | 524.86 |
| 23 | 387.67 | 208 | 367 | 183.2 | 482 | 270.57 | 612.75 | 378.4 | 272.86 | 245.82 | 479 | 227.5 | 502.57 |
| 63 | 236 | 288.43 | 290.6 | 361 | 285.43 | 378.43 | 353.2 | 538.88 | 308.77 | 277.88 | 291.88 | 328 | 454.31 |
| 4 | 188 | 211.71 | 206.29 | 206.4 | 204.89 | 262 | 384.67 | 487.45 | 218 | 204.33 | 191 | 185.6 | 305.64 |
| 6 | 300.67 | 276.29 | 396.5 | 262.86 | 300.67 | 276.29 | 692.77 | 845.64 | 352.73 | 246.91 | 300.67 | 276.29 | 396.91 |
| 11 | 260.44 | 262.55 | 383.43 | 362 | 260.44 | 262.55 | 424.62 | 436.62 | 358 | 291.86 | 273.75 | 257.25 | 420.73 |
| 33 | 252.31 | 238.31 | 284 | 239 | 285.69 | 250.77 | 659.2 | 461.38 | 262.4 | 235.43 | 250.2 | 265.6 | 358.67 |
| 7 | 431.64 | 262.64 | 591.29 | 461.38 | 431.64 | 261.3 | 442.58 | 387.58 | 511.07 | 420.27 | 431.64 | 266.67 | 643.64 |
| 65 | 285.88 | 229.86 | 244.4 | 227 | 332.13 | 229.86 | 554.67 | 534.82 | 234.58 | 230.44 | 338 | 223 | 379.08 |
| 9 | 267.71 | 271 | 297.2 | 293.2 | 267.71 | 271 | 441.45 | 307.33 | 272.91 | 239.09 | 284.67 | 271 | 376.36 |
| 27 | 253.45 | 266.4 | 359 | 336.25 | 375.82 | 417.6 | 564.91 | 655 | 341.6 | 293.86 | 271.43 | 248.4 | 601.4 |
| 21 | 425.67 | 297 | 340 | 356.6 | 471.11 | 455.11 | 438.77 | 509.23 | 203.07 | 324.15 | 228.33 | 209.71 | 321.07 |
| 26 | 217.83 | 219.6 | 186 | 418.5 | 300.5 | 275 | 462.33 | 472.46 | 232.77 | 260.91 | 248.36 | 316.44 | 343.85 |
| 48 | 248.36 | 302.36 | 245.2 | 243.14 | 248.36 | 327.64 | 304.17 | 495.27 | 309.17 | 289.5 | 398.43 | 354.25 | 389.58 |
| 38 | 364.75 | 313.67 | 285.38 | 272.17 | 387.38 | 377.83 | 519.22 | 382.33 | 303.64 | 298 | 405.71 | 289.4 | 571.64 |