

Prepositional phrase attachment  
ambiguities in declarative and  
interrogative contexts: Oral reading  
data

*Tyler J. Peckenpaugh*

*2019-08-25*

# Contents

<b>Abstract</b>	<b>vii</b>
<b>Acknowledgements</b>	<b>viii</b>
About this draft . . . . .	viii
<b>1 Introduction and background</b>	<b>1</b>
1.1 Motivations for the current study . . . . .	4
1.2 Structural overview of the ambiguity relevant to this study . . . . .	7
1.3 Interrogativity . . . . .	14
1.4 Prosody of questions vs. declaratives . . . . .	17
1.5 Evidence that prosody can affect syntactic parsing . . . . .	19
1.6 Predictions for the current study . . . . .	23
<b>2 Methodology</b>	<b>26</b>
2.1 Materials . . . . .	27
2.2 Participants recruitment . . . . .	35
2.3 Location . . . . .	35
2.4 Equipment and software . . . . .	36
2.5 Versions of the experiment . . . . .	36
2.6 Procedure . . . . .	37
2.7 Measurements of utterance timing . . . . .	44

<i>CONTENTS</i>	ii
2.8 Prosodic judgments . . . . .	47
<b>3 Results and discussion</b>	<b>51</b>
3.1 Data for analysis . . . . .	52
3.2 Prosodic break patterns . . . . .	53
3.3 Discussion of prosodic break patterns . . . . .	65
3.4 Inter-reading time . . . . .	69
3.5 Discussion of IRT results . . . . .	72
<b>4 General discussion</b>	<b>76</b>
4.1 Behavioral correlate for the 2016 intuition and the current hypothesis?	77
4.2 On possible explanations for the intuition . . . . .	77
4.3 Conclusions . . . . .	82
4.4 Summary . . . . .	86
<b>A Experimental items</b>	<b>87</b>
<b>B Filler items</b>	<b>90</b>
<b>C Recruitment notice</b>	<b>94</b>
<b>D Instructions to participants</b>	<b>95</b>
<b>E Instructions for prosodic coding</b>	<b>97</b>
<b>References</b>	<b>99</b>

# List of Tables

2.1	Illustrative experimental item, constructed in four versions. . . . .	27
2.2	Illustrative filler items, constructed in two versions. . . . .	34
2.3	Distance in characters from fixation to disambiguation of experimental items for the current study. . . . .	43
2.4	EVS-adjusted character distance to disambiguation in experimental items. . . . .	44
2.5	Percent agreement between the original ratings and the second rater (inter-rater) or the second rating by the original rater (intra-rater). .	49
3.1	Number of participants per version-order combination. . . . .	52
3.2	Number of recordings analyzed, as a function of Speech Act and PP2 Status. . . . .	53
3.3	Percent occurrence of OBJ break (frequency of occurrence in parenthesis) as a function of sentence type and Reading. . . . .	54
3.4	Percent occurrence of PP1 break (frequency of occurrence in parenthesis) as a function of sentence type and Reading. . . . .	54
3.5	Percent occurrence of both breaks as a function of sentence type and Reading. . . . .	55
3.6	Mixed effects logistic regression model predicting OBJ break occurrence (FULL). . . . .	59

3.7	Mixed effects logistic regression model predicting OBJ break occurrence (REDUCED). . . . .	59
3.8	Mixed effects logistic regression model predicting PP1 break occurrence (FULL). . . . .	60
3.9	Mixed effects logistic regression model predicting PP1 break occurrence (REDUCED). . . . .	61
3.10	Mixed effects logistic regression model predicting OBJ break dominance (FULL). . . . .	61
3.11	Mixed effects logistic regression model predicting OBJ break dominance (REDUCED). . . . .	62
3.12	Mixed effects logistic regression model predicting PP1 break dominance (REDUCED). . . . .	62
3.13	Mixed effects logistic regression models predicting break dominance in Reading 2 (REDUCED). . . . .	63
3.14	Linear mixed effects regression model predicting wIRT by sentence type with interaction term (FULL). . . . .	71
3.15	Linear mixed effects regression model predicting wIRT by sentence type (REDUCED). . . . .	71
3.16	Linear mixed effects regression model predicting wIRT by Speech Act and PP2 head. . . . .	72
4.1	Percent occurrence of break patterns in Reading 2 as a function of sentence type. . . . .	79
A.1	Experimental items in four versions . . . . .	87
B.1	Filler items with trailing PPs . . . . .	90
B.2	Filler items with no trailing PPs . . . . .	93

D.1 Table of keyboard mappings . . . . .	96
--	----

# List of Figures

1.1	Illustrative syntactic tree of a ternary-branching VP. . . . .	9
1.2	Syntactic tree of an illustrative example sentence with an ambiguous PP1 and a modifier-PP2 (Mod). . . . .	10
1.3	Syntactic tree of an illustrative example sentence with an ambiguous PP1 and an argument-PP2 (Arg). . . . .	11
1.4	Illustrative syntactic tree of the basic configuration for Mod cases. . .	23
1.5	Illustrative syntactic tree of the basic configuration for Arg cases. . .	23
2.1	Diagram of 4-screen sequence presented for each item, showing the key presses triggering movement between successive screens. . . . .	38
3.1	Break pattern as a function of sentence type and Reading. . . . .	56
3.2	Percent break dominance occurrence as a function of sentence type and Reading. . . . .	57
3.3	Distributions of R1 delay and R2 delay . . . . .	64
3.4	Plot of pattern proportions as a function of sentence type. . . . .	65
3.5	Distribution of wIRT. . . . .	70
3.6	Mean IRT as a function of sentence type. . . . .	70

# Abstract

*Abstract is a work in progress.* This paper reports a study on the effect of interrogativity on the oral reading of temporarily ambiguous prepositional phrases (PPs). Specifically, it looks at sentences ending in a of two PPs, where the first is interpretable as the goal argument of the preceding verb, and the status of the second (PP2 Status) is manipulated to either necessarily be the goal argument of that verb (Arg), forcing reanalysis, or not (Mod), allowing the original parse to stand. No evidence is found that interrogativity impacts the difficulty of understanding the Arg-type sentences, despite an intuitive decrease in difficulty when those sentences are presented in an interrogative context. A double-reading protocol is employed, where participants are asked to read a sentence first without preview (Reading 1), and then after unlimited preview (Reading 2). A robust effect of PP2 Status is found for the prosodic phrasing of the target sentences, and an effect of interrogativity on the study time between Readings, Inter-Reading Time (IRT), is reported.



# Acknowledgements

TBD

## About this draft

This represents the document that will be defended on August 29th. The goal is to deposit before September 16, after incorporating whatever revisions are requested.

# Chapter 1

## Introduction and background

added text~~deleted text~~

This paper presents a study on human sentence processing, or parsing, and on the parsing of a particular sort of ambiguity. Parsing is assumed to be the projection of structure by a reader or listener over a string of words (which ~~obviously~~ lacks inherent structure). Following the ~~sort of~~ models of parsing developed~~put forward~~ by, e.g., Kimball (1973), Frazier & Fodor (1978), and Frazier & Clifton (1996), this study assumes that parsing is done online, (i.e., during listening to or reading the word string) with the aim~~and~~ that most material is~~must be~~ incorporated into the structure being built as soon as it is encountered. This can lead to mis-parses, where the parser has guessed wrong about how to incorporate a ~~given phrase with a~~ temporarily ambiguous phrase in the input, which becomes apparent when ~~subsequent structure, and then encounters~~ material is encountered which~~that~~ cannot be incorporated into the resulting~~current~~ structure. This sort of parser~~crash~~ is called a "garden path." When ~~it~~this happens, the parser must reanalyze the material that had so far been processed, in order to arrive at a~~come up a~~ structure that can accommodate both the new and the old material grammatically. ~~This sort~~

[1]: example of what markup means; by default changes are based on discussion–<sup>superscript</sup> initials "TJP"<sup>superscript</sup> indicate changes that have not been discussed

of parser crash is called a garden path.

In short: "Garden path" effects occur when a temporarily ambiguous sentence resolves in such a way that the structure initially preferred by the parser is incompatible with how the sentence actually continues. These parsing errors have traditionally been attributed to structurally-driven focused parsing preferences (Frazier, 1979; Frazier & Fodor, 1978; Kimball, 1973) which ~~that~~ ignore semantic content on the first pass. Frazier (1979) formulates several of these structural preferences, including the following two which are widely accepted in one form or another:

- (1) *Minimal attachment* Attach incoming material into the phrase-marker being constructed using the fewest nodes consistent with the well-formedness rules of the language under analysis (Frazier, 1979, p. 24)
- (2) *Late closure* When possible, attach incoming material into the clause currently being parsed (Frazier, 1979, p. 20)

Because these strategies ignore semantic and pragmatic plausibility and the parser typically does not know what material might occur ~~be~~ further on in the word string, mis-parses at temporarily ambiguous regions can occur, ~~resulting in garden paths~~. Minimal Attachment (MA) is important to the present ~~this~~ study and will be revisited later on ~~in Section 1.2~~.

An example is the commonly studied garden path sentence, "The horse raced past the barn fell" (Bever, 1970). Here, the initial parse incorrectly assumes that the matrix subject is the unmodified NP *the horse*, per Minimal Attachment, and takes the matrix verb to be *raced*, as in the sentence, *The horse raced past the finish line*.

- (3) The horse raced past the barn fell. (Bever, 1970)

a) [<sub>S</sub> [<sub>NP</sub> The horse] [<sub>VP</sub> raced past the barn]]    ???    [<sub>VP</sub> fell]

b) [<sub>S</sub> [<sub>NP</sub> The horse raced past the barn] [<sub>VP</sub> fell]]

An attempted parse resulting in structure (3 a) crashes, as it is not possible to incorporate the final word *fell* in a grammatical way. Reanalysis is required, with the grammatical parse being (3 b) where the matrix subject is *the horse raced past the barn*, a noun phrase (NP) in which *the horse* is modified by *containing* a reduced relative clause *raced past the barn*. ~~Then~~Thus *fell* can be incorporated as the matrix verb of the sentence, with a structure comparable to, *The horse (that was) raced past the barn was hungry*.

The study ~~being reported~~ in this dissertation is concerned with ~~certain~~ sentences that contain a temporarily ambiguous prepositional phrase (PP1), followed by another (PP2) which causes the initial~~expected~~ parse to crash. Specifically, it is expected that PP1 in an example such as (4) will initially be interpreted as the goal of *cram*, but that parse will fail when it is realized that PP2 (in his briefcase) cannot plausibly modify drawer. Instead, PP1 will have to modify *paperwork* so that PP2 can be the goal argument of *cram*.<sup>TJP</sup>

(4) He had planned to cram the paperwork [<sub>PP1</sub> in the drawer] [<sub>PP2</sub> in his briefcase].

(5) He had planned to cram the paperwork [<sub>PP1</sub> in the drawer [<sub>PP2</sub> of his filing cabinet]].

This contrasts with the ~~similar~~ sentence in (5), where PP2 can plausibly modify *drawer* and so the parse where the PP1 in the drawer is the goal argument is acceptable,ed and *of his filing cabinet* is incorporated as a modifier within itPP1.

## 1.1 Motivations for the current study

The current study was initially motivated by ~~a phenomenon~~~~an observation~~ discovered by Janet Dean Fodor and Dianne Bradley, and originally reported in Peckenpaugh (2016). That observation ~~was~~ that a garden path sentence like (4) repeated here as (6) is, for whatever reason, not as difficult to process when presented as an interrogative, as in (7), rather than ~~as~~ a declarative.

(6) He had planned to cram the paperwork [~~PP1~~ in the drawer][~~PP2~~ in his briefcase].

(7) Had he planned to cram the paperwork [~~PP1~~ in the drawer][~~PP2~~ in his briefcase]?

Peckenpaugh (2016) attempted to find a behavioral correlate of this intuition by ~~examining~~~~looking at~~ variation in reading time ~~for~~~~of~~ sentences like (6) and (7), but the results were inconclusive. The current study continues that line of research by ~~testing~~~~looking at~~ similar sentences, while also attempting to control certain ~~additional~~ factors that may have led to Peckenpaugh (2016)'s inconclusive results. One such concern is that (6) relies on the practical implausibility of a drawer within a briefcase, i.e., on real world knowledge, in order to disambiguated the appropriate attachment sites for PP1 and PP2. Real world knowledge and beliefs ~~about~~~~of~~ what is or is not plausible likely ~~varyies~~ between speakers and so may not always be ~~effective~~~~reliable~~ as a trigger for reanalysis. The current study makes use of carefully constructed sentences (the criteria by which they were constructed ~~are~~~~is~~ detailed in Section 2.1) which do not rely on plausibility or pragmatics to disambiguated the PP attachment sites, but instead make use of what will be referred to as<sup>TJP</sup> syntactic disambiguation by including a PP2 that cannot grammatically be incorporated into PP1. This is illustrated, as in (8) and (9). Syntactic disambiguation is the term chosen to refer to the sort of disambiguation used for the current study, although

various terms could be used, e.g., lexical or semantic disambiguation; it hinges on the lexical identity of the preposition that heads PP<sub>2</sub>, i.e., *into* instead of *in* and what syntactic position or semantic role such a PP can be assigned.<sup>TJP</sup>

(8) He had planned to cram the paperwork [<sub>PP<sub>1</sub></sub> in the drawer] [<sub>PP<sub>2</sub></sub> into his briefcase].

(9) Had he planned to cram the paperwork [<sub>PP<sub>1</sub></sub> in the drawer] [<sub>PP<sub>2</sub></sub> into his briefcase]?

**[2]:** added interrogative example

The change from pragmatically disambiguation in (6) and (7) to syntactic disambiguation in (8) and (9) creates an important distinction between the hypothesis tested~~original observation~~ in Peckenpaugh (2016) and the hypothesis to be tested ~~in what~~ the current study ~~is considering~~. It cannot be assumed that the intuition about (6) ~~and~~ vs. (7) necessarily extends to cases like (8) ~~and~~ (9), and so one of the questions the current study ~~addresses~~ asks is whether the greater ease of processing a question than a declarative~~intuition~~ can be shown to extend to syntactically disambiguated cases that are similar to the pragmatically disambiguated cases for which the observation was first made. In order to keep a clear focus on this ~~small~~ fine but potentially important difference, a convention will be adopted throughout this document where the intuited amelioration of the garden path effect in pragmatically disambiguated cases is referred to as the 2016 Intuition, while the possibility that the intuition extends to syntactically disambiguated cases is referred to as the Current Hypothesis.

(10) *The 2016 Intuition:* Certain pragmatically disambiguated prepositional phrase (PP) attachment ambiguities which are difficult to parse in the declarative are less difficult to parse when presented as yes-no interrogatives (e.g., *The nanny sat the cranky little boy on the stroller on the swing*, vs., *Did the nanny seat the cranky little boy on the stroller on the swing?*).

**[3]:** examples updated and pairs added

- (11) *The Current Hypothesis:* The 2016 Intuition may be extensible to PP attachment ambiguities that are syntactically disambiguated in addition to those that are pragmatically disambiguated (e.g., *He had planned to cram the paperwork in the drawer into his briefcase*, vs., *Had he planned to cram the paperwork in the drawer into his briefcase?*).

In addition to exploring the possible extensibility of the 2016 intuition, the current study ~~addresses~~ is interested in what property or properties of polar questions might lead to an easier parsing of garden paths, or at least to the perception that they are easier to parse when compared to declaratives. Because there are minimal differences between the polar question and ~~declarative version of~~ declarative versions of these sentences, it ~~seems likely~~ is fairly easy to assume that the cause lies in one of two domains: either the prosodic changes triggered by the use of question intonation, or the pragmatic and semantic properties that are not shared across the versions.

An obvious reflex of the former possibility is another question: how are the various versions of these sentences actually pronounced, prosodically? The reported study seeks to answer this, but while the recordings collected provide some insight, more work is likely needed to ~~satisfactorily~~ provide a fully satisfactory answer.

Likewise, the latter possibility raises the question: exactly ~~leads one to wonder:~~ what are the semantic and pragmatic differences between a polar question and its declarative counterpart? This question has been ~~can be~~ approached in the ~~a more~~ theoretical way, and ~~has been to some degree in the literature~~ (see, e.g. Fiengo (2007) ). Nevertheless ~~That said,~~ it remains to be determined how those properties could lead to an easier or more difficult parsing process, and ~~whether or not a~~ satisfactory explanation for the intuition at large can be pulled from these differences. }

## 1.2 Structural overview of the ambiguity relevant to this study

The 2016 intuition and the current hypothesis are both concerned with a temporarily ambiguous sequence of PPs at the end of a sentence. This section will discuss what the possible attachment sites for those PPs are and which structures ultimately do and do not prove viable work.

The example in (12) shows a pragmatically disambiguated sentence with an argument-PP2. The initial parse (a) is implausible, resulting in structural reanalysis to the preferable parse in (b):

(12) Jed crammed the newspapers under the sofa in the wastebasket.

- a) # ... [<sub>VP</sub> crammed [<sub>NP</sub> the newspapers] [<sub>PP1</sub> under [<sub>NP</sub> the sofa [<sub>PP2</sub> in the wastebasket]]]]\
- b) ✓ ... [<sub>VP</sub> crammed [<sub>NP</sub> the newspapers [<sub>PP1</sub> under [<sub>NP</sub> the sofa ] ] [<sub>PP2</sub> in the wastebasket]]] \ “#” indicates a structure with an implausible reading

The initial parse is expected to be (12 a) because of a bias (due to Minimal Attachment, or some variation thereof, see (1) above), which favors a structure where the first PP attaches into the verb phrase (VP) as an argument of the verb, i.e., [<sub>VP</sub> V NP PP1] with PP1 denoting the goal, which leaves nowhere for the second PP to attach but as a modifier of the noun phrase (NP) inside PP1 ([<sub>PP1</sub> under [<sub>NP</sub> the sofa [<sub>PP2</sub> in the wastebasket]]]). This initial parse (12 a) is pragmatically implausible, since as one does not generally find sofas are generally not found inside wastebaskets. Structural reanalysis is required to bring about the only plausible correct parse (12 b), where PP1 attaches as an NP modifier of the direct object and so allows PP2 to attach as a VP argument, resulting in a structure such as [<sub>VP</sub> V [<sub>NP</sub> N PP1] PP2], i.e., where it is the newspapers under the sofa that are being crammed



into the wastebasket.

For clarity in what follows~~It is important to be clear on how the sentences of~~  
~~concerned are structured. As a matter of terminology,~~ the current study categorizes  
 the sentence types being discussed into two groups based on the status of the PP2  
 they contain: (a) cases where PP2 is an argument of the verb are referred to as  
 Arg-type sentences, and (b) cases where PP2 is a modifier of the preceding noun are  
called Mod-type sentences ; see (13). ~~As noted, it is the~~In practice, Arg sentences  
~~that are generally reported to be~~ garden paths, because PP2 must fill the goal role  
 that PP1 is expected to have filled as just discussed. Mod sentences are not likely to  
result in a garden paths, because PP2 can modify the NP within PP1 toand become  
 part of the goal, and therefore need not disrupt PP1's attachment as from being the  
 goal.

### (13) **PP2 attachment status**

#### (a) *PP2 Argument* (Arg)

He had planned to cram the paperwork [<sub>PP1</sub> in the drawer] [<sub>PP2</sub> into his  
 briefcase].

#### (b) *PP2 Modifier* (Mod)

He had planned to cram the paperwork [<sub>PP1</sub> in the drawer [<sub>PP2</sub> of his  
 filing cabinet]].

In order for the differing parsing process for (13 a) and (13 b) to be explained by a  
 strictly structureally based model of parsing, certain assumptions have to be made  
 about the syntax. A simple way to achieve it~~get the explanation to work~~ is to assume  
 that all arguments of a verb are syntactic sisters to the verb, resulting in a three-way  
 branching VP for ditransitive verbs such as *cram*. In this case, in order to avoid  
 postulating extra nodes that would be required for PP1 to be a modifier of the object  
NP, Minimal Attachment dictates that PP1 should be assumed to fill the verb's goal

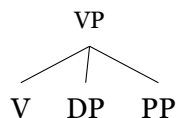


Figure 1.1: Illustrative syntactic tree of a ternary-branching VP.

argument slot. However, ~~t~~This is not how current~~modern~~ syntactic theory assumes the structure to belooks, as three-way branching is proscribed (see e.g., Chomsky (2014)).

The most pronounced structural difference between the two structures is that in the sentence with a modifier-PP2 (Mod) the major disjuncture comes fairly early, just after the object NP *the newspapers*. For the sentences with argument-PP2s (Arg), the major disjuncture is later, just after PP1 (*under the sofa*).

It bears mentioning that Minimal Attachment as defined by Frazier (1979) is somewhat at odds with recent developments in syntactic theory, e.g., obligatory binary branching (cf. Chomsky, 2014, p. 62). As originally postulated, Minimal Attachment relies on a verb with multiple internal arguments incorporating<sup>TJP</sup> each of those arguments as a sister (i.e., a ternary branching structure as in 1.1).

Within current theories of syntax, where binary branching is obligatory, two XPs (NP and PP) cannot both be syntactic sisters of the verb, and the structures are assumed to be as shown in the ~~trees presented earlier~~ (Figures 1.2 and 1.3 below<sup>1</sup>), so it becomes less clear that the VP attachment site for PP1 actually creates fewer nodes than the lower NP attachment site, as per Minimal Attachment. Nonetheless, studies by, the preference for VP attachment in these kinds of sentences is there, be it due to Minimal Attachment, a preference for arguments over non-arguments, or something else, as evidenced by experimental data from e.g., Rayner, Carlson, & Frazier (1983), Clifton, Speer, & Abney (1991), and others, show that a preference

<sup>1</sup>Note that when the internal structure of an NP is not relevant (no PP is within it) it is not drawn in these figures, i.e., [NP newspapers] is shorthand for [NP [N' [N newspapers]]].

for VP attachment of a PP in contexts such as the one the current study addresses exists, whether that preference is best explained by Minimal Attachment, or a preference for arguments over non-arguments, or some other mechanism<sup>TJP</sup>.

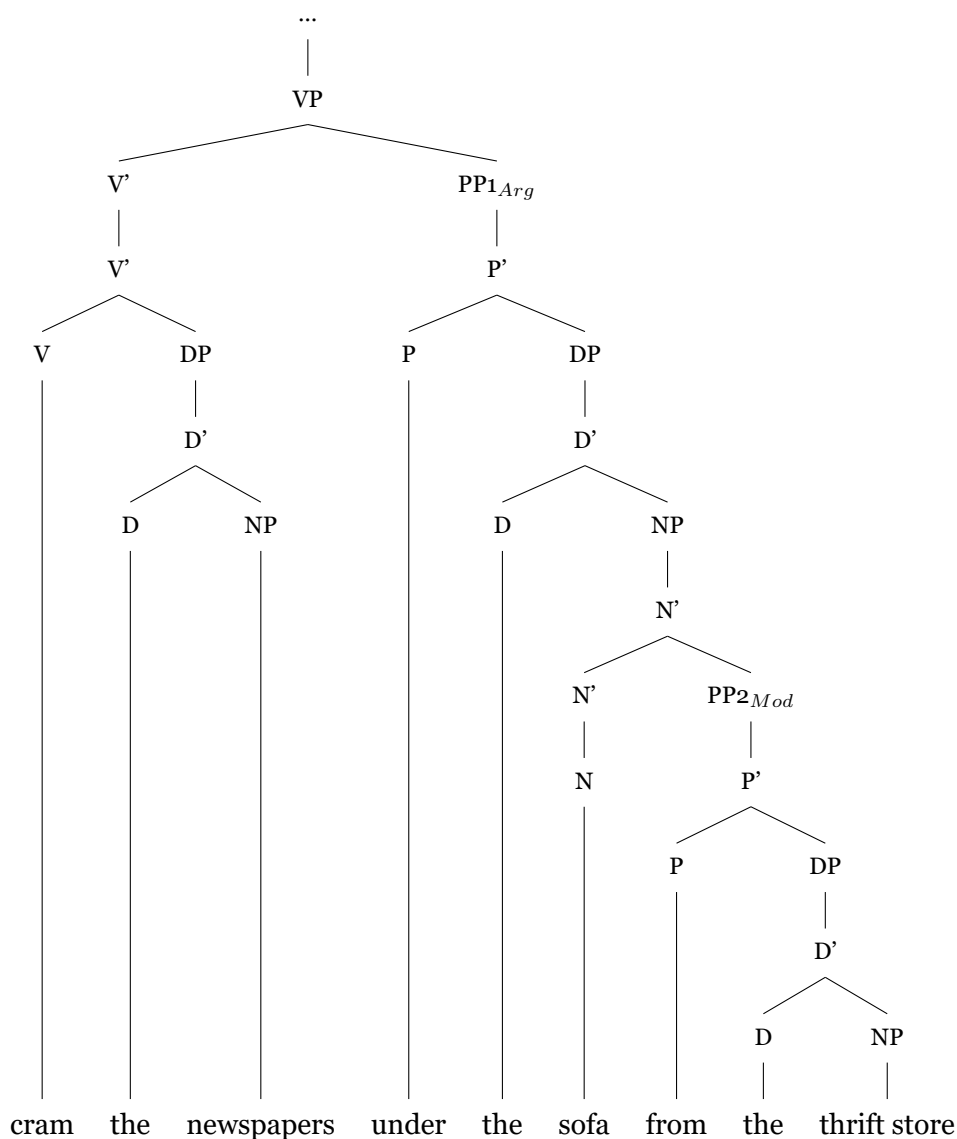


Figure 1.2: Syntactic tree of an illustrative example sentence with an ambiguous PP1 and a modifier-PP2 (Mod).

Setting aside rather than worry about the particularities of syntactic theory aside now, and it also is not necessary to rely on Minimal Attachment, an appeal can be made to a psycho-linguistic distinction added to parsing theory in *Construal* (Frazier & Clifton, 1996): that of primary vs. non-primary relations.

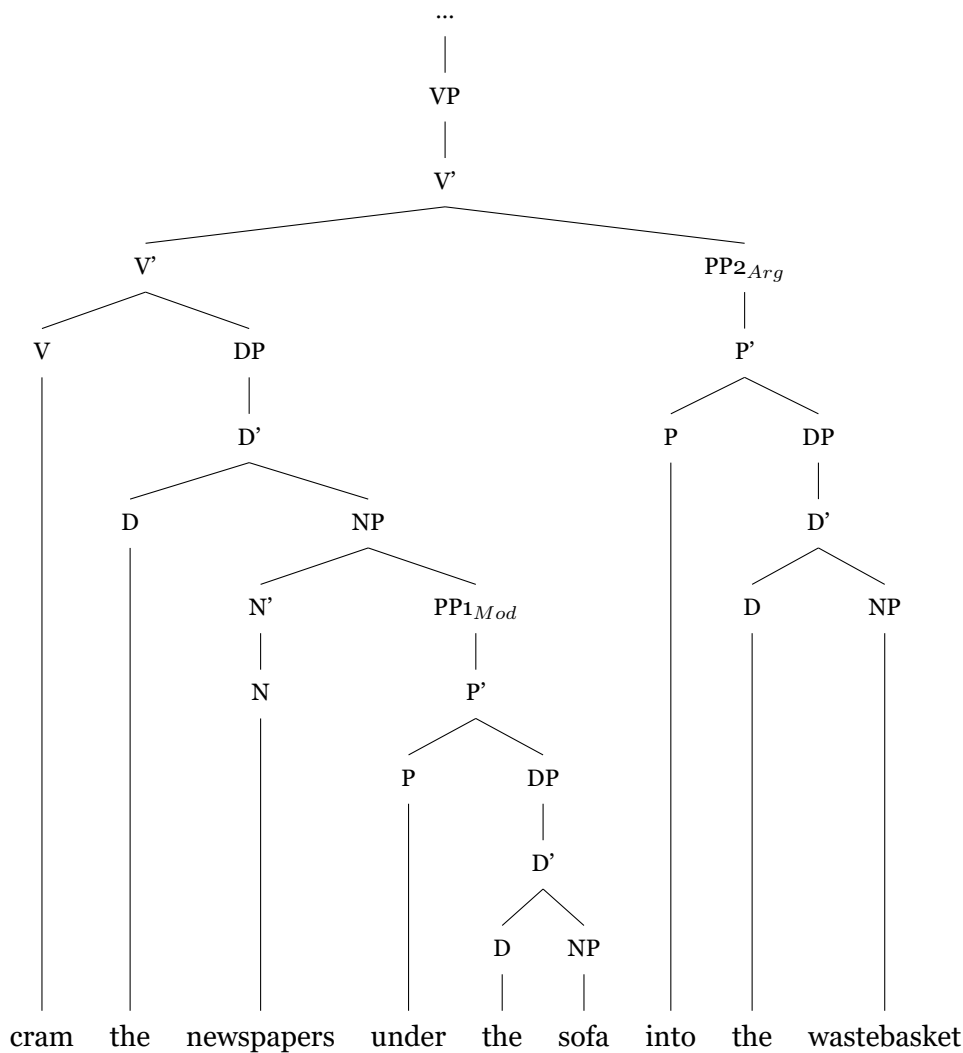


Figure 1.3: Syntactic tree of an illustrative example sentence with an ambiguous PP1 and an argument-PP2 (Arg).

(14) "Primary phrases and relations include

- a) The subject and main predicate of any (+ or -) finite clause
- b) Complements and obligatory constituents of primary phrases" (Frazier & Clifton, 1996, p. 41)

This additional contrast machinery is independently motivated: Frazier and Clifton while ~~structure-first decision making seems to hold for the parsing of many~~ structures, there are some that seem to follow different rules. \*Construal\* illustrates this by way of relative clause (RC) attachment in constructions like (15).

- (15) The journalist interviewed [<sub>NP1</sub> The daughter<sub>i</sub>] of [<sub>NP2</sub> the colonel<sub>j</sub>] [<sub>RC</sub> who<sub>i/j</sub> had an accident] (Frazier & Clifton, 1996, p. 71).

The RC in (15) can modify either NP1, *the daughter*, or NP2 *the colonel*. A ~~structure-first parsing system, together with the widely agreed upon structural parsing strategy~~ Late Closure(see, (2)) predicts a consistent preference for local attachment of the RC in (15), i.e., the structure where the RC modifies NP2. This is because Late Closure dictates that the NP2 phrase not be closed until material that cannot, for one reason or another, be incorporated into that phrase is encountered. Since the RC can modify NP2, Late Closure predicts that it will. Instead, what a number of empirical studies (e.g., Clifton, 1988; Cuetos & Mitchell, 1988) find is a pattern where the preferred structure depends on the relationship between NP1 and NP2. Frazier & Clifton (1996) describe five categories of relationship, and a gradient of preferred RC attachment, from NP1 preference to NP2 preference.

- (16) **RC Attachment by NP1-NP2 relation** (Frazier & Clifton, 1996, p. 31)

a) *Material*

The table of wood [<sub>RC</sub> that was from Galicia]

b) *Quantity*

The glass of wine [<sub>RC</sub> you liked]

c) *Relational* (friend, enemy, son, and other argument taking NPs, e.g., picture-NPs)

The son of the woman [<sub>RC</sub> that was dying]

d) *Possessive*

The car of the company [<sub>RC</sub> that was falling apart]

e) *Non-accompaniment* with

The girl with the hat [<sub>RC</sub> that looked funny]

Frazier & Clifton (1996) cite studies showing that (16 a-b) type configurations favor

NP1 RC attachment, (16 e) type configurations favor NP2 RC attachment, while (16 c-d) are intermediate. They argue that this gradient cannot be readily explained by structural parsing, and instead they make use of a mechanism they call association. RCs, rather than being immediately slotted into a tree in a specific way, are associated with a thematic domain, i.e., the maximal projection of whatever lexical item last assigned theta-roles, together with associated functional projections; in the case of the examples in (16), the last theta assigner is NP2, and its domain extends up to the DP that contains NP1. Association This is a later parsing decision that allows the syntactic structure to be decided after semantic information becomes available: the RC can ultimately modify whichever member of the thematic domain is appropriate.

~~The factor that determines when~~ structural association ~~or vs.~~ structural parsing is appropriate is the distinction ~~between~~ idea of primary vs. non-primary relations. Frazier and Clifton's<sup>TJP</sup> formalization of<sup>TJP</sup> this distinction was previously mentioned<sup>TJP</sup> ~~as in~~ (14) above<sup>TJP</sup>.

RC attachment undergoes association because relative clauses are by definition modifiers, and the relationship between a modifier and whatever is modified is a non-primary relation, ~~and a relative clause is by definition a modifier and not an argument.~~ Returning now ~~Circling back to the PP-attachment~~ construction that this study is concerned with, the argument vs. modifier distinction is precisely what distinguishes the two possible roles ~~statuses~~ of PP2 shown and illustrated in (13), and repeated here as (17).

**[4]:** Moved the formalization from here to above.

#### (17) **PP2 types**

- (a) *PP2 Argument* (Arg) He had planned to cram the paperwork [<sub>PP1</sub> in the drawer] [<sub>PP2</sub> into his briefcase].

- (b) *PP2 Modifier* (Mod) He had planned to cram the paperwork [<sub>PP1</sub> in the drawer [<sub>PP2</sub> of his filing cabinet]].

Without locking down the exact syntactic structures that (@pp2tr) represents, we can nonetheless say that the parser would seek to immediately incorporate PP1 into the tree in both cases. The infinitival (–finite) clause headed by *cram* is a primary phrase, and so its obligatory constituents hold primary relationships with *cram*. As such, association is not available as a mechanism and structural parsing must occur, without semantic information being available. Thus, PP1 is initially interpreted as the goal of *cram*. \*Cram\* takes an obligatory goal argument, so the parser cannot wait for semantic information to inform its association, it must make its best guess based on the principles of structural parsing, and attach it as an argument, as that property is what is forcing the immediate decision to be made. Then, in the case of (17 a), when PP2 is encountered, PP1 must be removed from its syntactic position and be reanalyzed as a modifier of the object NP in order to allow PP2 to fill the goal argument position reanalysis of the PP1 attachment will be required, due to the fact that PP2 must be the goal of \*cram\* and therefore takes the syntactic position that PP1 had been filling. In the case of (17 b), no reanalysis is necessary because PP2 can modify *the drawer* and so does not need to dislodge PP1.

### 1.3 Interrogativity

This section returns to the questions raised by the 2016 Intuition discussed in Section 1.1 about why certain interrogative garden paths of the sort just discussed might appear to be easier to parse in interrogative constructions than in otherwise than similar declarative ones (see examples (18) and (19) ).

- (18) He had planned to cram the paperwork [<sub>PP1</sub> in the drawer] [<sub>PP2</sub> in his briefcase].

**[TJP 1]:** on second thought, I think these need to stay

(19) Had he planned to cram the paperwork [<sub>PP1</sub> in the drawer] [<sub>PP2</sub> in his  
briefcase]?

~~W~~Specifically, what exactly differs between (18) and (19)? Syntactically, very little: the position of the subject *he* and the auxiliary *had* have been reversed, but the words that follow are identical.

The semantic, or perhaps more accurately pragmatic, differences between (18) and (19) lie with the presuppositions the sentences carry with them, which may vary and with the placement of focus, marked phonologically in spoken language. The exact phonetic properties of focus, and the semantic and syntactic consequences of it, are widely studied and debated. For an excellent overview of the topic, see<sup>TJP</sup> Ladd (2008), pp.213-23. The assumptions made about focus in this dissertation are about English only, and are more or less compatible with a Focus-to-Accent model as described by Ladd, which distinguishes "the semantic/pragmatic notion 'focus' from the phonetic/phonological notion 'accent' and – crucially – [...] allows focus to apply to portions of utterances larger than individual words"<sup>TJP</sup> (Ladd, 2008, p. 217). This results in sentential stress being split into two related but independent components: where focus lies within the sentence, and how it is conveyed with regard to the location of accent.

The declarative in (18) has few presuppositions beyond the existence of the actors and objects involved (the referents of *he*, *paperwork*, *drawer* and *briefcase*), and that these actors and objects can be involved in *cramming*. The presuppositions of (19) are different from a super-set of those of (18): for instance, a yes/no question additionally presupposes that the listener does or may know the answer to the question, ~~for one,~~ \added{and also may not presuppose the referent of a focused element. Further presuppositions may might exist, depending on where the focus lies within the sentence.



Focus in a declarative like (18) is typically widebroad<sup>TJP</sup>, i.e., meaning no element is having special<sup>TJP</sup> attention called to it. A polar question like (19), however, will often ~~typically~~ receive narrower<sup>TJP</sup> focus on one word or phrase~~element~~<sup>TJP</sup>, so that when uttered, that one word or phrase~~element~~<sup>TJP</sup> is more prominent than the others. The focused element becomes the part of the sentence that the question is about. Focus can fall on any of the lexical or referential elements in the sentence (subject, auxiliary verb, matrix verb, infinitival verb, object NP, or the NP of either PP1 or PP2) or any part of one of these of the sentence, or the auxiliary verb. Words in all capital letters indicates verbal emphasis signifying focus. ~~Bold facing indicates verbal emphasis in order to locate focus~~

- (20) HAD he planned to cram the paperwork [PP1 in the drawer] [PP2 into his briefcase]?
- (21) Had HE planned to cram the paperwork [PP1 in the drawer] [PP2 into his briefcase]?
- (22) Had he PLANNED to cram the paperwork [PP1 in the drawer] [PP2 into his briefcase]?
- (23) Had he planned to CRAM the paperwork [PP1 in the drawer] [PP2 into his briefcase]?
- (24) Had he planned to cram the PAPERWORK [PP1 in the drawer] [PP2 into his briefcase]?
- (25) Had he planned to cram the paperwork [PP1 in the DRAWER] [PP2 into his briefcase]?
- (26) Had he planned to cram the paperwork [PP1 in the drawer] [PP2 into HIS briefcase]?
- (27) Had he planned to cram the paperwork [PP1 in the drawer] [PP2 into his BRIEFCASE]?

In (20), with focus on the auxiliary, the question encompasses~~is about~~ the entire proposition, and asks whether or not it is true. In this case, there are not any additional presuppositions when compared to the declarative counterpart of the sentence. In (21), with focus on *he*, the question is asking about whether the referent of *he* is the actor who performed the action described; in this case, the entire predicate is presupposed: someone *had planned to cram the paper in the drawer into his briefcase*, but was it *him*? ~~In~~ Skipping ahead to (25), with focus on

*drawer*, the question is instead about which paperwork this is all happening to: the paperwork that is in the drawer, or some other stack of paperwork? In this case, it is presupposed that the referent of *he* was the one who had planned to cram some paperwork into his briefcase, and only the exact referent of *the paperwork* is not presupposed. For each other location of focus, the Presuppositional content is similarly complementary to whichever element is focused and therefore being asked about.

This set of pragmatic differences between (18) and (19) might very well be the source of the 2016 Intuition that (19) is easier to comprehend than (18), but it is not entirely clear why, and it~~that~~ is not the only possibility. Another significant difference between the two ~~sentences~~<sup>Speech Acts</sup> is the ~~rhythm~~<sup>prosody</sup> and intonational melody (which, taken together, make up the prosody)<sup>TJP</sup>. While dialects of English may differ prosodically, there is typically a difference in melody between a declarative and a question, and in many American English dialects, the interrogative is pronounced with a final rise, while the declarative exhibits just a series of down-steps. This difference is the one that the current study explores, to see if it can be shown to correlate with a difference in processing difficulty in syntactically disambiguate declarative and interrogative PP-attachment garden paths, and by extension lend insight into the 2016 Intuition~~readily explain the intuitive difference in processing difficulty.~~

## 1.4 Prosody of questions vs. declaratives

In pursuing the possibility that it is the ~~intonation~~ and prosody of polar interrogatives which creates the 2016 Intuition that motivated this study it is important to~~we must~~ consider the details of~~what~~ question intonation actually sounds like. It is generally agreed that in American English, the intonation of a

polar question has exhibits the property of a final rise. ~~T~~Indeed, this has been confirmed in corpus studies such as Hedberg, Sosa, & Görgülü (2017) who found that 79.8% of the 410 American English yes/no questions in their study (ten-minute phone conversations from the CallHome Corpus of American English and the Fisher English Corpus) had a “low-rise nuclear contour” ( $L^*H-H\%$ ,  $L^*H-\uparrow H\%$ , or  $L^*L-H\%$ )<sup>2</sup>. In their ToBI notation, a tone T is coded either L for low or H for high; T\* is anchored to the stressed syllable, and T- and T% are boundary tones (intermediate phrase boundary and intonational phrase boundary respectively). See, e.g., *Guidelines for ToBI labeling* (Beckman & Ayers, 1997) for a more thorough explanation of ToBI. An additional 10.7% of the Hedberg et al. (2017) data had a “high-rise nuclear contour” (the authors categorizes the following tunes as high-rise nuclear contours:  $H^*H-H\%$ , or  $!H^*L-L\%$ ). That leaves only 9.5% spread across the other 5 categories (High-fall, Rise-fall, Low-fall, Fall-rise, and Level). Only 5.6% of the total data showed any sort of falling contour. According to the authors’ analyses, final high-rise these contours occur on the final main stress of a sentence and thereafter. In the case of the types of sentences examined in the current study, that would result in a rising contour on the head noun of the final PP as in (28), regardless of whether the final PP is a Mod or an Arg.

(28) Did Jed cram the newspapers under the sofa in the [ $L^*H-H\%$  guestroom].

The need to prepare for that final rising tone might make a prosodic break before the PP more likely, and thus ease reanalysis or even encourage a different prosodic chunking which might encourage argument attachment. This possibility is revisited and more fully explained in Section 4.2.1.

The prosodic structures observed found in the data collected for the current study are discussed in 3.2. For discussion of what constitutes a prosodic boundary, see

<sup>2</sup>Hedberg et al. (2017) use  $\uparrow$  to indicate an up-step, which is not standardly transcribed with ToBI. **[5]:** This has been moved up

e.g., Streeter (1978) and Salverda, Dahan, & McQueen (2003).

## 1.5 Evidence that prosody can affect syntactic parsing

A number of studies have shown that in listening to speech, prosodic cues appear to help reduce the frequency with which incorrect parsing (i.e., a garden path) occurs. For example, Kjelgaard & Speer (1999) conducted a study using digital manipulation of recorded speech to create three versions of sentences containing a temporary ambiguity which could result in a garden path. They recorded speakers saying sentences with natural prosody, such as the following pair (not bracketed in presentation to the speakersparticipants):

(29) [When Roger leaves] the house is dark. (Early closure)

(30) [When Roger leaves the house] it's dark. (Late closure)

They then cross-spliced these together to make several versions. One version had prosodic cues which cooperated with the intended reading of the sentence; another attempted to have “neutral” prosody; and the third used intentionally misleading prosody which favored the garden path. The initial fragment of each (the portion from the beginning of the sentence to the word *house* in (29 - 30)) was then presented to participants who (the portion from the beginning of the sentence to the word *\*house\** in (@ee - @le) and they were asked to agree or disagree with whether a visually presented word, either *is* or *it's* was likely to be the next word in the sentence. Participants gave more accurate and speedier judgments when the visually presented word was compatible with the structure indicated by the prosodic cues when the prosodic cues lined up with the correct parsing. The results of this study, as well as a growing body of literature, suggest that ~~that~~ prosodic information

[6]: Section title has been updated

can indeed be used by the parser in making processing decisions.

Consider the analysis by Fodor (2002) invoked prosody in an explanation of differing relative clause attachment preferences of relative clause attachment preferences in English an example of the differing RC attachment preferences across languages first pointed out by Frazier & Clifton (1996) and discussed above in Section 1.2. This concerns sentences such as English (31) and Spanish (32):

(31) Someone shot the servant<sub>N1</sub> of the actress<sub>N2</sub> [<sub>RC</sub> who was on the balcony].

(32) Alguien disparó contra la criada de la actriz que estaba en el balcón.

The relative clause in(RC) (31) *who was on the balcony*, can attach either locally (modifying N<sub>2</sub>), making it *the actress who was on the balcony*, or non-locally (so that it is *the servant who was on the balcony*). While Late Closure predicts local/low attachment in these structures of sentences, Cuetos & Mitchell (1988) found a 60% preference for local/low attachment for their English materials speakers, but only a 40% preference for local/low attachment in Spanish speakers. In apparent violation of the general preference for local/low attachment, some languages, like Spanish (and French and Russian, but not Romanian or Brazilian Portuguese, so this is not a general feature exclusive to Romance languages), prefer to attach relative clauses higher, while others more often obey Late Closure as in English (e.g., Swedish and, Egyptian Arabic, and English). Interestingly, the This non-local preference in languages like Spanish is weakened in cases where the ambiguous RC is short (one prosodic word), which is compatible with its being attributable to prosodic phrasing; see below. Fodor (2002) maintains asserts that these tendencies exist both in listening to spoken sentences, with and without words (under conditions where a particular parse is not favored by the explicit prosody) and in silent reading.

OFodor notes that other researchers, e.g., Maynell (1999), have shown that the

presence or absence of a prosodic break before the RC in sentences like (~~@servant~~) encourages high or low attachment respectively. Fodor leverages this in order to explain the differences in RC attachment site tendency between languages. She argues that the cross-language differences~~phenomenon~~ can be neatly explained ~~account for~~<sup>TJP</sup> by linking attachment site preference to the likelihood of a prosodic break before the RC. This difference in prosodic tendency, in turn, can be explained using a constraints-based approach. Consider Selkirk's (1986) alignment constraints:

- (33)  $\text{Align}(\alpha\text{Cat}, E; \beta\text{Cat}, E)$
- a.  $\text{Align}(\text{GCat}, E; \text{PCat}, E)$
  - b.  $\text{Align}(\text{PCat}, E; \text{GCat}, E)$
  - c.  $\text{Align}(\text{PCat}, E; \text{PCat}, E) \setminus \text{GCat ranges over morphological and syntactic categories; PCat ranges over prosodic categories; } E = \text{Right or Left}$   
(Selkirk, 1986, p. 6)

Truckenbrodt (1999) provides a prose-based formalization of this general~~the same~~ idea for English.

(34) **Align-XP/R**

For each XP there is a PP such that the right edge of the XP coincides with the right edge of the PP, where XP is a maximal projection and PP is a Phonological Phrase. This constraint represents the end based mapping assumption for Major Phonological Phrases in English, whose right end is supposed to align with the right end of Maximal Projections (Truckenbrodt, 1999, p. 223).

Essentially, ~~@Fodor2002-io~~ argues that Therelative ranking of alignment constraints for the left edge of phrases (*Align-XP/L*) with those for the right edge of phrases (*Align-XP/R*) can impact the distribution of prosodic breaks. These

alignment constraints dictate that the edges of prosodic units (and thus the location of prosodic breaks) should align with the edges of syntactic constituents. Because the prosodic break that encourages high attachment is one which aligns with the left edge of the RC in examples like (32), postulating that *Align-XP/L* is postulated to be ranked above *Align-XP/R* in languages like Spanish French that prefer high attachment since an account for that preference (remember that a prosodic break in that place has been shown to encourage a high attachment interpretation) (Maynell, 1999). In languages where low RC attachment is preferred, we can assume that *Align-XP/R* can be postulated to be ranked higher, and thus a prosodic break is more likely to occur after the RC than before it.

The same sort of argument can explain the difference in tendency between long and short RCs. Consider Selkirk's (2011) *BinMin* defined below.

(35) **BinMin**( $\phi$ )

A  $\phi$  (phonological phrase) must consist of at least two  $\omega$  (phonological words).

If we assume, in Optimality Theoretic (Prince & Smolensky, 1993) terms, that a constraint like *BinMin* is ranked above *Align<sub>L</sub>*, then it seems quite reasonable to assume that a prosodic break before a short RC (which would encourage high attachment) is much less likely than before a long RC. That is, when the RC is short, its left edge is prevented from aligning with the beginning of a prosodic phrase (it violates *Align<sub>L</sub>*) by the higher ranked *BinMin*, because it needs more material in order to achieve a length of at least two phonological words, and so must borrow some for the preceding NP. Longer RCs can have their left edge align with the start of a prosodic phrase, and thus can have the high-attachment encouraging prosodic break, because they have enough phonological content to stand alone.

## 1.6 Predictions for the current study

The experimental study presented here addresses ~~This study is concerned with a~~ number of issues. First: is syntactic phrase attachment in any way encoded in the speech signal? I hypothesize, following e.g., Schafer, Speer, Warren, & White (2000) ~~as noted above that it is, and therefore that we can use~~ prosody can be used to diagnose attachment site. Consider the basic configuration in Figures 1.4 and 1.5 below.

[7]: not all changes to this section are well represented by annotation; we should read it together

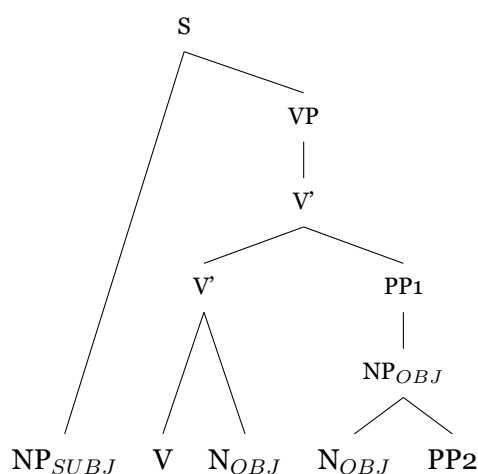


Figure 1.4: Illustrative syntactic tree of the basic configuration for Mod cases.

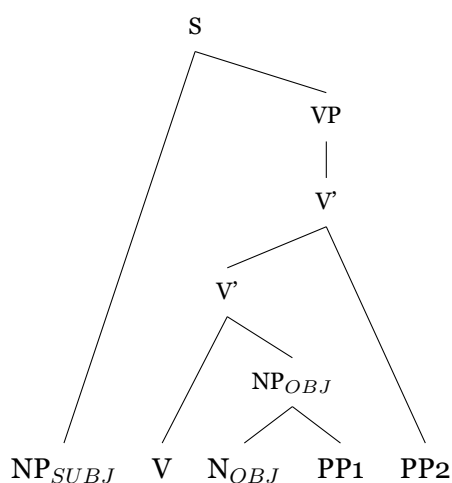


Figure 1.5: Illustrative syntactic tree of the basic configuration for Arg cases.

The predictions made here rely on there being an ideal prosodic structure for the sentences tested to compare them with. Following @Fodor2002-io, a major disruption to the syntax like the one after PP1, ~~I suggest that argument attachment of PP2<sup>TJP</sup> as in Figure 1.5, should~~ will be marked by a prosodic boundary between PP1 and PP2. On the other hand, m Modifier attachment of PP2, ~~on the other hand,~~ will lack any substantial boundary marking in that position unless it is necessitated for length reasons; instead, a break after the object is expected. These expectations are shown in (36) and (37), where ”|” a less prominent or no break and ”||” represents a linguistically motivated prosodic break.



(36)

... stick the letter      OBJ Break      PP1 Break  
                                  |||      |  
                                  in the mailbox      of the proper stack  
                                  PP1      vice president.  
                                  Direct object

(37)

... stick the letter      OBJ Break      PP1 Break  
                                  |      ||  
                                  in the mailbox      onto the proper stack.  
                                  PP1      PP2  
                                  Direct object

AnThe second issue that is addressed is: what factors impact immediate on-line parsing, and what factors only affect later, post-parse considerations? To address this, the study will employ the double reading paradigm of Fodor, Macaulay, Ronkos, Callahan, & Peckenpaugh (2019) (more on double reading in the methods section). For example, if first-pass parsing ignores semantic information, then semantically implausible parses should be more frequent in Reading 1 of a sentence with argument-PP2 than in a second reading of the same sentence.

(38) *Hypothesis 1* A first reading of a sentence where PP2 is a goal argument (Arg) will exhibit less natural prosody (more hesitation at and within the PP2 region) than:

- a. A first reading of a sentence where PP2 is a modifier (Mod)
- b. A second reading of a sentence where PP2 is a goal argument (Arg)

(39) *Hypothesis 2* A first reading of a sentence where PP2 is a goal argument (Arg) will more often be produced with prosodic structure that represents an implausible or ungrammatical parse of the string (PP2 incorrectly attached as a modifier), whereas a second reading of that sentence will more often be pronounced with the prosodic structure that represents the intended parse (argument attachment of PP2).

*Hypotheses 1* and *2* together make a third prediction: readers should struggle more on the cold reading of a sentence with an argument PP2 (Arg) to obtain a plausible structure, and thus the appropriate prosody, than on a second ~~previewed~~ reading.

(40) *Hypothesis 3* Reading 1 of a declarative sentence with an argument-PP2 will exhibit less natural prosody (more hesitation at and after the disambiguating region) and be more likely to be produced with prosodic structure that represents an implausible or ungrammatical parse of the string than a Reading 1 of an interrogative sentence with an argument-PP2.

Finally, the question of whether the 2016 Intuition described in Section 1.1 can be shown to extend to the syntactically disambiguated sentences tested in the present used in this study is formalized as~~returned to with~~ *hypothesis 4*. This hypothesis presumes that IRT reflects processing difficulty.

(60) *Hypothesis 4* The inter-reading time (IRT) will be longer for Arg sentences that **[8]: this is new** are declarative than for: a. Arg sentences that are interrogative b. Mod sentence that are interrogative or declarative

Note that these hypotheses are not applicable~~\*hypothesis 3\* can't be applied~~ in cases where the reader fails to successfully and fluently produce the sentence.

These hypotheses are returned to in Chapter 3. The findings from the current study that will be presented in what follows~~in the current study~~ do not successfully settle all of the issues raised here, but it is hoped that they will help guide future work in directions that may do so.

# Chapter 2

## Methodology

This ~~chapter describes~~section outlines the methodology employed for the reported study reported in Chapter 3. The protocol outlined is referred to as the *Double Reading Procedure* and was first implemented by Fodor et al. (2019). Under this protocol, participants are asked to read aloud a visually presented sentences twice, once without taking any time to preview sentence content (Reading 1), and then again after unlimited preview (Reading 2).

Fodor et al. (2019) aimed to investigate the extent to which preview impacted the prosodic phrasing of center embedded sentences, as well as whether or not readers would find the doubly center embedded sentences more comprehensible after preview (or, comprehensible at all, as the doubly center embedded sentences often were not on first attempt) under the assumption that in order to pronounce a sentence with the optimal prosody, it's necessary for the speaker to understand the sentence. In the prosody literature up to this point, preview has largely been ignored as a factor in reading aloud tasks. Fodor et al. (2019) found that preview did indeed impact both the prosodic grouping that readers used, suggesting that comprehension was improved on the second readingand ~~comprehensibility~~.

While the questions being addressed in the current study ~~here~~ are different from those of Fodor et al. (2019), it is~~we are~~ still concerned with the prosody that is produced, as well as the degree of difficulty the reader experiences in parsing a sentence in order to read it aloud. This experimental paradigm eliminates the possible uncertainty~~noise~~ of not knowing whether a given pronunciation represents a naive or considered ~~or naive~~ attempt to read a sentence aloud.

## 2.1 Materials

In total there were 16 experimental items each constructed in 4 versions, and 32 filler items each in two versions. The design decisions are discussed in detail in this section.

### 2.1.1 Experimental items

The basic experimental items were created in a 2 x 2 design with one factor being Speech Act (interrogative/Q vs. declarative/D) and the other being PP2 Status, i.e., PP2 was either a PP1 which must be an argument of the verb (Arg) or else one which must~~can~~ be a modifier (Mod) of the preceding phrase (PP1). A full list of experimental items is available in Appendix A.

Table 2.1: Illustrative experimental item, constructed in four versions.

Version	Sentence
D Arg	He had decided to stick the large check in the envelope into her wallet.
Q Arg	Had he decided to stick the large check in the envelope into her wallet?
D Mod	He had decided to stick the large check in the envelope from her church.
Q Mod	Had he decided to stick the large check in the envelope from her church?

As previously mentioned, the current study was motivated by an observation first reported in Peckenpaugh (2016). The experimental stimuli used in the current

study were based on the materials from that study ~~an earlier pilot study exploring this same a similar phenomenon @qp2~~, with several adjustments made to accommodate the objectives of the current study. The sequence of parts for each of the basic items was always the same, shown in (41). Note that the material starting with the infinitival verb, e.g., *cram* until the end of sentence will be referred to as the construction throughout this and later chapters, as labeled in the example below.

(41)

Introductory material			Construction			
Subject	Auxiliary	Matrix verb	Infinitival verb	Object	PP1	PP2
Order shown for D versions (reversed in Q)					always ambig- uous	Disam- big- uation of PP1

All four versions of any given quadruple used the same introductory material, the only difference arising through the necessary inversion of auxiliary and matrix subject, as required by the Speech Act factor. Across quadruples, subjects alternated between *she* and *he*, with half using one and half using the other; the auxiliary was always *had*. The matrix verb did not vary within a quadruple, but did vary between quadruples; for any given quadruple, the matrix verb was one of four verbs of mental state (*decide*, *intend*, *plan*, or *want*). The rationale for the use of these mental state verbs are discussed later in this section.

The verb within the construction did not vary within a quadruple, but a given quadruple could have one of four verbs: *cram*, *put*, *set* or *stick*. The construction verb form was always infinitival. Each construction verb appeared in four different quadruples, and was paired once with each matrix verb, to create 16 unique pairings of matrix verb to construction verb. Thus, for matrix verb *decide*, for example, *decided to cram*, *decided to put*, *decided to set*, and *decided to stick*; and for construction verb *cram*, *decided to cram*, *intended to cram*, *wanted to cram*, and *planned to cram*.

The word order and content of the construction was the same across all versions of a quadruple, with the exception of the content of PP2 which varied across the PP2 Status factor: The Arg versions of a quadruple had a PP2 which was headed by *into* or *onto*, while the Mod versions had a PP2 which was headed by *of* or *from*.

PP1 was the same across versions of a given quadruple, e.g., *cram the paperwork in the drawer...* (see Table 2.1's illustrative example). That is, PP1 was identical (and temporarily ambiguous) in every version of a given quadruple, being interpretable as either the goal argument of the construction verb or as a modifier of the object NP. However, in Arg versions of a quadruple, the argument interpretation of PP1 cannot be sustained once PP2 is encountered. In those cases PP2 must fill the goal argument slot and PP1 must be a modifier. The working assumptions about parsing discussed earlier, i.e., that the parser will initially postulate assume PP1 to be the goal argument due to the primary status of arguments, predicts assumes that (all and only) Arg versions of a quadruple require will reanalysis. Between quadruples, the preposition that headed PP1 varied, but was always one which was compatible with it being a goal argument or a modifier of the object: *in* in 8 cases(8), *on* in 7 cases(7), and in one case, *under*.

AOne benefit of using a complex verb cluster (auxiliary + matrix participle + infinitive; see ?? above) rather than a single verb<sup>1</sup> was that the differences between declarative and interrogative versions of a quadruple were it isolated the differences across the versions of a quadruple triggered by the Speech Act factor to the left extremity of the introductory material of the sentence: i.e., only the position of the subject and the auxiliary were affected, meaning that the construction itself and several words prior to it were was completely untouched by the Speech Act this

<sup>1</sup>Note that the use of an auxiliary also eliminates length differences across D vs. Q versions of a quadruple: if an auxiliary verb were not present, interrogative versions of a basic item would have an extra word, the result of so-called *do*-support, that would not appear in the declaratives (e.g., *he crammed ...* vs. *did he cram ...?*)

manipulation. The construction is in a sense buffered from changes triggered by Speech Act manipulation, and is only effected by the PP2 Status manipulation.

The purpose of including introductory matrix verbs (e.g., *intended*) was to reduce the oddity of the polar interrogative versions of each quadruple. It seems odd to ask, “Did Mary put the jelly beans in the window onto a fancy dish?” because, when it is clear that the speaker already knows so much about the situation, it becomes difficult to imagine a pragmatically plausible context where such a question would be asked. Such sentences might well be described as “prosecutorial<sup>2</sup>.” Arguably, this is somewhat mitigated by the addition of a verb like *intended*: rather than asking about facts that we already seem to know, we are instead asking about an actor’s mental state with regard to those facts. Even if we know the facts of the situation, we do not necessarily know, for instance, whether it was the result of a decision, some third party’s action, or mere happenstance. Another adjustment made in order to make the polar interrogative versions of each quadruple more pragmatically acceptable limited the amount of detail in the experimental sentences, so that fewer adjectives and adverbs were included compared to the items employed in Peckenpaugh (2016), and subjects were always third person nominative pronouns (*he* or *she*).

Importantly, the construction verb was always one which demanded a goal argument. Where some of the verbs used in the items employed by Peckenpaugh (2016) only optionally took a goal, the current study used only verbs which require a goal argument: *cram*, *put*, *set* or *stick*. Verbs that optionally take a goal (e.g., *hide*) might result in a parse where PP1 is not immediately incorporated as the goal argument, which would mean that PP2 would not necessarily force reanalysis. Consider the contrasting sub-categorization of the verbs in (42) and (43):

---

<sup>2</sup>Thank you to Dr. Dianne Bradley for making this observation, and for the very clever “prosecutorial” descriptor.

(42) **Optional goal** (*hide*) \ *The gangsters had hidden the shotguns in a U-Haul truck.* \ ✓ *The gangsters had hidden the shotguns.*

(43) **Obligatory goal** (*put*) \ *The gangsters had put the shotguns in a U-Haul truck.* \ \* *The gangsters had put the shotguns.*

A verb like *hide*, as in (42), can take a goal, but can also be used without one. A verb like *put*, on the other hand, as in (43), really must have a goal even if not fully detailed (e.g. *put down*). The use of verbs that require[^reqGoal] a goal argument in the current study maximized the likelihood of a robust garden path effect in the Arg versions, when PP2 triggered reanalysis. The four construction verbs used in this study were: *cram*, *put*, *set* and *stick*. While certain constructions containing these verbs do exist where no goal is needed (e.g., "the student had needed to *cram* all night"; "the narrator had set the scene"; "the disgruntled worker decided to stick with the program"; and others), it is arguable whether such instances should be considered the same lexical item as the ones being used here, and in any case all experimental items used in the current study would be rendered ungrammatical by the omission of the goal argument, as shown in (43).

Another important consideration was ensuring that the Arg versions had a PP2 which definitively disambiguated the attachment site of PP1 from the goal argument role to being a modifier of the object; i.e., that reanalysis was forced.

(44) She had decided to put the child [<sub>PP1</sub> on the rocking horse] [<sub>PP2</sub> on the see-saw].

(45) She had decided to put the child [<sub>PP1</sub> on the rocking horse] [<sub>PP2</sub> onto the see-saw].

In (44), PP2 is implausible as a modifier of *rocking horse*, but not strictly impossible syntactically, and the sentence is grammatical with PP2 modifying it. On the other hand, the use of *onto* in (45) completely disallows the modifier interpretation of PP2



at the syntactic level: a PP headed by *onto* cannot grammatically modify the preceding NP.

Where Peckenpaugh (2016) relied on plausibility to force reanalysis, the current study uses syntactic disambiguation, such that the Arg versions always have a PP2 headed by *into* or *onto* which cannot head a PP2 that modifies the NP of PP1. This avoids any inconsistency in the resultsnoise that might result from discrepancies between individuals' real world knowledge or beliefs. For the Mod items, the head preposition of PP2 was always either *from* or *of*, which are compatible with a parse where PP1 is the goal argument and PP2 is modifying the NP within PP1.

It is worth noting that some linguists (e.g., Den Dikken (2006)) believe *of* is not a preposition in the same sense as *from*, *on*, or *in*, etc., in that it appears to be serving a strictly grammatical or functional purpose, without real lexical content.

Importantly, it is also only 2 characters long, whereas *into*, *onto*, and *from* (the other possible heads of PP2) are all 4 characters. This is revisited and its possible impact is explored in the results section (section 3.4.3).

To sum up, the experimental items were designed to have limited detail, with either *he* or *she* as the matrix subject. A complex verb cluster, e.g., *had decided to cram* was used to facilitate subject-auxiliary inversion without *do*-support in the interrogatives and limit the difference between items, as well as provide a verb of mental state (*decide*, *intend*, *plan*, or *want*) to support more pragmatically plausible questions. PP1 was always interpretable as either the goal argument or a modifier of the object. PP2 differed across the PP2 Status factor, but not across the Speech Act factor. In the two Arg versions of a quadruple, PP2it was headed by *into* or *onto* and was intended to force reanalysis, under the assumption that PP1 had been incorporated into the parse as an argument, since a PP headed by *into* or *onto* must be interpreted as the goal argument, which is the position that PP1 would have

presumably been occupying in the ongoing parse. For the two Mod versions of a quadruple, PP2 was headed by *from* or *of* and therefore was not expected to require reanalysis, as *from*- and *of*-headed PPs can attach as modifiers of a preceding NP (in this case, the NP within PP1), allowing PP1 to stay in the goal argument slot.

### 2.1.2 Fillers

There were 32 filler items that ranged in complexity, e.g., some contained embedded finite or non-finite clauses, some contained reduced relative clauses or full relative clauses, and some were simple matrix clauses. Of these 32, 16 were designed to end in a sequence of two PPs, to mirror the experimental items (+PP), while the other half contained no final PPs (-PP). The +PP fillers were unrelated to the -PP fillers. All fillers were designed in two versions: declarative (D) and interrogative (Q). For the +PP fillers, PP1 was an argument in 5 of 16 cases, a modifier of the object in 6 cases, and a modifier of the verb phrase in 5 cases. The distribution of attachment sites for PP2 in the +PP filler items was the same (e.g., there were 5 cases where PP2 was an argument), except there were 6 that modified the NP embedded in PP1 instead of 6 that modified the object. The purpose of this decision was to avoid any discrepancy that might result from any difference in processing difficulty of certain configurations of attachment sites for the PPs. An even distribution of attachment site for both PPs ensures that any analysis over the filler items should not be unduly impacted by PP attachment site when that is not what the analysis is concerned with, and it makes available an analysis of those attachment sites, although that analysis is not pursued here. A full list of fillers is available in Appendix B.

All filler items had the same sort of introductory material as the experimental items (*he/she + had + past participle verb of mental state*). The past participle was either one of the four mental state verbs used for the experimental items (*decide, intend, plan, and want*), or one of four additional verbs of mental state: *forgot, mean, need,*

Table 2.2: Illustrative filler items, constructed in two versions.

Version	Sentence
D +PP	He had forgotten to try the famous pastry in the restaurant of the fancy hotel.
Q +PP	Had he forgotten to try the famous pastry in the restaurant of the fancy hotel?
D -PP	She had forgotten to report that the clerk was ignoring her request.
Q -PP	Had she forgotten to report that the clerk was ignoring her request?

or *remember*, with each of the 8 past participles being used twice in the +PP fillers and twice in the -PP fillers, for a total of 4 times each. This means that a participant would see 6 instances each of *decide*, *intend*, *plan*, and *want*, (i.e., twice in experimental items and 4 times in filler items) but only 4 instances of the filler-only mental state verbs. Fillers used both mental state verbs from the experimental items as well as others was in order to prevent the experimental items as being identifiable by which mental state verb was used, and to avoid extreme amounts of repetition for any given lexical item.

### 2.1.3 Length

Sentence Length was tightly controlled across items. For experimental quadruples, all sentences were between 66 and 75 characters long, and between 13 and 15 words long. The length within a quadruple never varied across the D vs. Q factor. Across the PP2 Status factor, given that the content of PP2 differed within a given quadruple, there was a maximum length difference of one character. Two quadruples varied in word length across PP2 Status by one word. Across all quadruples an equal number were longer (word- and character-wise) in the Arg condition than in the Mod condition. The experimental items ranged from 18 to 22 syllables.

Control over filler pair length was slightly less stringent. They ranged from 63 to 79 characters and 12 to 14 words. Length was never different within a filler pair, since

only the Speech Act factor was implemented in the construction of fillers.

## 2.2 Participants recruitment

All participants in the current study were undergraduate students enrolled at Queens College in Psychology 101<sup>3</sup> who participated for course credit. Self-reported age ranged from 18 to 25 years. Participants were recruited via a software system designed for university participant pools. Students saw a recruitment notice on the system website (see Appendix C), and were able to schedule their own appointment time within the hours offered.

The 35 participants recruited were self-identified native and primary speakers of American English. One participant was disqualified post-hoc after producing a Caribbean English pronunciation pattern; one further participant was excluded post-hoc due to an extremely disfluent reading cadence. A final participant was excluded due to a technical issue. All excluded participants were still awarded class credit for participating.

## 2.3 Location

All data were collected in a private room with only the experimenter and participant present. While every effort was undertaken to ensure a quiet environment, intrusive noise from passersby or neighboring rooms were sometimes unavoidable. This resulted in some unusable or partially unusable recordings (detailed in section 3.4.1 of the results chapter).

---

<sup>3</sup>IRB approval number: 2018-0072

## 2.4 Equipment and software

The experiment was presented on a laptop running Windows 10 with stickers on the keyboard labeling relevant keys: the left shift key was labeled *START*, right shift was labeled *NEXT*, and the touch-pad was labeled *DONE*.

The presentation of items and instruction<sup>4</sup> was done using the Open Sesame software (Mathôt, Schreij, & Theeuwes, 2012) which provides a graphical user interface, scripting language, and interpretation of Python code. The system was capable of 10-20 millisecond accuracy, with the display's 60Hz refresh rate being the limiting factor. Key input had a latency of about 10ms.

Recording used a Blue Yeti USB microphone positioned near the participant's left hand and angled to point at the space in front of the participant's mouth. The angle was adjusted for each participant's height. Audio was recorded at 44.1kHz single-channel quality.

## 2.5 Versions of the experiment

The experiment was presented in 4 basic versions, with split-half ordering (where the first 24 of the items presented to one group was the second 24 presented to the other) for a total of 8 groups. Each version contained 7 practice items, 3 of which were overt practice and 4 of which were covert practice, as well as one version of each of the 16 experimental and 32 filler items. No version contained more than one version of a given experimental quadruple, or a given filler pair, and each version contained one member of every experimental quadruple and filler pair. Each participant saw the same number of each type of experimental quadruple: 4 D Arg, 4 Q Arg, 4 D Mod and 4 Q Mod. The experimental items were presented in

---

<sup>4</sup>Instructions were also provided verbally and via printout, see Appendix D for the latter.

pseudo-random order, interspersed with 1 to 3 fillers. Ignoring fillers, the same version of a different quadruple never occurred in sequence (e.g., after encountering a D Arg, the next experimental item was never another D Arg).

## 2.6 Procedure

Participants were given a verbal overview of the experimental procedure and then asked to read a one page review of the procedure (see Appendix D) before signing a consent form. Participants then sat at the computer and were again walked through instructions before the first practice item was presented. The instructions made clear that the task consisted only of reading each sentence twice in succession, under different timing conditions as specified.

Participants completed 3 practice items, then consulted with the experimenter before beginning the main portion of the study. The study also contained 4 covert practice items that were not included in any analyses, to allow participants to settle into the procedure before any results were recorded.

Participants used keyboard button presses to navigate the experimental presentation. Each such key-press terminated the current screen, and initiated display of the screen that was programmed to follow. The succession of 4 screens constituting the presentation of any item was participant-paced, as was the progress from item to item. Between items, the display defaulted to a fixation screen showing a line of ten pluses aligned with the left edge of the to-be-revealed sentence. This was designed to direct the participant's attention to the beginning of the sentence, and thus minimize unintended look-ahead (the issue of potential look-ahead is discussed at greater length in section 2.6.1). Items were uniformly presented without line breaks.

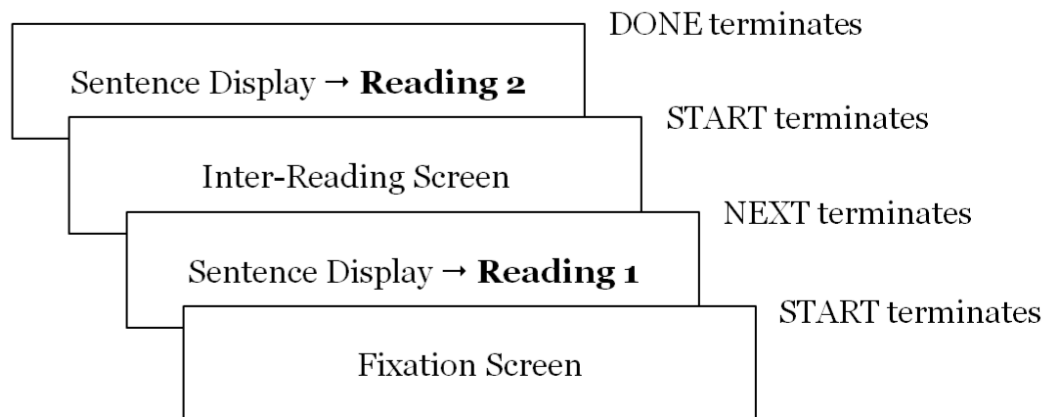


Figure 2.1: Diagram of 4-screen sequence presented for each item, showing the key presses triggering movement between successive screens.

The first *START* key-press that terminated the fixation screen and initiated the first display of a given item also began the first of the 2 recordings collected for each item. RThat recording continued through the presentation of the inter-item instruction slide and was terminated by the press of *START* that terminated the inter-item instruction slide and initiated the second display of the item.

The second display of an item was displayed in black font on a pale blue background. All other screens were displayed in black font on a pale green background (i.e., fixation screen, first display of an item, and inter-item instruction screen, as well as the initial instructions).

The inter-item instruction slide which was displayed after the first display of an item was terminated contained the following text:

Your first reading is complete.

Press *START* to begin the second reading.

The shifting of required key presses and the changing background color were

intended to prevent accidental double-presses of any given button from having unintended side effects, and to help the participant track where in the protocol a given screen was located. It took some time for the participants to adapt to the procedure, but generally the necessary habits were acquired before the first item of the experiment proper was presented.

### 2.6.1 On look-ahead

An advantage of that the Double Reading Procedure has is that it allows for certain assumptions to be made about Reading 2 for the purposes of analysis that otherwise would be unclear to the researcher: Reading 2 certainly represents a considered reading of the sentence. The reader not only has necessarily read the sentence and (heard it read) in producing Reading 1, but has the reader had ample time to examine the sentence, ~~but has necessarily read it and heard it read in producing Reading 1.~~ This means Reading 2 can plausibly be thought to represent a considered prosodic and syntactic structure, at least more so than an entirely naive reading, and should not reflect any processing issues; a parse should have already been developed during Reading 1, or during subsequent study of the sentence prior to Reading 2.

The nature of Reading 1 is less clear. Because there is variability in the delay between the display of the sentence and the onset of phonation, it is possible that Reading 1 is not entirely delivered without preview. The properties of these Reading 1 delays are discussed at length in a later section, but for now it suffices to say that the very limited preview is possible during a delay that typically falls in the 0.2 to 2.7s range (median = 1s, SD = 0.4). As an example of common reading rates, Ashby, Yang, Evans, & Rayner (2012) reported faster readers as averaging 328 words per minute (wpm), and slower readers 228 wpm, in silent reading. That study found that reading time is slower for reading aloud, and that the availability of parafoveal information (i.e., the difference between 1 word and 3 word windows) is less



impactful for that reading mode. Given that the experimental items range from 13 to 15 words, most of the R1 delays would not allow even a fast reader to read the entire sentence: the median R1 delay was 1s ~~which R1 delay~~ would allow a fast reader time to read very few words. In fact, ~~keep in mind that~~ the window is even shorter, because in addition to just reading, the subject is also performing ~~handling~~ several other cognitive processes (e.g., visual processing, lexical access, issuing motor commands, etc.). For most recordings, therefore, ~~t~~The utterance of Reading 1 should, ~~therefore,~~ contain within it any behavioral reflex of whatever online parsing difficulty the reader has, ~~for most recordings.~~

In order to clearly understand the results of this double reading study, it is important to understand the mechanics of reading. Specifically, we would want to know at what point during the reading of a temporarily ambiguous sentence the participant will become aware of the existence of a disambiguating PP2, since this is when it will be realized that the initial parse may well crash. The work of several decades on the time course of reading ~~this subject~~ is thoroughly summarized in Rayner, Pollatsek, Ashby, & Clifton (2012). They describe reading as consisting of a series of fixations, during each of which ~~when~~ foveal vision takes in a small region of the visual field, and saccades, where the eyes move ahead ballistically (i.e., on a planned trajectory that cannot be interrupted). As a consequence of the ballistic property of saccadic movement and the additional finding that landing sites (fixations) are not random, it ~~we~~ can be inferred that at least some look-ahead is available, i.e., a reader must know something about what is coming in order to plan a suitable landing site. The primary predictor of fixation point seems to be the character length of a word, meaning that the presence of characters and word boundary information (represented orthographically by spaces in languages like English) at least are necessary at the periphery of attention, i.e., within the perceptual span. Some details on the perceptual span, or the information that can

be accessed by the eyes at any given time, is discussed in brief here, with special attention to its relevance for the study at hand.

Rayner et al. (2012) discuss a number of studies that explore the size and properties of this span, the most fruitful of those studies being based on a gaze-contingent moving-window technique. In this technique, text is presented on a video monitor while the reader is also hooked up to eye-tracking equipment. A computer constantly samples the position of the reader's eyes and updates the display accordingly. ~~T~~Using this elaborate system, and the mutilation of text outside a window of clear text creates, a so-called moving window around the reader's point of fixation is created. By manipulating the size of this window, it has been~~was~~ found that reading speed is maximized when about 15 characters to either side of the fixation site is accessible to the reader~~available~~ (it turns out this is actually asymmetric, and the window need only extend~~go~~ as far as the start of the currently fixated word in the direction of what has already been read, i.e., to the left for English readers).

In a study by McConkie & Rayner (1975), in order to determine what information was available at the periphery of the perceptual span, the amount of information outside a window of clear text known to be smaller than the ideal (e.g., 21 characters, 10 to either side) was manipulated. When all other characters and spaces were replaced with X, essentially destroying all information outside the window, reading was slower than when character spaces were maintained, but all other information was obscured. Improvements in reading speed also occurred when the original characters outside the moving window were replaced with characters that had similar shape (i.e., the same pattern of ascenders and descenders) as the character they replaced, with and without spaces. Using these techniques and manipulating the size of the window, McConkie and Rayner~~they~~ were able to determine that it is

only word boundary information that is available at the extreme edge of the perceptual span; character shape (ascenders and descenders) is available about 10 characters out from the fixation point, and character identity is available more or less only for the fixated word.

The relevant question for the study at hand is as follows: how much of the sentence will the reader have seen and processed when a given word is being spoken? A typical item from the current study is displayed in (46), with the words expected to be fixated underlined, numbered by presumed fixation sequence, and labeled. The number of characters (including spaces) intervening before the start of the disambiguating region (the left edge of PP2) is displayed below each label. These counts are calculated from the initial character of the fixated word to the initial character of the disambiguating region; the actual fixation site is likely to be closer to the center of the word, meaning the distance would be shortened by a few (1-4) characters, depending on the length of the fixated word.

(46)

He had	<u>intended</u>	to	<u>stick</u>	the	<u>letter</u>	in the	<u>mailbox</u>	onto the	<u>proper stack.</u>
	1-INITIAL		2-VERB		3-OBJ		4-PP1		DISAMBIGUATION-PP2
	45		32		22		7		0

Table 2.3 ~~presents~~describes these distances averaged across items all experimental items. Note that these values do not vary across condition, because the initial fixation is locatedcounting starts after both the subject and auxiliary verb, and the counts ends before PP2. ~~T~~, and the only changes across versions are subject-auxiliary inversion and the content of PP2 which fall outside the string of material over which these counts were calculated.

From the initial fixation point, the distance to disambiguation ranges from 45 to 50 characters, with a median of 46 characters. Rif we recall that word boundary

**[9]:** table caption revised

Table 2.3: Distance in characters from fixation to disambiguation of experimental items for the current study.

	1-INITIAL	2-VERB	3-OBJ	4-PP1
Median	46	34.5	25.5	7.5
Maximum	50	38.0	27.0	9.0
Minimum	45	32.0	21.0	5.0

information is available only 15 to 18 characters to the right of fixation, thus the we can be certain that the disambiguating region is far out of view until several fixations in.

When does the reader become aware of the existence of PP2? When fixated on the direct object head noun, the range of distance is 21 to 27 characters, with a median of 25.5: PP2's content is still outside of view, even in the case of the smallest distance, and adjusting it to be a few characters smaller to account for the fact that fixation is likely to occur closer to the center of a word rather than on its first character. At most, the presence of the first few characters of PP2's preposition may be available, but certainly not the character space after it. The distance from the fourth PP1 fixation point (the head noun within PP1 that PP) to the disambiguating region, PP2, ranges from 5 to 9 characters, with a median between 7 and 8 characters. Thus, we can say with some certainty that the reader of a sentence such as (46) will be aware that another phrase, one which starts with a 4-character word, remains to be incorporated into the parse sometime after processing of the direct object, and before processing of PP1.

There is yet another factor ~~pieee~~ to consider: the so-called eye-voice span (EVS), and the fact that the readers in this study are reading aloud rather than silently.

According to Laubrock & Kliegl (2015), when reading aloud the voice is typically behind the eyes by some 10-20 characters (M = 16.2 characters, SD = 5.2 characters). Adjusting Table 2.3 by subtracting 16 from each cell, we can

approximate the position of the voice when the disambiguating region comes within the perceptual span. These values are shown in Table 2.4.

Table 2.4: EVS-adjusted character distance to disambiguation in experimental items.

	1-INITIAL	2-CONSTRUCTION VERB	3-OBJ	4-PP1
Median	30	18.5	9.5	-8.5
Maximum	34	22.0	11.0	-7.0
Minimum	29	16.0	5.0	-11.0

It is likely, then, that an oral reader's voice would actually still be on the object when the eyes' fixation begins to provide information of some kind about the existence of PP2, and will still be pronouncing PP1 when the eyes are first fixated on PP2. This raises a question about any prosodic breaks produced after the object, because it is difficult to distinguish between an intentional prosodic break at that point, and one arising from the reader using a natural position for a break tofor hesitation due to related to the garden path effect of discovering the disambiguating PP2. This property of oral reading calls into question the status of OBJ breaks reported in 3 with regard to whether they are linguistically motivated, or represent the processing difficulty experienced when PP2 is finally observed by the reader.

## 2.7 Measurements of utterance timing

The elicitation protocol described above asked participants to read each sentence twice, once with no preview at all (Reading 1), and then again without any time pressure (Reading 2). Reading 1 (R1) delay is the elapsed time after a sentence is first displayed and when the participant begins speaking. Reading 2 (R2) delay is the same measure, but from the start of the second recording, which begins after the key press that terminates the inter-item instruction slide. Inter-reading time (IRT) is a measure of the time elapsing between when a participant stops speaking after R1 and when speaking resumes for R2. IRT encompasses but is not synonymous

with R2 delay, because IRT also includes the elapsed time after the participant stops speaking and the end of the first recording, because the experiment was self-paced, and the participant might spend time studying the sentence after their first reading but before advancing to the next frame. For this reason~~In this way~~, IRT is measured across both recordings.

The process for measuring makes use of Voice Activity Detection (VAD) software, which reports whether a given interval in a sound file contains speech-like noise. It i's worth making clear that while VAD is employed, most of the measurements of interest are actually the inverse, i.e., the amount of time in a recording that does not contain speech-like noise. For each recording, the amount of time elapsed from the beginning of the recording to phonation onset and offset was found using VAD; then, R1 delay, R2 delay, utterance length and IRT were calculated as a function of each recording's length and the VAD-reported onset and offset of phonation.

The specific software used included a homemade Python script and Google's WebRTC VAD. The recordings were 44.1kHz WAV files down-sampled to 8kHz via SOX<sup>5</sup>. Google's VAD system used Gaussian Mixture Models to make probabilistic decisions as to whether a given audio frame was speech or noise (see Falk & Chan (2006) for a complete description). Google's implementation takes one parameter called aggressiveness: a 4-tier setting for the level of confidence necessary to call a given interval speech. The implementation codes this setting on a 0-3 scale, where 0 is the most lenient (most likely to label a frame as speech) and 3 is the most stringent (most likely to label a frame as noise).

The recordings vary in the volume of the speaker's voice and the amount of background noise present. An algorithm was constructed to allow for the most stringent (highest VAD aggressiveness) measurement of the least modified data that

---

<sup>5</sup>Google's VAD API only accepts WAV files with sample rates that are a multiple of 8kHz. It ultimately down-samples all files to 8kHz, regardless of the input sampling rate.

gave plausible measurements. Specifically, each file was measured using the highest possible aggressiveness for the VAD algorithm and no modification of the recording. If the timings detected were not plausible, the timings were re-measured with the same rejection rate, but after the recording had undergone a 200Hz high-pass filter<sup>6</sup> (HPF). If that still failed, a 400Hz HPF was used. After a further failure, the VAD aggressiveness was lowered, with each HPF value tried again (0, 200Hz, 400Hz); and that process was itself repeated until the lowest possible rejection rate was tried of the four possible settings. The majority of measurements were collected using the highest aggressiveness (85.4%), with more than half requiring no HPF (59.6%) and most of the remaining recordings requiring a 200Hz HPF (40.1%).

A plausible set of measurements was required to meet the following criteria:

A. *Utterance length*: An utterance length between 2s and 10s, where utterance timing is the longest contiguous span in the recording that VAD reports as phonation, with breaks in phonation of less than 1s not breaking contiguity, because as Goldman-Eisler (1961) found that a large majority (82.5 to 87%) of pauses in fluent speech are less than 1s. Stimuli range from 18-22 syllables in length. If we assume a speech rate of 3 to 7 syllables per second (Jacewicz, Fox, & Wei, 2010) we would expect utterances between 2.5s and 7.3s. Conservative thresholds higher and lower than the expected were used, especially on the higher end, to allow for possibly very slow readings of the admittedly difficult items being tested in the current study ~~any difficulties processing or fluency that might have lead to longer reading times.~~

B. *Minimum leading silence*: A leading silence (“delay”) of more than 120ms. Even a very fast human reaction time should not permit a delay shorter than 120ms, so a shorter delay likely means an inaccurate set of measurements has been reported.

C. *Maximum edge silence*: A maximum trailing and leading silence length of less

---

<sup>6</sup>The exact algorithm is available on github (URL: [bit.ly/2uMrerG](https://bit.ly/2uMrerG))

than 95% of the file's length was also used, in order to filter out recordings that do not represent a valid trial. Very long silences less than this very conservative threshold that impact the IRT are dealt with in the data clean-up rather than via phonation detection, as described in the results section of this paper (section 3.4.1).

With 32 participants reading 48 items (experimental and filler) twice each, there are an expected number of 3072 recordings; due to technical issues at the time of data collection, 71 recordings are missing. Of the 3001 recordings subjected to this treatment, 2976 resulted in plausible timings[^handset]. A review of those that did not result in plausible timings found 9 recordings that were too noisy for computer analysis, but still usable, and those timings were recorded by hand.

To verify the accuracy of the computer measurement, timings were collected by hand for 240 recordings. There was a significant positive correlation between hand-measured and computer-measured timings ( $r(118)=0.87$ ,  $p < 0.001$ ), with a median difference of  $0.4s^7$  ( $SD = 1.5$ ).

## 2.8 Prosodic judgments

A trained linguist informant naive to the research being conducted listened to the 978 recordings of experimental items (note that 46 recordings were missing or omitted, as discussed in Section 3.1) and reported the presence or absence of breaks in certain regions of each the sentence, as well as several other judgments (e.g., the relative strength of these breaks, whether or not the reader struggled, and whether or not the reader used final-rise). Analyses of some of these judgments (e.g., the presence or absence of the OBJ and PP1 break) are reported in Chapter 3, while others lacked interesting results and were not reported (e.g., the presence or absence of the V break). The raterShe was instructed to familiarize herself with a

---

<sup>7</sup>Hand measurement was done to the nearest half second, so a fair amount of error is to be expected.



speaker's speech patterns before rating any recordings by listening to 6 filler item recordings from that speaker. She was given a diagram of the sentences as in (47), as well as full plain-text lists of all items.

(47) She had wanted to set % the textbooks % on the top shelf % into the file box.  
V OBJ PP1  
V OBJ PP1 PP2

The rater She was asked to report on whether or not she heard a prosodic boundary directly following the region labeled **OBJ**, and directly after the following the region labeled **PP1**. Instruction regarding how to define The following definition of prosodic break was provided by way of the following prose:

Please work with the assumption that “prosodic boundary” in what follows is any subset of the following features, clustered in such a way as to trigger your intuition that a new prosodic element (of any size) is beginning: pitch change, volume change, segmental lengthening, or pause.

The judgments requested also included whether or not the speaker struggled, where that struggle began, whether or not the speaker used question intonation, and which break(s) were stronger or more prominent than which other break(s) in the same utterance.

Detailed instructions on the order in which items should be listened to, both within speaker and across speakers, were also provided. The result was that she never listened to both readings of a sentence in sequence; she never listened to 2 Reading 1 versions of different sentences in sequence; and she never listened to the sentences in the same order for a given participant as she did for the previous one.

Details on the instructions given and the judgments collected can be found in Appendix E.

This strict procedure was implemented to hinder the informant from recognizing any patterns in the data, e.g., a systematic difference between Readings 1 and 2. It also mitigated any ordering effects that might occur in the data or as a result of the informant's own process. The familiarization process via filler items allowed the informant to judge the existence of breaks relative to the typical cadence and fluency of a given speaker, prior to exposure to any of the experimental items for that speaker.

### 2.8.1 Reliability

A second trained linguist, also naive to the purpose of the research, repeated the task over 120 recordings selected from 8 participants (two from each group, one per ordering). Even number experimental items were used from 4 participants, and odd numbered from the other 4. There were 8 recordings missing from the 128 selected, so the reliability task resulted in judgments over 120 recordings. The first informant also blindly re-rated those 120, with the recording name obscured and instructions not to revisit her original ratings. Reliability scores (percent of recordings agreed upon) are reported in Table 2.5.

Table 2.5: Percent agreement between the original ratings and the second rater (inter-rater) or the second rating by the original rater (intra-rater).

	OBJ	PP1	Break strength
<b>Inter-rater</b>	65.0% K = 0.17** (z = 2.61)	78.3% K = 0.09 . (z = 1.86)	54.2% K = 0.25*** (z = 3.99)
<b>Intra-rater</b>	77.5% K = 0.52*** (z = 5.73)	85.0% K = 0.52*** (z = 5.82)	72.5% K = 0.44*** (z = 5.70)

*Note:*

\*\*\* p < 0.001; \*\* p < 0.01; \* p < 0.05, . p < 0.1

The lower intra-rater agreement for relative break strength was likely impacted by

the method of reporting: because the informant was actually asked to provide judgments over three break locations (the third, V, is omitted throughout this report because it was extremely rare, occurring in just over 8% of recordings). As such, disagreement on that break and the fact that break strength is actually a compilation of two judgments (weakest and strongest break) amplified the noise to some extent.

# Chapter 3

## Results and discussion

This ~~chapter~~section reports various descriptions and analyses of the recordings obtained, and the relevance of those findings to the research questions motivating this study. The reported results include the effect of Speech Act (declarative/D vs. interrogative/Q) and PP2 Status (argument/Arg vs. modifier/Mod) on the location of prosodic breaks, as well as on time spent by a participant considering reflecting upon a sentence between readings, which will be called~~I-call~~ inter-reading time (IRT). In order to evaluate the extent to which participants adhered to the protocol as intended, i.e., began to read immediately for Reading 1 as opposed to producing a considered reading in Reading 2, the delay for which a sentence is displayed before a participant begins to read it is compared for Reading 1 (R1 delay) vs. Reading 2 (R2 delay). The prosodic patterns for participants with especially fast and especially slow R1 delays are presented as a way of investigating the extent to which individual differences might impact those patterns, and as a further exploration of the success of the protocol instructions in producing the intended behavior. A finding on the apparent processing cost of interrogative context when compared to declarative context among the filler sentences is also reported.

### 3.1 Data for analysis

Data for 32 total participants were analyzed. Given 4 versions of the experiment and 2 possible orderings there would ideally be 4 participants per version-order combination. Ultimately, 3 participants had to be excluded for different reasons, resulting in the distribution ~~is~~ as shown in Table 3.1<sup>1</sup>. Participants were removed for the following reasons: one for use of a non-standard dialect, one for extremely disfluent oral reading, and one who was missing more than half of the expected recordings because of a system crash during the procedure.

Table 3.1: Number of participants per version-order combination.

	Order		Sum
	1	2	
Version 1	5	4	9
Version 2	4	4	8
Version 3	4	4	8
Version 4	2	5	7
Sum	15	17	32

Some of the expected 3072 recordings (32 participants x 2 readings x 48 items, 16 experimental and 32 filler) were not used due to intrusive noise during the recording session. Additionally, data were also excluded from analysis if any (Reading 1/Reading 2 pair) was missing; there were 9 such incomplete pairs excluded. Without analyzable data from both members of a pair, it is difficult to determine the extent to which the elicitation protocol was executed as intended (i.e., the extent of preview for Reading 1 vs. Reading 2).

**[10]:** parenthetical revised

For experimental items, 978 recordings were subjected to prosodic analysis, constituting 95.6% of the utterances elicited. Because IRT data is a property of pairs of considered utterances in pairs (Reading 1 and/ Reading 2) rather than single

<sup>1</sup>The two 5-count cells include 2 additional participants whose data were collected in pursuit of another full set (i.e., towards an expansion to 40 participants) that was not completed due to a lack of participant sign-ups.

recordingsseparately, the database for response timing took in 489 data points.

Table 3.2: Number of recordings analyzed, as a function of Speech Act and PP2 Status.

	D	Q
Mod	246	248
Arg	244	240

**[11]:** this table has had the order of the rows reversed

## 3.2 Prosodic break patterns

ThisThe section will report the prosodic phrasings found in the recordings collected, and the extent to which those patterns are or are not influenced by the design parameters of the study (Speech Act and PP2 Status), as well as which reading (Reading 1 or Reading 2) the recording represents. These data areThis reported first descriptively (i.e., in terms of frequency), and then using regression models to calculate the statistical significance of whatever effects are found. Finally, a summary of findings and their implications for the hypothesis motivating this study is provided.

In what follows, the distribution of OBJ breaks and PP1 breaks are reported as a function of the four sentence types created by the materials design (D/Q x Arg/Mod), for each of Reading 1 and Reading 2. Then, the patterns of breaks over the two positions are considered, before moving to statistical analysis. Note that while breaks after the infinitival construction verb were reported, these breaks were exceptionally rare and occurred in only 8% of recordings, so they have been set aside. The break locations are indicated with a % symbol in (48).

			V			OBJ				PP1				
	She	had	wanted to set	%	the textbooks	%	on the top shelf	%	into the file box.					
(48)			V		OBJ		PP1		PP2					

As noted in section 2.8, the results reported are based on the subjective judgments of a trained linguist who was naive to the purposes and hypotheses underlying the

research. A second linguist provided judgements over a subset of the recordings, and a comparison between their judgments is available in Section 2.8.1 above.

### 3.2.1 Individual break patterns

The presence of the OBJ break was sensitive to both Speech Act and reading, with Reading 2 showing a different distribution across the D vs. Q distinction than the Reading 1 recordings.

**[12]:** this table has had its formatting changed so that freq is in parenthesis and percentage is not, and the order of the rows have been reversed

Table 3.3: Percent occurrence of OBJ break (frequency of occurrence in parenthesis) as a function of sentence type and Reading.

	Reading 1		Reading 2	
	D	Q	D	Q
Mod	77.2% (95)	76.6% (95)	84.6% (104)	72.6% (90)
Arg	57.4% (70)	56.7% (68)	73.0% (89)	74.2% (89)

The PP1 break was almost always present for cases where PP2 was an argument; and it was present substantially less often, but still there a majority of the time, for cases where PP2 could be interpreted as a modifier. Speech act and reading did not appear to impact the overall distribution of the PP1 break.

**[13]:** this table has been updated to have the order of the rows reversed

Table 3.4: Percent occurrence of PP1 break (frequency of occurrence in parenthesis) as a function of sentence type and Reading.

	Reading 1		Reading 2	
	D	Q	D	Q
Mod	68.3% (84)	68.5% (85)	84 (68.3%)	83 (66.9%)
Arg	99.2% (121)	99.2% (119)	121 (99.2%)	117 (97.5%)

### 3.2.2 Combined break patterns

When looking at both breaks together, a sentence could have one of four patterns: both the OBJ and PP1 break present; only OBJ present; only PP1 present; or neither

break present. There were only 5 cases where neither was present, and those were omitted in all subsequent reports the tables of prosodic patterns.

Table 3.5: Percent occurrence of both breaks as a function of sentence type and Reading.

	Reading 1				Reading 2			
	Mod		Arg		Mod		Arg	
	D	Q	D	Q	D	Q	D	Q
<b>OBJ only</b>	31.7%	30.9%	0.8%	0.8%	31.1%	31.4%	0.8%	2.5%
<b>Both</b>	45.5%	46.3%	56.6%	55.8%	54.1%	43.0%	72.1%	71.7%
<b>PP1 only</b>	22.8%	22.8%	42.6%	43.3%	14.8%	25.6%	27.0%	25.8%

**[14]:** this table has been corrected; the wrong data were labeled Reading 1 and Reading 2 (they were reversed)

There was very little difference across readings, with the following generalizations of the data shown in Table 3.5 holding for both Reading 1 and Reading 2. Table 4.5 shows that For Mod sentences the OBJ-only pattern is relatively frequent (31.1% in declaratives, 31.4% in interrogatives), whereas for Arg sentences there are very few instances with the OBJ-only pattern (0.8% in declaratives, 2.5% in interrogatives). The pattern with both breaks is less common for Mod (54.1% in declaratives, 43.0% in interrogatives) sentences than for Arg sentences (72.1% in declaratives, 71.7% in interrogatives). The PP1-only pattern occurs at about the same rate in Mod interrogatives (25.6%) as in Arg declaratives (27%) and Arg interrogatives (25.8%), but is noticeably less common for Mod declaratives (14.8%). These proportions are visually represented in figure 3.1.

### 3.2.3 Break Dominance

The relative strength of the PP1 and OBJ breaks was also collected. Figure 3.2 incorporates this information, where “PP1 dominance” means that the PP1 break was reported to be stronger than the OBJ break; “OBJ” dominance means the opposite; and “Equal strength” means that neither break was reported to be stronger than the other (the 5 instances with no breaks were again omitted).



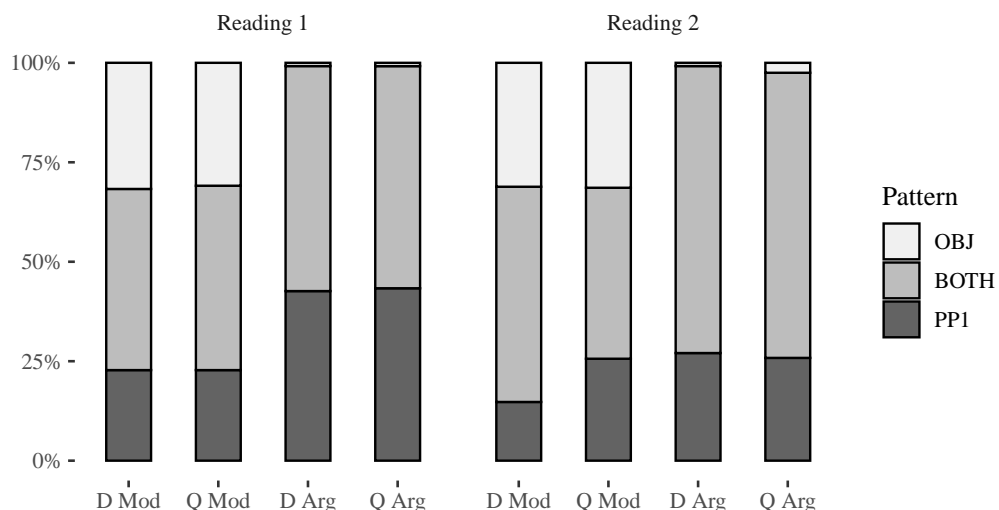


Figure 3.1: Break pattern as a function of sentence type and Reading.

For both the analysis where relative strength is ignored, and the analysis of break dominance ~~When looking at the combined break patterns, the data can be thought of as existing in of one can think of there being three bins: the PP1 (only or dominant) bin, the OBJ (only or dominant) bin, and a neutral (both breaks present or neither break dominant) bin between them.~~ In Section 3.2.2, the neutral bin containing instances of both breaks occurring is robust. ~~The~~ This break dominance analysis in Figure 3.2 distributes most of those cases that have both breaks into either the OBJ or PP1 bin, depending on which break is more prominent. When the breaks are of equal strength, they remain in the middle bin, but there are much fewer such cases of "neither dominant" for the strength-sensitive data ~~when looking at~~ than of "both breaks present" for the ~~instead of simple occurrence data.~~

Figure 3.2 clearly shows a robust effect of PP2 Status (Mod or Arg) on break dominance, and little to no impact of Reading or Speech Act (Q or D). A dominant break after PP1 is frequent for all Arg sentences, both declaratives and questions, and in both Reading 1 and Reading 2. It is less frequent for Mod sentences in both

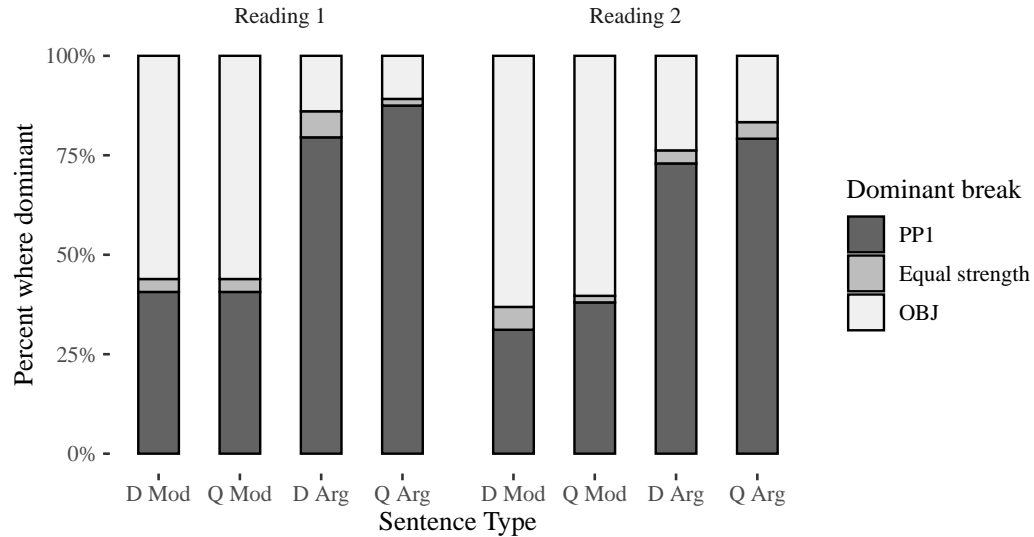


Figure 3.2: Percent break dominance occurrence as a function of sentence type and Reading.

Reading 1 and Reading 2. The higher relative frequency of PP1 break dominance in Arg sentences compared to Mod sentences is expected, since it represents a linguistically motivated prosodic break in the Arg cases. For the Mod sentences, there is no linguistic motivation for the PP1 break, so it makes sense for it to be less frequent and less dominant in Mod sentences. It is noteworthy and puzzling that the patterns do not differ greatly across the two readings, as one would expect difficulty getting the prosody right on the first try for the difficult sentences being tested in the current study.<sup>TJP</sup>

### 3.2.4 Regression models of prosodic break patterns

A number of mixed effects logistic regression models support the general observations above. Models predicting PP1 break, OBJ break, PP1 break dominance and OBJ break dominance are reported. All models include crossed random effect intercepts (participant and item), but due to convergence errors, no random slopes for any predictors are included.

The intercept always represents the Mod sentence type, which is not expected to present any particular difficulty to the reader, since the Mod PP2 Status is compatible with what is assumed to be the running parse when it is encountered (i.e., PP1 has been interpreted as the goal argument of the verb, and PP2 does not disrupt that interpretation). For those models where Speech Act is included in the model, the intercept represents the declarative sentence type. In this way, the more complex sentence types are compared to the simplest available in the model. If Reading is included in a model, the intercept represents Reading 1.

For each analysis, a reduced model and the full model (i.e., the model containing all predictors of interest) is reported. In each case, the reported reduced model is the one with the lowest reported Akaike Information Criterion<sup>2</sup> (AIC) from the set of models that include any subset of the following predictors: Speech Act, PP2 Status, Reading, and the interactions between Speech Act, PP2 Status and Reading. This method of model selection is consistent with the proposal of Wax & Kailath (1985).

Model comparisons did not always find significant differences between the more complex models, but in each case, the selected model was compared to a minimal model where where fixed effect variables were removed, leaving only an intercept, and all reported models represent improvement over the minimal model to a statistically significant degree. That comparison is reported for each model. All regression models were run using the lme4 R package (Bates, Maechler, Bolker, & Walker (2019)), with p-values calculated via the lmerTest R package (Kuznetsova, Bruun Brockhoff, & Haubo Bojesen Christensen (2019)).

---

<sup>2</sup>AIC is a representation of the amount of information lost by using a regression model to estimate data points. It is a measure that balances both the goodness of fit of a model and the simplicity of a model, guarding against over fitting and under fitting the data involved.

### 3.2.4.1 Break occurrence

In the full model predicting OBJ break occurrence, shown in Table 3.6, only the estimate for D Mod Reading 1 (the intercept) and the effect of PP2 Status show statistical significance.

Table 3.6: Mixed effects logistic regression model predicting OBJ break occurrence (FULL).

Outcome: OBJ break (FULL)	Estimate	Std. Error	p
D Mod, Reading 1 (Intercept)	0.70	0.09	< 0.001
Q	0.11	0.11	0.34
Arg	-0.28	0.11	< 0.05
Reading 2	0.07	0.05	0.15
Q:Arg	-0.14	0.16	0.39
Q:Reading2	-0.11	0.07	0.12
Arg:Reading2	0.08	0.07	0.27
Q:Arg:Reading2	0.13	0.10	0.19

Table 3.7 shows a reduced model, predicting the occurrence of an OBJ break with estimates for the coefficients of the fixed effects of Reading 2, PP2 and the interaction between Reading and PP2 Status. The removal of other predictors allowed the Reading x PP2 Status to become a significant predictor. A comparison between the reported model and a minimal one found that the reported model was better with a high level of confidence ( $AIC_{MIN}=1068.0$ ,  $AIC_{BEST}=1031.6$ ,  $\chi^2(2)=30.5$ ,  $p < 0.001$ ).

Table 3.7: Mixed effects logistic regression model predicting OBJ break occurrence (REDUCED).

Outcome: OBJ break (REDUCED)	Estimate	Std. Error	p
D Mod, Reading 1 (Intercept)	1.39	0.45	< 0.01
Reading 2	0.11	0.23	0.62
Arg	-1.98	0.50	< 0.001
Reading 2 x Arg	0.81	0.32	< 0.05

The log odds<sup>3</sup> of an OBJ break for Mod Reading 1 is 1.39 (std. error = 0.45,  $p < 0.01$ ).

<sup>3</sup>Log odds is, in this case, the natural log of the odds ratio, so the log odds of A is  $\log_e(P(A)/P(\neg A))$ .

The log odds of that break increased in Reading 2 but the increase was not statistically significant. PP2 arguments reduced the log odds of an OBJ break compared to PP2 modifiers by a robust amount, but less so in Reading 2 than in Reading 1.

The OBJ break is expected to occur more often in Mod cases, because that break marks the argument attachment (and therefore a change in branching direction) of PP1.

The full model for predicting PP1 also showed significance only for the intercept and the effect of PP2 Status.

Table 3.8: Mixed effects logistic regression model predicting PP1 break occurrence (FULL).

Outcome: PP1 break (FULL)	Estimate	Std. Error	p
D Mod, Reading 1 (Intercept)	0.68	0.07	< 0.001
Q	0.01	0.09	0.89
Arg	0.31	0.09	< 0.001
Reading 2	0.00	0.04	1.00
Q:Arg	0.00	0.13	0.97
Q:Reading2	-0.01	0.06	0.82
Arg:Reading2	0.00	0.06	1.00
Q:Arg:Reading2	0.00	0.08	0.97

The best model for predicting the occurrence for the PP1 break was one where only PP2 Status was included as a predictor. The chosen model was again significantly better than the minimal model ( $AIC_{MIN}=855.6$ ,  $AIC_{BEST}=629.6$ ,  $\chi^2(1)=228.0$ ,  $p < 0.001$ ).

Sentences with argument PP2s had greatly increased log odds of a PP1 break compared to ones with modifier PP2s. This is again expected, because the PP1 break is indicating the change in branching direction for argument attachment of PP2.

A log odds of 1.39 translates to an odds ratio of 4.01:1 ( $1.39^e=4.01$ ) and a probability of 80% ( $4.01/(1+4.01)=0.80$ ).

Table 3.9: Mixed effects logistic regression model predicting PP1 break occurrence (REDUCED).

Outcome: PP1 break (REDUCED)	Estimate	Std. Error	p
Mod (Intercept)	0.96	0.30	< 0.01
Arg	4.12	0.44	< 0.001

That Speech Act is not a relevant predictor is evidence against a prosodic explanation of the motivating intuition for this study; we would expect both a main effect of Speech Act and definitely an interaction between Speech Act and PP2 Status, if the prosody were more (or less) different across the PP2 Status factor for interrogatives than for declaratives.

### 3.2.4.2 Break dominance

Models were also run for predicting break dominance. The full model predicting OBJ break dominance is shown in Table 3.10.

Table 3.10: Mixed effects logistic regression model predicting OBJ break dominance (FULL).

Outcome: OBJ dominance (FULL)	Estimate	Std. Error	p
D Mod, Reading 1 (Intercept)	0.50	0.09	< 0.001
Q	0.03	0.12	0.82
Arg	-0.45	0.12	< 0.001
Reading 2	0.07	0.05	0.22
Q:Arg	-0.03	0.17	0.86
Q:Reading2	-0.03	0.07	0.68
Arg:Reading2	0.03	0.07	0.71
Q:Arg:Reading2	0.00	0.11	0.97

Table 3.11 reports the best model for predicting OBJ break dominance. The best model was one with fixed effects for reading and PP2 Status. There was no statistically significant effect of Speech Act on OBJ break dominance.

PP1 break dominance and OBJ break dominance are not entirely complementary, because it is possible for both breaks to have equal prominence. As such, models

Table 3.11: Mixed effects logistic regression model predicting OBJ break dominance (REDUCED).

Outcome: OBJ dominance (REDUCED)	Estimate	Std. Error	p
Mod, Reading 1 (Intercept)	-0.16	0.32	0.62
Reading 2	0.40	0.16	< 0.05
Arg	-2.32	0.18	< 0.001

predicting PP1 were also explored.

The full model predicting PP1 break dominance failed to converge, so only the reduced model is reported.

Table 3.12 reports the best model for predicting PP1 break dominance. Unlike the model for predicting OBJ break dominance, the best model for predicting PP1 break dominance includes Speech Act as a predictor. The best model is one with fixed effects for reading, Speech Act, and PP2 Status.

Table 3.12: Mixed effects logistic regression model predicting PP1 break dominance (REDUCED).

Outcome: PP1 dominance (REDUCED)	Estimate	Std. Error	p
D Mod, Reading 1 (Intercept)	-0.19	0.33	0.57
Reading 2	-0.38	0.15	< 0.05
Q	0.31	0.15	< 0.05
Arg	2.20	0.17	< 0.001

This model was better than a minimal model ( $AIC_{MIN}=1290.4$ ,  $AIC_{BEST}=1078.8$ ,  $\chi^2(3)=217.59$ ,  $p < 0.001$ ). PP1 break dominance was much more likely for sentences with argument PP2s than sentences with modifier PP2s, with interrogatives having slightly increased log odds of PP1 break dominance. Log odds of PP1 break dominance were slightly less in Reading 2 than Reading 1. There were no significant interaction terms.

Because reading was a significant predictor for 3 of the 4 models reported, and there are theoretical reasons to believe that Reading 2 is more representative of the

natural or intended prosody of the reader, models were also run predicting PP1 dominance and OBJ dominance for Reading 2 data only. In both cases, the best model had the same structure: fixed effects of Speech Act and PP2 Status, with no interaction term.

Table 3.13: Mixed effects logistic regression models predicting break dominance in Reading 2 (REDUCED).

(Reading 2 only)	Outcome: OBJ Dominance			Outcome: PP1 Dominance		
	Estimate	Std. Err	p	Estimate	Std. Err	p
D Mod (Intercept)	0.66	0.24	< 0.01	-0.97	0.27	< 0.001
Q	-0.30	0.22	0.16	0.35	0.22	0.1
Arg	-2.07	0.24	< 0.001	2.15	0.24	< 0.001

For both OBJ dominance and PP2 dominance, the main effect of Speech Act is non-significant, but its inclusion marginally improves the fit of each model. Even when limited to only Reading 2 data, Speech Act does not interact with PP2 Status, which is again supportive of a non-prosodic explanation for the motivating intuition. That there is a robust effect of PP2 Status is reassuring evidence that prosody is sensitive to syntax, and that the study's item construction is motivating the intended parse.

### 3.2.5 On Reading 1 delay

Reading 1 (R1) delay is the amount of time between the initial display of a sentence and the start of phonation. Participants' median R1 delay ranged from 0.6s to 1.6s with a standard deviation of 0.25s. The distribution of R1 delay was notably different than that of R2 delay, shown in Figure 3.3 which indicates that participants were adhering to the protocol at least most of the time.

As a way of analyzing the protocol, and the extent to which participants performed as expected, participants were categorized based on their median R1 delay. In what



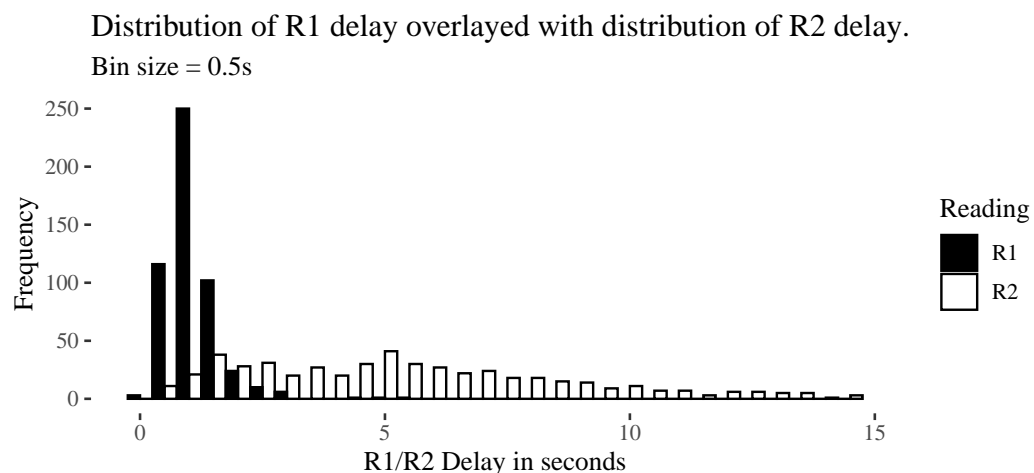


Figure 3.3: Distributions of R1 delay and R2 delay

follows, a fast median R1 delay was shorter than or equal to 0.9s, and a slow one was longer than 1.05s, resulting in 12 participants per category. Ten participants had R1 delays between those values, categorized as “normal,” and set aside. The calculations for categorizing participants were done over Reading 1 of experimental items ( $n = 489$ ). Note that while R1 delay category (i.e., fast or slow) is a property of R1 delay, data for both readings is nonetheless explored within these categories.

There is a statically significant difference between the number of cases where both breaks were produced across the fast (44) vs. slow (65) category for Reading 1 ( $\chi^2(1) = 4.05$ ,  $p < 0.05$ ), but not for Reading 2 ( $\chi^2(1) = 1.86$ ,  $p = 0.17$ ). There was also a statically significant difference in the occurrence of both breaks for Arg compared to Mod sentences across PP2-Status for Reading 2 within the slow R1 delay category ( $\chi^2(1) = 3.97$ ,  $p < 0.05$ ), but comparisons across other factors represented in Table ?? did not yield significant results. The maximum count per cell in the table is 96 (12 participants per category, 8 items per PP2 Status), ignoring missing items.

In cases where R1 delay was small, readers were more likely to produce only one break (PP1 or OBJ) than if R1 delay was fast. A possible explanation is that for

[15]: big ugly table has been omitted

[16]: table 4.15 omitted

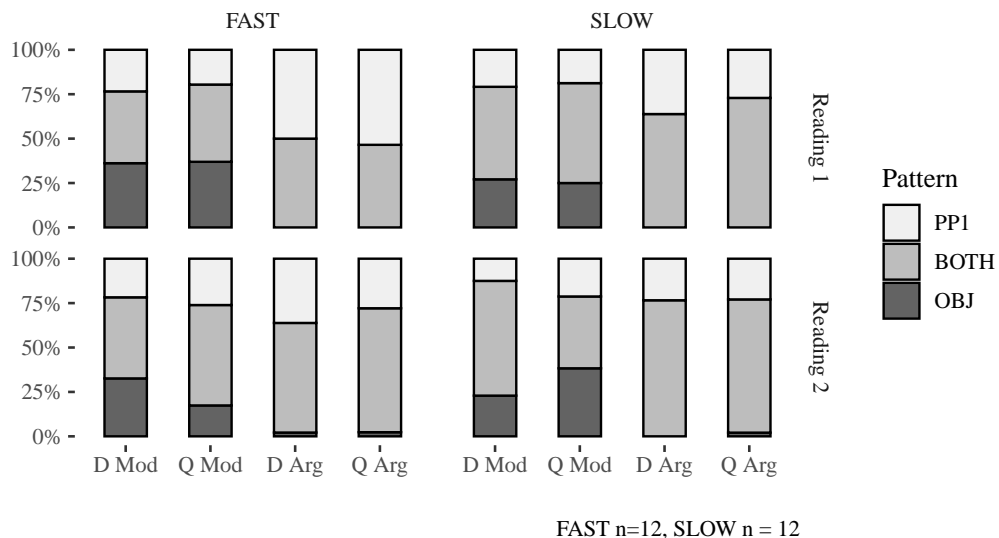


Figure 3.4: Plot of pattern proportions as a function of sentence type.

recordings in the slow category starters readers were more prone to hesitation in general, or perhaps, contrary to instructions, were using both the delay time as well as and their extra break time created via hesitation to look ahead. The significant effect of PP2 Status for the slow category starters in Reading 2 could be due to the readers slow starters not fully understanding the Arg sentences prior to Reading 2, thus increasing their likelihood of hesitation; i.e., where the fast category represents confident readers, the slow category represents less confident readers. As a result of the reader's increased difficulty in Reading 1 for the recordings in the slow category, the effect of processing difficulty is not limited to Reading 1 but in fact spills over into Reading 2. The slow category represents readers who are never able to fully comprehend the sentence, and thus hesitate more frequently.

### 3.3 Discussion of prosodic break patterns

Throughout the analysis of break patterns, PP2 Status was the most robust predictor of OBJ and PP1 break occurrence and their relative strengths (see Section ??). The

OBJ break was more frequent and more frequently dominant for sentences with a PP2 that was an argument than those with a PP2 that could be a modifier; conversely, the PP1 break was more frequent and more frequently dominant for sentences with a PP2 that was interpretable as a modifier than those with a PP2 that was an argument.

Returning to the predictions discussed in ??, recall the contrast made there between linguistically motivated prosodic breaks and other breaks or pauses. The Arg and Mod sentences differ with regard to where the linguistically motivated breaks fall. These expectations are shown in (49) and (50), where ”|” a less prominent or no break and ”||” represents a linguistically motivated prosodic break.

(49)

...	stick	the letter	OBJ Break 	in the mailbox	PP1 Break 	of the proper stack
		Direct object		PP1		vice president.

(50)

...	stick	the letter	OBJ Break 	in the mailbox	PP1 Break 	onto the proper stack.
		Direct object		PP1		PP2

The break patterns observed in the data just discussed mostly do reflect the expectations laid out in (49) and (50): PP1 is the dominant break in a majority of Arg cases for both Readings, and OBJ is dominant in less than half of Arg cases for both Readings (see Figure 3.2 above). Conversely, for Mod cases, OBJ is the dominant break in a majority of cases and PP1 in less than half, across both Readings.

That Reading 2 is a significant predictor of both OBJ break dominance and PP1 break dominance in 3 of the 4 analyses where its inclusion is possible (see Section ??) supports, at least provisionally, hypotheses 1 and 2 from Section ??, repeated below.

- (51) *Hypothesis 1* \\* A first reading of a sentence where PP2 is a goal argument (Arg) will exhibit less natural prosody (more hesitation at and within the PP2 region) than:

- a. A first reading of a sentence where PP2 is a modifier (Mod)
- b. A second reading of a sentence where PP2 is a goal argument (Arg)

(52) *Hypothesis 2* \\* A first reading of a sentence where PP2 is a goal argument (Arg) will more often be produced with prosodic structure that represents an implausible or ungrammatical parse of the string (PP2 incorrectly attached as a modifier), whereas a second reading of that sentence will more often be pronounced with the prosodic structure that represents the intended parse (argument attachment of PP2).

The present study did not specifically attempt to assess~~It is difficult to know~~ whether a given reading represents more or less natural prosody for these constructions,~~but g~~Given that there is a difference between readings, it seems most likely that Reading 2 is the more natural of the two since it represents a considered reading, rather than a ~~hurried~~ one without as much preview. *Hypotheses 1-2* are supported only on their~~that~~ assumption that this is so~~is~~ accepted.

(53) *Hypothesis 3* Reading 1 of a declarative sentence with an argument-PP2 will exhibit less natural prosody (more hesitation at and after the disambiguating region) and be more likely to be produced with prosodic structure that represents an implausible or ungrammatical parse of the string than a Reading 1 of an interrogative sentence with an argument-PP2.

*Hypothesis 3* is not supported by the evidence just reported: there is no statistically significant interaction between Speech Act and PP2 Status for any of the prosodic patterns. A possible explanation is that~~t~~It is surprising that the effect of PP2 Status is generally lessened in Reading 2 when compared to Reading 1, but this can likely be explained as an epi-phenomenon. There is no way to distinguish between prosodic breaks that are intentional and syntactically motivated as compared to those that represent hesitation, a need for a breath, or other factors. It is likely that

some of the effect of PP2 Status is actually an increase in hesitation after PP1, and therefore more or longer pauses at that position, which is mitigated in Reading 2. If some readers are, in general, simply producing a break after every phrase, but happen to produce what is perceived as a dominant break after PP1 for the Arg sentences when they are confused, that increase in the dominance of the PP1 break as an effect of PP2 Status will go away in Reading 2 once they have had time to figure the sentence out. This might mean that the noise caused by readers that are simply breaking phrase-by-phrase is actually amplified in Reading 2.

That a prosodic break also frequently occurs between phrases when not linguistically motivated, i.e., the PP1 break in Mod versions of sentences or the OBJ break in Arg versions of sentences, ~~there is no change in branching direction~~ is mitigated somewhat by the fact that such breaks are usually weaker than the ones that do represent such a change. It is likely that these breaks are actually there for non-syntactic reasons; the end of a phrase represents a reasonable time for the speaker to take a breath or pause briefly for phonological length~~processing~~ reasons. It is also likely that some readers are simply producing a break after each phrase as a strategy for dealing with difficult to comprehend sentences.

Speech act is a significant predictor of PP1 break dominance ( $\beta=0.31$ , std. error = 0.15,  $p < 0.05$ , see Section ??), but not of any of the other prosodic outcomes ~~just discussed~~. It is plausible that the PP1 break is more likely to be dominant in questions than in declaratives because of the need to begin the sentence final rise of question intonation. That there is no observed~~never an~~ interaction between Speech Act and PP2 Status is discouraging for the hypothesis that what it's the \*prosody\* of questions that makes the Arg cases seem easier in the interrogative cases compared to the declarative \added{is the prosody of questions}.

## 3.4 Inter-reading time

Inter-reading time is the amount of time after the completion of Reading 1 and before the beginning of phonation of Reading 2. The details of how this was measured and defined can be found in section 2.7. IRT is meant to provide an estimate~~be a measure to some extent~~ of how much difficulty the reader has in processing a given sentence. If a reader spends more time studying a sentence prior to reading it aloud a second time, the IRT will be longer, which can be taken as~~I take that~~ an indicator of processing load on the first reading.

Since~~Importantly~~, IRT is a measure across pairs of recordings (Reading 1/Reading 2), so the number of data analyzed in this section are half as many as those in the prosody analyses.

### 3.4.1 Data cleanup

IRTs below 0.25s ( $n = 2$ ) and above 25.0s ( $n = 5$ ) were assumed to be implausible and omitted from the analyses reported below. Experimental data were then Winsorized by participant to bring data below the 2.5% and above the 97.5% threshold to the value at those thresholds. The resulting measure is referred to as wIRT and is distributed as shown in figure 3.5 ( $n = 489$ ).

Overall mean for wIRT was 6.5s ( $sd = 3.8$ ). The longest wIRT was 22.2s and the shortest was 0.7s. Median wIRT was 6.1s.

### 3.4.2 Analysis of IRT data

Figure 3.6 shows the mean IRT as a function of sentence type.

The two slopes are only very slightly divergent. Notably, both Speech Act and PP2 Status appear to have main effects on wIRT, with interrogatives attracting longer

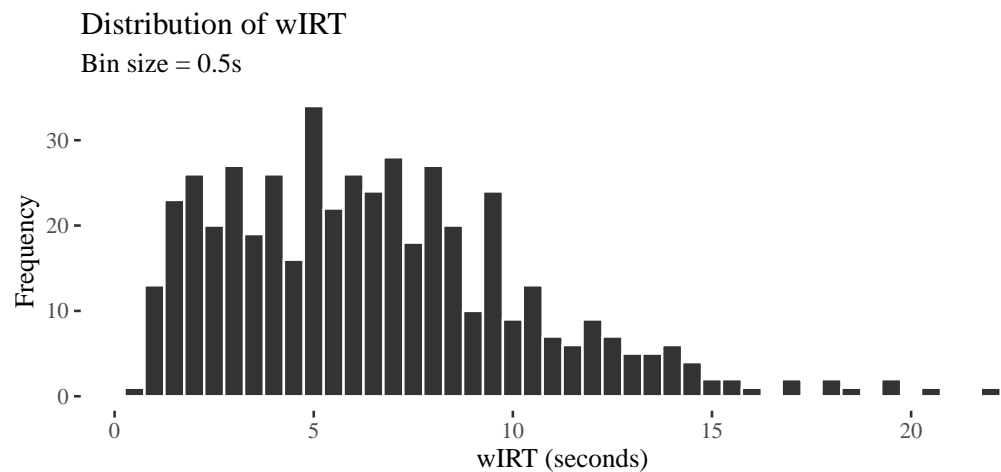


Figure 3.5: Distribution of wIRT.

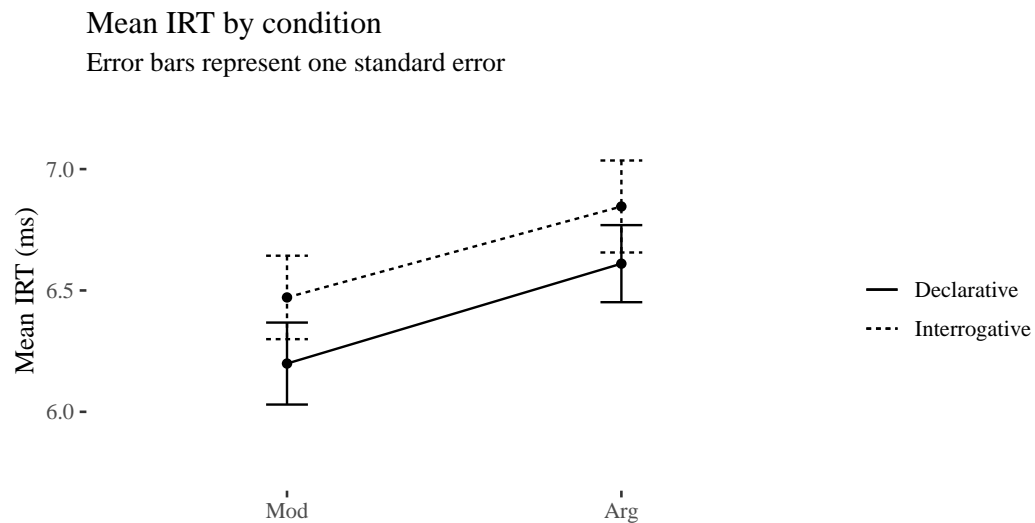


Figure 3.6: Mean IRT as a function of sentence type.