

Computational models of acquisition for islands

Lisa Pearl & Jon Sprouse

1. Why look at language acquisition?

Though it is not always directly stated, the debate at the center of this volume is in many ways driven by language acquisition considerations. Long-distance dependencies are themselves relatively complex, as they involve context-sensitive grammatical operations (e.g., *wh*-movement or slash-passing). The existence of context-sensitive operations alone increases the complexity of the hypothesis space of possible grammars that must be considered by children during the acquisition process. If island effects are indeed the result of grammatical constraints, then the hypothesis space increases yet again, as the grammar must also contain complex constraints on context-sensitive operations. A common hypothesis in the generative syntax literature is that this level of complexity (constraints on context-sensitive grammatical operations) cannot be learned directly from the input that children receive (i.e., this is a poverty of the stimulus problem). As such, many generative syntacticians have postulated the existence of innate domain-specific knowledge about the form that such constraints must take. In other words, the grammatical approach to island effects is often correlated with a nativist, or Universal Grammar (UG) based, view of language acquisition. In this way, a reductionist approach to island effects could be seen as a type of simplifying approach to the grammar, as it could eliminate the need for one set of innate constraints on the shape of human grammars. Because of this, it seems to us that discussions of “parsimony” and “simplification” in the

reductionist literature either directly or indirectly concern the presumed problem that occurs during language acquisition.

Given the amount of research that has been conducted on the debate between grammatical and reductionist approaches to island effects, it seems important at this stage to determine exactly what type of innate knowledge (if any) would be necessary to learn the grammatical constraints that give rise to island effects, given the input that children receive during language acquisition. Such an investigation will help determine exactly what is at stake in this debate. If grammatical island constraints cannot be learned from the input available to children without innate domain-specific knowledge (UG), then this debate has direct implications for the language acquisition process. However, if grammatical island constraints can be learned from the input available to children without UG-like knowledge, then this debate is simply one empirical question among the hundreds that must be answered in order to have a complete theory of language.

In this chapter, we examine child-directed speech input in order to formalize the apparent induction problem that has been claimed by linguists. We then explore a statistical learning model of island constraints that is based upon the frequency of certain abstract structures in the input. The model is tested on input derived from child-directed speech (from CHILDES: MacWhinney (2000)) as well as input derived from adult-directed speech (Switchboard section of Treebank-3: Marcus et al. 1999) and adult-directed text (Brown section of Treebank-3: Marcus et al. 1999). We use this statistical model to investigate the types of learning biases that are necessary to learn these constraints from the input, with the goal of determining whether any innate domain-specific biases (i.e., UG) are necessary. Our results suggest that a learner only requires the following biases to learn

syntactic island constraints from child-directed input, none of which are necessarily specific to the nativist/UG approach to language acquisition, though they do raise difficult questions about how these particular biases arise in the learner (see also Pearl and Sprouse *submitted*):

- (i) perceive the input with a phrase-structure-based representation of sentences (i.e., a parser)
- (ii) characterize dependencies as sequences of phrase structure nodes
- (iii) track the frequency of sequences of three phrase structure nodes (trigrams of phrase structure nodes), and their associated probability of occurring
- (iv) construct a longer dependency by combining trigrams of phrase structure nodes, and assess that dependency's grammaticality based on that combination

The fact that syntactic island constraints can potentially be learned from realistic child-directed and adult-directed input without any clearly nativist/UG-specific abilities suggests that the grammatical versus reductionist debate has no implications for the debate between nativists and non-nativists, but is instead just one question among many required to fully understand the human language system.

2. The induction problem

Investigating the learning of syntactic island effects requires a formally explicit definition of the target state beyond the asterisks/no-asterisks that are typically used to delineate

unacceptable sentences in syntactic articles. To that end, we decided to explicitly construct the target state from data from Sprouse et al. (2012), who collected formal acceptability judgments for four island types using the magnitude estimation task: Complex NP islands (1), Subject islands (2), Whether islands (3), and Adjunct islands (4). Sprouse et al. (2012) used a factorial definition of island effects for each island type (see Sprouse (*this volume*) for discussion of the value of the factorial definition of island effects). For our purposes, this simply means that each island type was defined by four sentence types (4 island types x 4 sentence types = 16 sentence types). An example of each sentence type and the resulting container node sequence is given in (1)–(4): (a) MATRIX gap, NON-ISLAND structure, (b) EMBEDDED gap, NON-ISLAND structure, (c) MATRIX gap, ISLAND structure, (d) EMBEDDED gap, ISLAND structure.

(1) Complex NP islands

- | | | |
|----|--|-----------------------|
| a. | Who _ claimed that Lily forgot the necklace? | MATRIX NON-ISLAND |
| b. | What did the teacher claim that Lily forgot _? | EMBEDDED NON-ISLAND |
| c. | Who _ made the claim that Lily forgot the necklace? | MATRIX ISLAND |
| d. | *What did the teacher make the claim that Lily forgot _? | EMBEDDED ISLAND |

(2) Subject islands

- | | | |
|----|---|-----------------------|
| a. | Who _ thinks the necklace is expensive? | MATRIX NON-ISLAND |
| b. | What does Jack think _ is expensive? | EMBEDDED NON-ISLAND |
| c. | Who _ thinks the necklace for Lily is expensive? | MATRIX ISLAND |
| d. | *Who does Jack think the necklace for _ is expensive? | EMBEDDED ISLAND |

(3) Whether islands

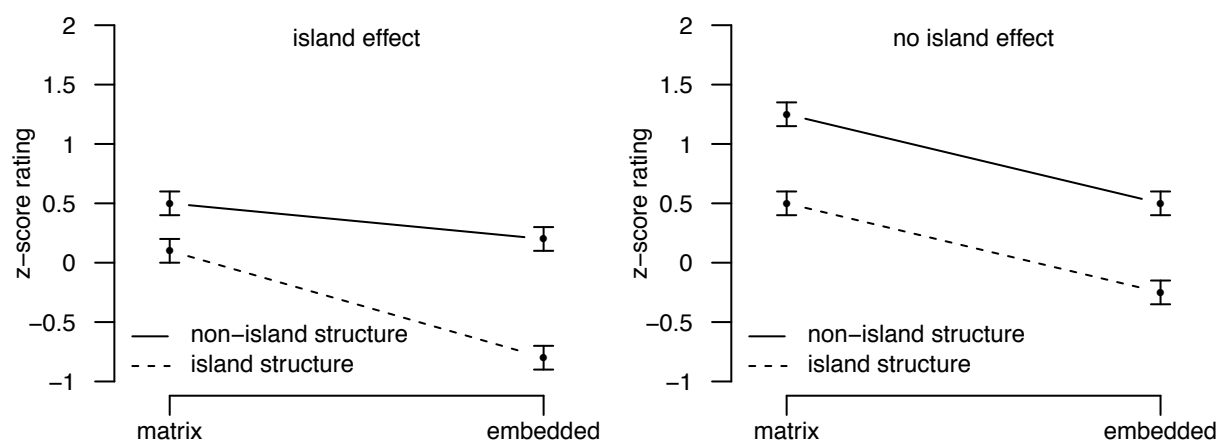
- | | | |
|----|---|-----------------------|
| a. | Who __ thinks that Jack stole the necklace? | MATRIX NON-ISLAND |
| b. | What does the teacher think that Jack stole __ ? | EMBEDDED NON-ISLAND |
| c. | Who __ wonders whether Jack stole the necklace? | MATRIX ISLAND |
| d. | *What does the teacher wonder whether Jack stole __ ? | EMBEDDED ISLAND |

(4) Adjunct islands

- | | | |
|----|---|-----------------------|
| a. | Who __ thinks that Lily forgot the necklace? | MATRIX NON-ISLAND |
| b. | What does the teacher think that Lily forgot __ ? | EMBEDDED NON-ISLAND |
| c. | Who __ worries if Lily forgot the necklace? | MATRIX ISLAND |
| d. | *What does the teacher worry if Lily forgot __ ? | EMBEDDED ISLAND |

The factorial definition of island effects makes the presence of an island effect visually salient: If we plot the acceptability of the four sentence types in a configuration known as an interaction plot, the presence of an island effect shows up as two non-parallel lines, which indicates a statistical interaction of the two factors in the definition (the left panel of Figure 1); the absence of an island effect shows up as two parallel lines, which indicates no interaction of the two factors in the definition (the right panel of Figure 1).

Figure 1: Example graphs showing the presence (left panel) and absence (right panel) of island effects using the factorial definition (see also Sprouse (*this volume*)).



Sprouse et al. (2012) found that adult judgments demonstrated an island effect for all four island types, which means that knowledge of these syntactic islands is indeed necessary to acquire.

To assess a child's input for constraints on *wh*-dependencies (and, specifically, the data in the input directly relevant for generating the judgments in Sprouse et al. 2012), we examined child-directed speech samples to determine the frequency of the structures used as experimental stimuli in Sprouse et al. (2012). While the CHILDES database has many corpora that are annotated with syntactic dependency information (Sagae, Davis, Lavie, MacWhinney, & Wintner, 2010), it is difficult to automatically extract the kind of *wh*-dependency information we needed to identify. For this reason, we selected five well-known corpora of child-directed speech from the CHILDES database (MacWhinney, 2000) to annotate with phrase structure tree information: the Adam, Eve, and Sarah corpora from the Brown data set (Brown, 1973), the Valian corpus (Valian 1991), and the Suppes corpus (Suppes 1974). We first automatically parsed the child-directed speech utterances using a

freely available syntactic parser (the Charniak parser¹), yielding the basic phrase tree structures. However, due to the conversational nature of the data, there were many errors. We subsequently had the parser's output hand-checked by two separate annotators from a group of UC Irvine undergraduates who had syntax training, with the idea that errors that slipped past the first annotator would be caught by the second.² We additionally hand-checked the output of our automatic extraction scripts when identifying the frequency of *wh*-dependencies used as experimental stimuli in Sprouse et al. (2012) in order to provide a third level of error detection.

The data from these five corpora comprise child-directed speech to 25 children between the ages of one and five years old, with 813,036 word tokens total. In all the utterances, 31,247 contained a *wh*-word and a verb, and so were likely to contain a syntactic dependency. Table 1 shows the number of examples found containing the structures and dependencies examined in Sprouse et al. (2012).

Table 1. The corpus analysis of the child-directed speech samples from CHILDES, given the experimental stimuli used in Sprouse et al. (2012) for the four island types examined. The syntactic island condition (which is ungrammatical) is italicized.³

¹ Available at <ftp://ftp.cs.brown.edu/pub/nlparser/>.

² This work was conducted as part of NSF grant BCS-0843896, and the parsed corpora are available at <http://www.socsci.uci.edu/~lpearl/CoLaLab/TestingUG/index.html>.

³ Note that the number of MATRIX | NON-ISLAND data are identical for all four island types since that control structure was identical for each island type (a *wh*-dependency linked to the subject position in the main clause, with the main clause verb (e.g., *thinks*) taking a tensed subordinate clause (e.g., *Lily forgot the*

	MATRIX NON-ISLAND	EMBEDDED NON-ISLAND	MATRIX ISLAND	EMBEDDED ISLAND
Complex NP	7	295	0	0
Subject	7	29	0	0
Whether	7	295	0	0
Adjunct	7	295	15	0

From Table 1, we can see that these utterance types are fairly rare in general, with the most frequent type (LONG | NON-ISLAND) appearing 0.9% of the time (295 of 31,247). Secondly, we see that being grammatical doesn't necessarily mean an utterance type will occur in the input. Specifically, while both the MATRIX | NON-ISLAND and MATRIX | ISLAND utterance types are grammatical, they rarely occur in the input (7 for MATRIX | NON-ISLAND, 15 for the Adjunct MATRIX | ISLAND type). This is problematic from a learning standpoint, if a learner is keying grammaticality directly to input frequency. Unless the child is very sensitive to small frequency differences (even 15 out of 31,247 is less than 0.05% of the relevant input), the difference between the frequency of grammatical MATRIX | ISLAND or MATRIX | NON-ISLAND utterances and that of ungrammatical EMBEDDED | ISLAND utterances is very small for Adjunct island effects. It's even worse for Complex NP, Subject, and Whether island effects, since the difference between grammatical MATRIX | ISLAND utterances and ungrammatical EMBEDDED | ISLAND structures is nonexistent. Since neither utterance type

necklace)). Similarly, the number of EMBEDDED | NON-ISLAND data are identical for Complex NP, Whether, and Adjunct islands since that control structure was identical for those island types (a *wh*-dependency linked to the object position in the embedded clause, with the main clause verb taking a tensed subordinate clause).

appears in the input, how would this learner classify one as grammatical and the other ungrammatical? Thus, it appears that child-directed speech input presents an induction problem to a learner attempting to acquire adult grammatical knowledge about syntactic islands.

The existence of an induction problem then requires some sort of learning bias in order for children to end up with the correct grammaticality judgments. We note that this induction problem arises when we assume that children are limiting their attention to direct evidence of the language knowledge of interest (something Pearl & Mis (submitted) call the *direct evidence assumption*) – in this case, utterances containing *wh*-dependencies and certain linguistic structures. One useful bias may involve children expanding their view of which data are relevant (Foraker et al., 2009; Pearl & Mis, 2011; Perfors, Tenenbaum, & Regier, 2011), and thus including *indirect positive evidence* (Pearl & Mis, submitted) for syntactic islands in their input. We explore this option in the learning algorithm we describe in the next section.

3. A statistical learning algorithm for syntactic islands

Though there appears to be an induction problem for syntactic islands, children clearly must utilize some learning procedure in order for them to become adults who have the acceptability judgments observed in Sprouse et al. (2012). The essence of the acquisition process involves applying learning procedures to the available input in order to produce knowledge about language (Niyogi & Berwick 1996, Yang 2002, among many others). Pearl

& Lidz (2009) suggest that the complete description of the acquisition process must contain at least the following:

- (i) a specification of the child's representation(s) of the hypothesis space
- (ii) a representation of the input that is available to children (the *intake* (Fodor 1998a))
- (iii) the updating procedure that is used to navigate the hypothesis space

In a modeled learner, we can (and must) precisely specify each component of the acquisition process, including whether a bias is present and what the bias does to the hypothesis space, the input, and/or the update procedure. For example, almost all theories assume that children must have a bias to represent their hypotheses about linguistic structures as abstract phrase structure trees. Nativist/UG-based theories may go even further and assume an even more abstract hypothesis space, perhaps in the form of primitives necessary for innate syntactic constraints (e.g., bounding nodes for the Subjacency condition (Chomsky 1973)). Similarly many theories assume that children have a bias to use probabilistic reasoning to update their beliefs about which structures are grammatical (e.g., Tenenbaum & Griffiths 2001, Griffiths & Tenenbaum 2005, Gerken 2006, Xu & Tenenbaum 2007, Frank et al. 2009). Nativist/UG-based theories may again go even further by assuming that a single occurrence of a given structure is enough to instantiate a given grammar (e.g., triggers (Lightfoot 1991, Gibson & Wexler 1994, Niyogi & Berwick 1996, Fodor 1998a, Dresher 1999, Lightfoot 2010, among others)). Formally modeling these allows us to see the effect of any given learning bias on acquisition, and determine which biases are necessary. Once we have that, we can then investigate the nature of the

necessary biases to determine if they qualify as unique to nativist/UG-based approaches to acquisition, or are shared by non-nativist theories of acquisition.

We will use the three components mentioned above to organize the presentation of our learning algorithm, albeit in a slightly different order: the representation of the input, the representation of the hypothesis space given the input, and the updating procedure given the input. We describe the performance of this learning strategy based on realistic input in section 4. We postpone discussion of the nature of the components of the learning strategy until section 5.

3.1 The representation of the input

Turning first to the input representation, we suggest that children may be tracking the occurrence of structures that can be derived from phrase structure trees. To illustrate, the phrase structure tree for “Who did she like?” can be represented with the bracket notation in (5a), which depicts the phrasal constituents of the tree. We also assume that the learner can extract one crucial piece of information from this phrase structure tree: all of the phrasal nodes that dominate the gap location, which we will metaphorically call its “container nodes.” A simple way to identify the container nodes is simply those phrasal constituents currently unclosed (opened with a left bracket), given the understood position of the dependencies. Since container nodes play an integral role in all syntactic formulations of island constraints, they therefore seem like a necessary starting point for constructing such constraints. Furthermore, the sentence-processing literature has repeatedly established that the search for the gap location is an active process (Crain & Fodor 1985, Stowe 1986) that tracks the container nodes of the gap location (see Phillips

2006 for a list of real-time studies that have demonstrated the parser’s sensitivity to island boundaries). In this way, our assumption that the learner can in principle extract this information from the phrase structure trees is actually a well-established fact of the behavior of the human sentence parser (though there is a difference between having access to information and actually using that information, which we discuss in detail in section 5). For (5a), the container nodes would be the sequence in (5b), where the gap location of the displaced NP *who* is dominated by the matrix VP and then the matrix IP. We can represent this dominance information as a sequence of container nodes, as in (5c). Another example is shown in (6a)-(6c), with the utterance “Who did she think the gift was from?” Here, the gap position of the displaced NP *who* is dominated by several nodes (6b). This can be represented by the container node sequence in (6c).

(5) a. [CP Who did [IP she [VP like [NP _]]]]?

b. IP VP

c. IP-VP

(6) a. [CP Who did [IP she [VP think [CP [IP [NP the gift] [VP was [PP from _]]]]]]]]?

b. IP VP CP IP VP PP

c. IP-VP-CP-IP-VP-PP

In order to track container node sequences, the learning algorithm must also specify the set of possible container nodes. For the current algorithm, we assume phrase structure nodes that are relatively universal across syntactic theories (e.g., NP, VP, IP, CP). However,

the definition of island effects and the corpus study in section 2 make it clear that CP nodes must be subcategorized in order to successfully learn syntactic islands. For example, without subcategorizing the CP node, the container node sequence for the grammatical EMBEDDED | NON-ISLAND sentence in the Whether island design would be identical to the ungrammatical EMBEDDED | ISLAND condition: IP-VP-CP-IP-VP. In order to separate these two conditions, the algorithm must track the lexical item that introduces the CP (*that* versus *whether*): IP-VP-CP_{that}-IP-VP versus IP-VP-CP_{whether}-IP-VP. This is an empirical necessity; however, we discuss potential empirical motivation for this assumption in section 5.

3.2 The representation of the hypothesis space

Given this input representation, we propose that the hypotheses concern which container node sequences are grammatical and which are not. That is, one hypothesis might be something like “The container node sequence IP-VP is grammatical”. Children’s acquisition then consists of assigning some probability to each hypothesis, explicitly or implicitly. We propose a learning algorithm below that implicitly assigns a probability to each hypothesis like this, based on the form of the container node sequence. In order to represent the hypothesis space this way, children need only to represent the input in terms of these container node sequences, which comes from being able to parse and track dependencies in a given utterance.

The learning algorithm we propose involves the learner tracking the frequency of smaller sub-sequences of container node sequences, as encountered in the input. In particular, we suggest that a learner could track the frequency of container node trigrams

(i.e., a continually updated sequence of three container nodes) in the input utterances.⁴ For example, the container node sequences from (5c) would be represented as a sequence of trigrams as in (7c), and the container node sequences from (6c) would be represented as a sequence of trigrams as in (8c):

- (7) a. [CP Who did [IP she [VP like [NP _]]]]?
- b. IP VP
- c. start-IP-VP-end =
start-IP-VP
IP-VP-end

- (8) a. [CP Who did [IP she [VP think [CP [IP [NP the gift] [VP was [PP from _]]]]]]]]?
- b. IP VP CP_{null} IP VP PP
- c. start-IP-VP-CP_{null}-IP-VP-PP-end =
start-IP-VP
IP-VP-CP_{null}
VP-CP_{null}-IP

⁴ Note that this means a learner is learning from data containing dependencies besides the one of interest. For example, a learner deciding about the sequence IP-VP-CP_{that}-IP-VP would learn from IP-VP dependencies that the trigram *start-IP-VP* appears. This is an implicit learning bias that expands the relevant intake set of the learner – all dependencies are informative, not just the ones being judged as grammatical or ungrammatical.

CP_{null}-IP-VP

IP-VP-PP

VP-PP-end

3.3 The updating procedure

The learner generates the probability of a given container node trigram based on the observed data. Then, to gauge the grammaticality of any given container node chain (such as an island), the learner calculates the probability of observing that sequence of container node trigrams, which is simply the product of the trigram probabilities.⁵ For example, in (3), the sequence IP-VP would have a probability equal to the product of the trigram *start-IP-VP* and the trigram *IP-VP-end*.

All other things being equal, this automatically makes longer dependencies less probable than shorter dependencies since more probabilities are multiplied together for longer dependencies, and those probabilities are always less than 1. Note, however, that the frequency of the individual trigrams comprising those dependencies still has a large effect. In particular, a shorter dependency that includes a sequence of very infrequent trigrams will still be less probable than a longer dependency that contains very frequent trigrams. Thus, the frequencies observed in the input temper the detrimental effect of dependency length. The learning algorithm and calculation of grammaticality preferences

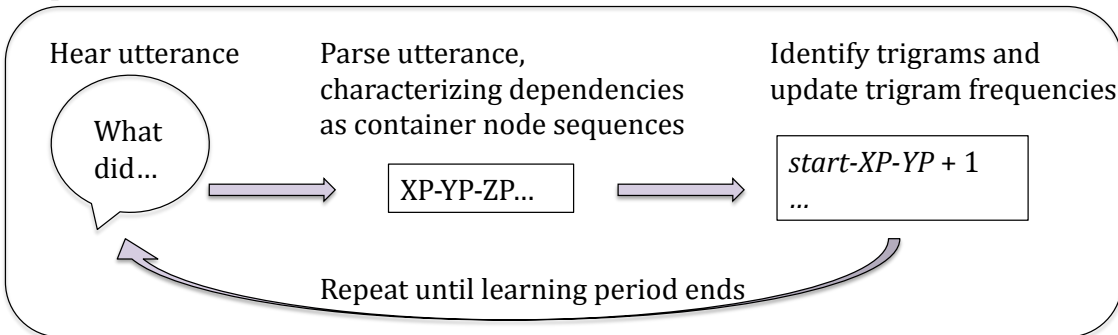
⁵ We note that the learner uses smoothed trigram probabilities (using Lidstone's Law (Manning & Schütze 1999) with smoothing constant $\alpha = 0.5$), so unobserved trigrams have a frequency slightly above 0.

Specifically, the learner imagines that unobserved trigrams have been observed α times, rather than 0 times, and all other trigrams have been observed α + their actual observed occurrences.

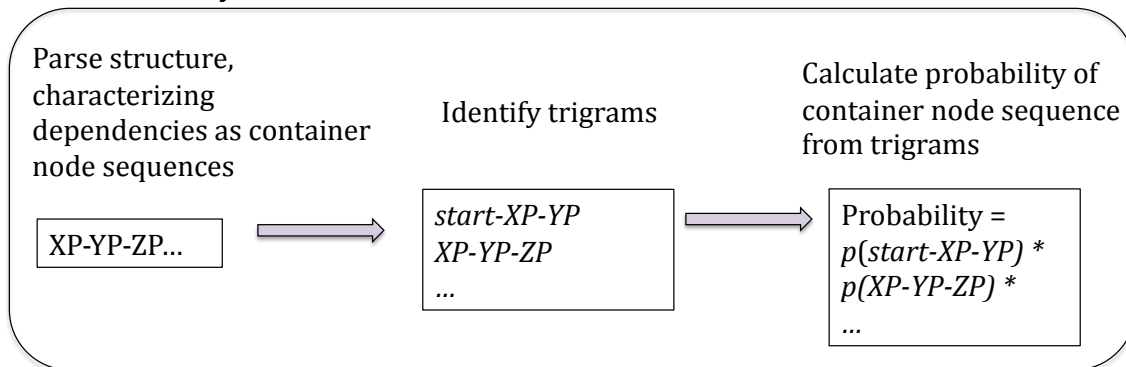
are schematized in figure 2, and two examples of grammaticality preferences are shown in (9) and (10).

Figure 2. Steps in the acquisition process and calculation of grammaticality preferences.

Acquisition Process



Grammaticality Preferences



(9) “Where does he think Jack stole from?”

[_{CP} Where does [_{IP} [_{NP} he] [_{VP} think [_{CP} [_{IP} [_{NP} Jack] [_{VP} stole [_{PP} from _]]]]]]]?”

IP VP CP_{null} IP VP PP

Sequence: start-IP-VP-CP_{null}-IP-VP-PP-end

start-IP-VP

IP-VP-CP_{null}

VP-CP_{null}-IP

CP_{null}-IP-VP

IP-VP-PP

VP-PP-end

Probability(IP-VP-CP_{null}-IP-VP-PP) =

p(start-IP-VP)*p(IP-VP-CP_{null})*p(VP-CP_{null}-IP)*p(CP_{null}-IP-VP)*p(IP-VP-PP)*p(VP-PP-end)

(10) *"Who does Jack think the necklace for is expensive?"

[_{CP} Who does [_{IP} [_{NP} Jack] [_{VP} think [_{CP} [_{IP} [_{NP} the necklace [_{PP} for _]] [_{VP} is expensive]]]]]]]?

IP VP CP_{null} IP NP PP

Sequence: start-IP-VP-CP_{null}-IP-NP-PP-end

start-IP-VP

IP-VP-CP_{null}

VP-CP_{null}-IP

CP_{null}-IP-NP

IP-NP-PP-

NP-PP-end

Probability(IP-VP-CP_{null}-IP-NP-PP) =

p(start-IP-VP)*p(IP-VP-CP_{null})*p(VP-CP_{null}-IP)*p(CP_{null}-IP-NP)*p(IP-NP-PP)*p(NP-PP-end)

Given this learning algorithm, a child can generate a grammaticality preference for a given dependency at any point during learning, based on the input previously observed, by calculating its probability from the frequency of the trigrams that comprise it (see Figure 2). Similarly, a relative grammaticality preference can be calculated by comparing the probabilities of two dependencies' container node sequences. This will allow us, for example, to compare the inferred grammaticality of dependencies spanning island structures versus dependencies spanning non-island structures.

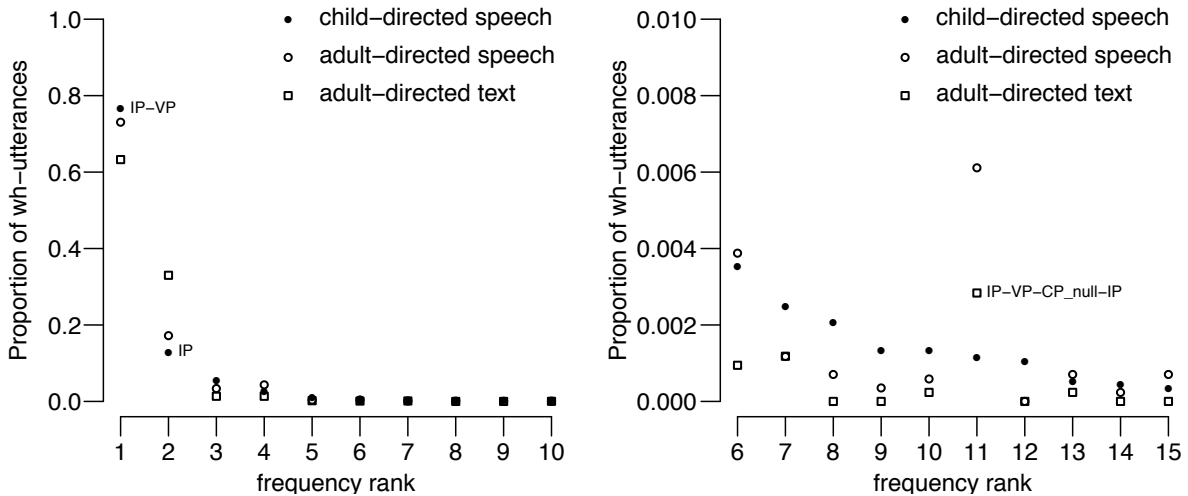
4. Learning about islands from realistic input

We turn now to specific case studies of learning preferences about structural dependencies. First, we consider the input to our modeled learners. If we are modeling how children acquire their grammaticality preferences, we should look at child-directed speech. If we are instead interested in how adults acquire their preferences (perhaps because we have empirical data from adults), then we may be interested in a mix of adult-directed speech and adult-directed text. Table 2 describes the basic composition of three corpora types: child-directed speech from the Adam and Eve corpora from Brown (Brown 1973), the Valian corpus (Valian 1991), and the Suppes corpus (Suppes 1974), adult-directed speech from the Switchboard section of the Treebank-3 corpus (Marcus et al. 1999) and adult-directed text from the Brown section of the Treebank-3 corpus (Marcus et al. 1999). Figure 3 provides a compact representation of the distribution of the types of *wh*-dependencies in each corpus.

Table 2: Basic composition of the child-directed and adult-directed input corpora.

	Child-directed: speech	Adult-directed: speech	Adult-directed: text
total # utterances	101838	74576	24243
total <i>wh</i> -dependencies	20923	8508	4230

Figure 3. The 15 most frequent *wh*-dependency types in the three corpora types. The left panel displays the 10 most frequent *wh*-dependency types for each of the three corpora types, with IP-VP and IP dominating all three corpora types (IP-VP: rank 1, IP: rank 2). The right panel displays the 6th-15th most frequent *wh*-dependency types on a smaller y-axis scale (0-.01) in order to highlight the small amount of variation between corpora types for these dependency types.



Notably, two sequences dominate the input, no matter what the corpus: IP-VP and IP, corresponding to main clause object and main clause subject dependencies, respectively. Interestingly, child-directed speech doesn't seem to differ much from adult-directed speech with respect to the proportional frequency of these two sequences (child-

directed: 78.3%/11.7%, adult-directed (Switchboard): 73.0%/17.2%). Adult-directed written text tends to be biased slightly more towards main clause subject dependencies, though main clause object dependencies are still far more prevalent (IP-VP: 63.3% to IP: 33.0%). Also, we note that overt complementizers are rare in general. This will become relevant when we examine the learned grammaticality preferences for dependencies involving the complementizer *that*.

We can test our modeled learners by comparing their learned grammaticality preferences to empirical data on adult grammaticality judgments available in Sprouse et al. (2012) (see also Sprouse (*this volume*)). Recall that Sprouse et al. (2012) examined four island types, using a factorial definition of island effects for each island type. The resulting container node sequence for each type is given in (11)–(14): (a) matrix gap, non-island structure, (b) embedded gap, non-island structure, (c) matrix gap, island structure, (d) embedded gap, island structure.

(11) Complex NP islands

a.	IP	MATRIX NON-ISLAND
b.	IP-VP-CP _{that} -IP-VP	EMBEDDED NON-ISLAND
c.	IP	MATRIX ISLAND
d.	*IP-VP-NP-CP _{that} -IP-VP	EMBEDDED ISLAND

(12) Subject islands

- | | | |
|----|-------------------------------------|-----------------------|
| a. | IP | MATRIX NON-ISLAND |
| b. | IP-VP-CP _{null} -IP | EMBEDDED NON-ISLAND |
| c. | IP | MATRIX ISLAND |
| d. | *IP-VP-CP _{null} -IP-NP-PP | EMBEDDED ISLAND |

(13) Whether islands

- | | | |
|----|-------------------------------------|-----------------------|
| a. | IP | MATRIX NON-ISLAND |
| b. | IP-VP-CP _{that} -IP-VP | EMBEDDED NON-ISLAND |
| c. | IP | MATRIX ISLAND |
| d. | *IP-VP-CP _{whether} -IP-VP | EMBEDDED ISLAND |

(14) Adjunct islands

- | | | |
|----|---------------------------------|-----------------------|
| a. | IP | MATRIX NON-ISLAND |
| b. | IP-VP-CP _{that} -IP-VP | EMBEDDED NON-ISLAND |
| c. | IP | MATRIX ISLAND |
| d. | *IP-VP-CP _{if} -IP-VP | EMBEDDED ISLAND |

Recall also that the factorial definition of island effects makes the presence of an island effect visually salient: If we plot the acceptability of the four sentence types in a configuration known as an interaction plot, the presence of an island effect shows up as two non-parallel lines, which indicates a statistical interaction of the two factors in the

definition (the left panel of Figure 1); the absence of an island effect shows up as two parallel lines, which indicates no interaction of the two factors in the definition (the right panel of Figure 1).

To evaluate the success of our learners, we can plot the predicted grammaticality preferences in a similar interaction plot: If the lines are non-parallel, indicating an interaction, similar to the graph in the left panel of Figure 1, then the learner has acquired island constraints; if the lines are parallel, indicating no interaction, similar to the graph in the right of Figure 1, then the learner did not acquire island constraints.

To ground the learning period for our modeled learners, we can draw on empirical data from Hart & Risley (1995) and assume children hear approximately 1 million utterances between birth and 3 years of age. If we assume our learners' learning period is approximately 3 years (perhaps between the ages of 2 and 5 years old, if we're modeling children's acquisition), we can estimate the number of *wh*-dependencies they hear out of those one million utterances. Given child-directed speech samples from Adam and Eve (Brown 1973), Valian (Valian 1991), and Suppes (Suppes 1974), and estimating the proportion of *wh*-dependencies given the total number of utterances (20%), we set the learning period to 200,000 data points. So, our learners will encounter 200,000 data points containing dependencies, drawn randomly from a distribution characterized by the corpora in Table 2 and Figure 3.

All our modeled learners will follow the learning algorithm and grammaticality preference calculation outlined in Figure 2. In particular, they will receive data incrementally, identify the container node sequence and trigrams contained in that sequence, and update their corresponding trigram frequencies. They will then use these

trigram frequencies to infer a probability for a given *wh*-dependency, which can be equated to its judged grammaticality – more probable dependencies are more grammatical, while less probable dependencies are less grammatical. Though the inferred grammaticality can be generated at any point during learning (based on the trigram frequencies at that point), we will show results only from the end of the learning period.

Because the result of a grammaticality preference calculation is often a very small number (due to multiplying many probabilities together), we will calculate the log probability. This allows for easier comparison of grammaticality judgments. All of the log probabilities are negative. The more positive numbers (i.e. closer to zero) represent “more grammatical” structures while more negative numbers (i.e., farther from zero) represent “less grammatical” structures.⁶ To make a direct comparison of these log probabilities with acceptability judgments, Figure 4 plots the experimentally obtained judgments for the dependencies from Sprouse et al. (2012), while Figure 5 shows the model-derived log probabilities of the dependencies, based on child-directed input and Figure 6 shows the model-derived log probabilities of the dependencies, based on adult-directed input.

Figure 4: Experimentally derived acceptability judgments for all four island types from Sprouse et al. (2012) (N=173).

⁶ This measurement is similar to *surprisal*, which is traditionally defined as the negative log probability of occurrence (Tribus 1961) and has been used recently within the sentence processing literature (Hale 2001, Jaeger & Snider 2008, Levy 2008, Levy 2011a). Under this view, less grammatical dependencies are more surprising.

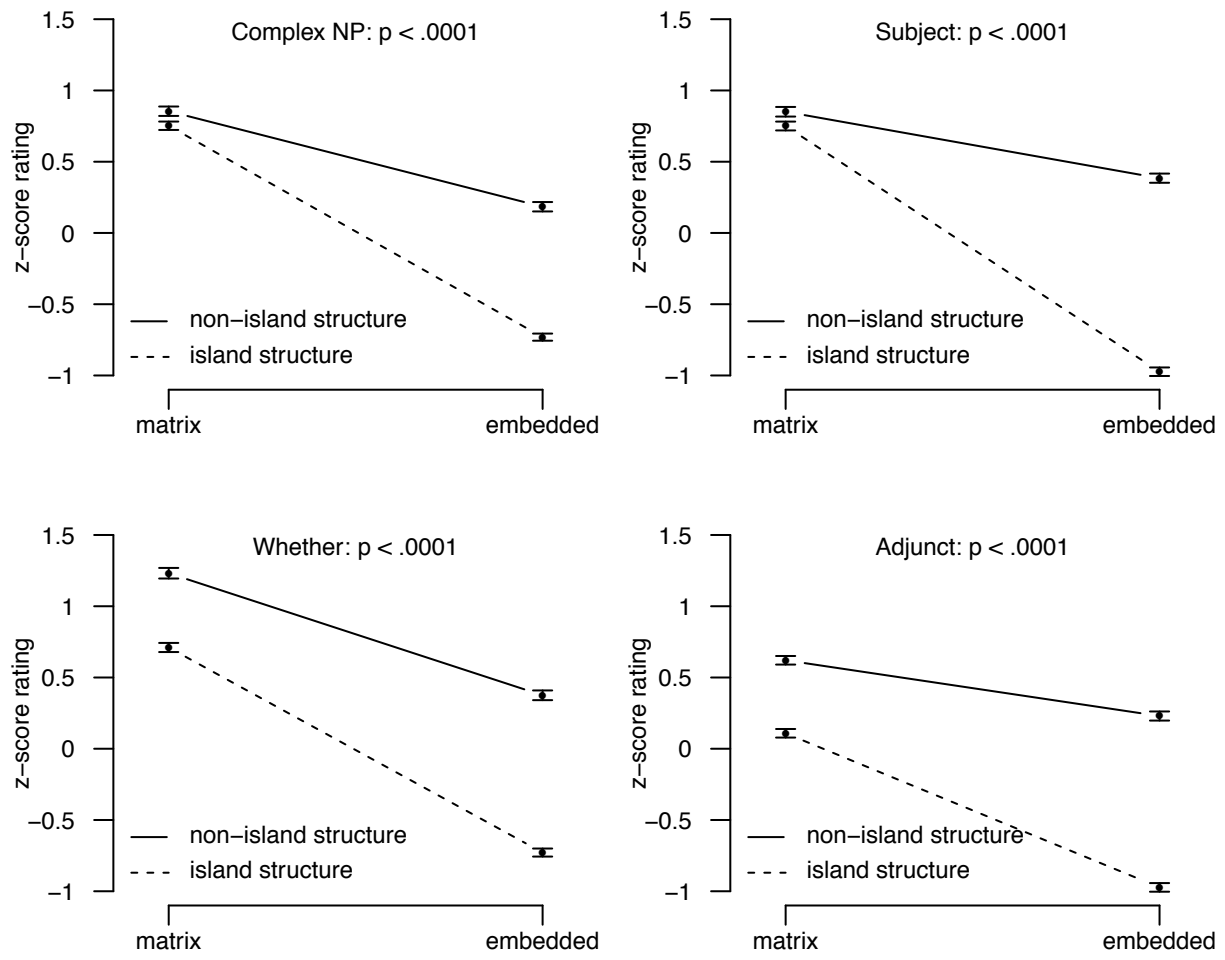


Figure 5: Log probabilities derived from child-directed speech.

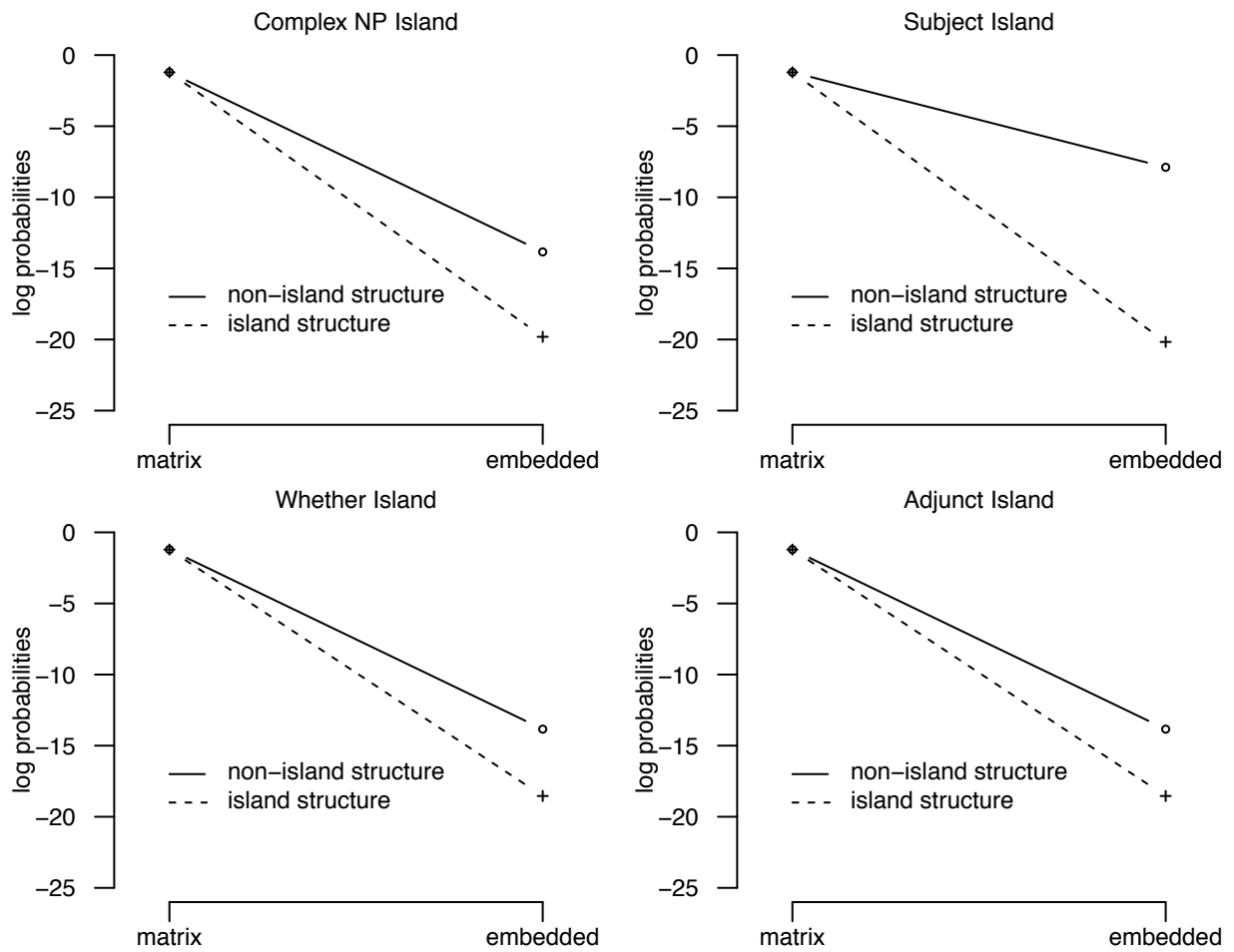
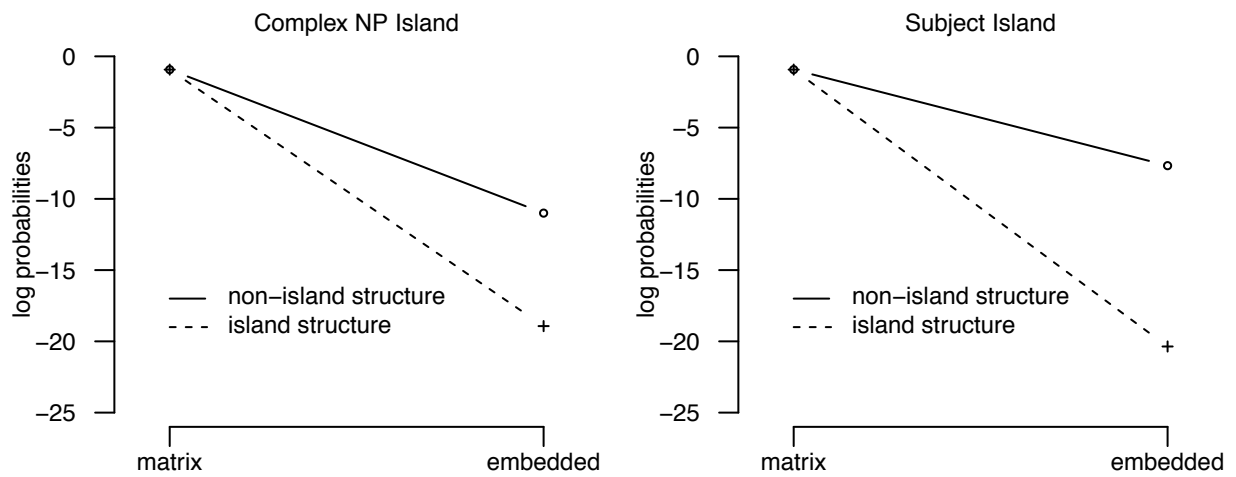
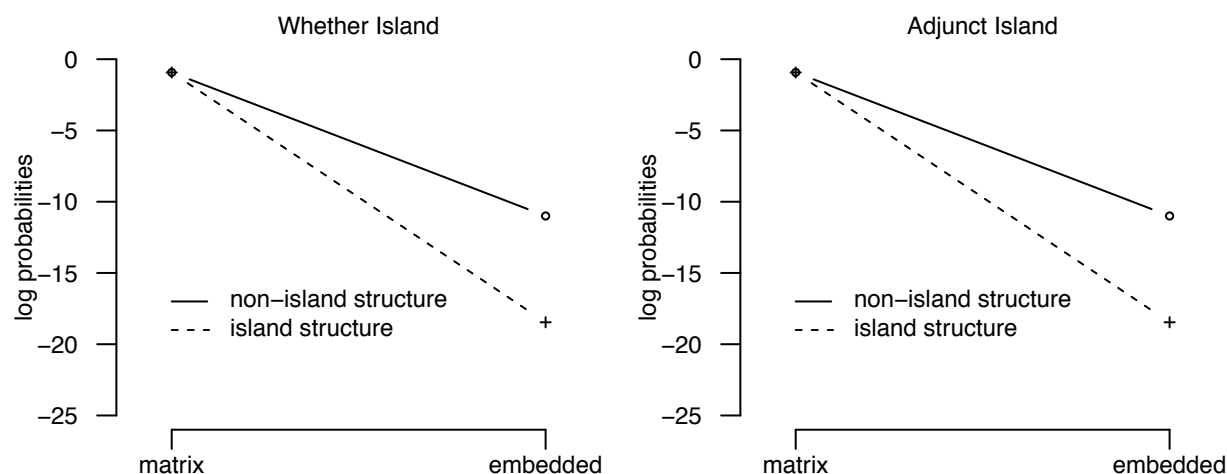


Figure 6: Log probabilities derived from adult-directed speech and text.





We see in Figures 5 and 6 that a learner using either child-directed data or adult-directed data would end up with the correct grammaticality preferences for all four islands (compare these figures to Figure 4).

To sum up, we find that a learner that tracks the probabilities of certain abstract representations of *wh*-dependencies in the input is able to reproduce adult judgments about the (un)grammaticality of islands. In order to capture adult judgments about all four islands investigated, the learning model requires adult-directed input and a certain level of specification in the representation. The proposed algorithm does require relatively sophisticated biases, such as (i) the parsing of sentences into phrase structure trees, (ii) the extraction of sequences of container nodes for the dependencies, (iii) the tracking of the frequency of trigrams of container nodes, and (iv) the calculation of the probability of the complete container node sequence for the dependency, based on its trigrams. In the next section, we discuss the nature of these component biases, and how they might actually arise in the learner.

5. The nature of the necessary learning biases

The question of whether a given learning bias is nativist or non-nativist in nature is actually quite a bit more complex than is often assumed in the syntactic literature. For example, there are at least three dimensions to learning biases that may be relevant (Pearl & Mis 2011, submitted):

- (i) Are they *innate* (and so part of the human biological endowment) or *derived* from prior experience (probably prior experience with language data)?
- (ii) Are they *domain-specific* (and are only used for learning language) or *domain-general* (and are used when learning anything)?
- (iii) Are they about *the hypothesis space* (and so may restrict the learner's hypotheses explicitly) or about *the learning mechanism* (and so may restrict the learner's hypotheses implicitly)?

Clearly, learning biases could involve any logically possible combination of these dimensions. For example, a more abstract representation of linguistic structure could be derived from phrase structure trees, which themselves may be derived from distributional properties of the linguistic input by using probabilistic learning. This might then be classified as a *derived, domain-specific* bias about the representation of *the hypothesis space*. Probabilistic learning, in contrast, might be classified as an *innate, domain-general* bias about *the learning mechanism*. Note that only learning biases that are both *innate* and *domain-specific* are candidates for UG. For example, an explicit constraint against syntactic

islands would be just this kind of bias, since it would be *innate* (it's explicitly built in) and *domain-specific* (it applies only to language). In addition, we could likely classify it as a bias about *the hypothesis space*, since it explicitly constrains the hypothesis space of the learner to exclude islands. Our learning strategy does not use this bias, but, as mentioned above, it does use a number of fairly sophisticated learning biases. We discuss each in turn with a particular focus on (i) the empirical motivation for each bias and (ii) the potential classification of each bias according to the framework above.

5.1 Parsing sentences into phrase structure trees

One of the most basic components of the proposed learning algorithm is that it operates over input that has been parsed into phrase structure trees. In order to represent the input this way, children need the ability to parse and track dependencies in a given utterance. Work by Fodor (Fodor 1998a, Fodor 1998b, Sakas & Fodor 2001, Fodor 2009) suggests that this ability may be useful for learning many different kinds of syntactic structures. This component assumes that both syntactic category information and phrase structure information have already been acquired by the learner (or are in the process of being acquired). We do not have too much to say about this assumption because basic syntactic phenomena like syntactic categories and phrase structure parsing are required by nearly every syntactic phenomenon. We would likely consider this ability to be a learning bias that is *domain-specific* since it applies to language data, and a bias about *the hypothesis space* since it involves representing the input in a particular way. It is possible that the process of chunking data into cohesive units is *domain-general* and *innate* (e.g., parsing visual scenes into cohesive units), though it is possible that the particular units that are being chunked

(i.e., phrasal constituents) can be *derived* from distributional properties of the input (for recent work investigating the acquisition of syntactic categories from child-directed input, see Mintz (2003) and (2006), and for recent work investigating the acquisition of hierarchical structure given syntactic categories as input, see Klein & Manning (2002)). Nonetheless, it may be the case that the acquisition of syntactic categories or phrase structure requires at least one innate, domain-specific bias, in which case every syntactic phenomenon, including syntactic islands, would (strictly speaking) require such a bias. Nonetheless, this would not be a fact that is specific to syntactic islands, but rather a general fact of every syntactic phenomenon. We are specifically interested in the consequences of syntactic islands for learning theories, rather than the consequences of every syntactic phenomenon.

5.2 Characterizing dependencies as sequences of container nodes

Identifying which units are potential container nodes is very important for this learning algorithm to be psychologically plausible. The bias to track sequences of container nodes appears relatively neutral at first glance; after all, syntactic island effects are constraints on dependencies, and therefore the algorithm should track information about the dependencies. However, this raises the question of how it is that the algorithm knows to track container nodes rather than some other piece of information about a dependency (e.g., number of nouns, number of verbs, etc.). It is true, as mentioned in section 3, that the fact that the parsing of long-distance dependencies is an active process means that the sequence of container nodes is information that is likely available to (and salient for) the language system, but *availability* is distinct from *attention*. The current algorithm is biased

to attend to container nodes instead of all of the other logically possible types of information about dependencies that are potentially available. This bias is likely *domain-specific*, as long-distance dependencies (and their constraints) have not been clearly demonstrated in any other domain of cognition. It is also likely a bias about the *hypothesis space*, since it involves the learner characterizing the dependencies in the hypothesis space a particular way. However, it is an open question whether this bias is also innate, or whether it can be derived from other biases. Nonetheless, it seems to be the case that any theory of syntactic islands that postulates a structurally-defined constraint will likely track container nodes, and therefore will be confronted with this difficult question.

In addition to a bias to heed container nodes, the proposed algorithm has a bias to track subcategories of CP based on the lexical item that introduces the CP (e.g., *that*, *whether*, *if*, and the null complementizer). Similar to the container node bias, this is empirically necessary: An algorithm that treats all CPs identically will fail to learn Whether islands and Adjunct islands, because the only difference between Whether and Adjunct violations and their non-island control conditions is in the type of CP (*that* versus *whether*, and *that* versus *if*). Again similar to the container node bias, this raises the question of how the algorithm knows what the proper set of container nodes to track is. It is logically possible to subcategorize any number of maximal projections, or none at all, or even to count intermediate projections (e.g., N') as a container node. The fact that CPs *can* be subcategorized is relatively straightforward. Different CPs introduce different types of clauses, with substantial semantic differences: *that* introduces declarative clauses (which are semantically propositions), *whether* introduces questions (which are semantically sets of propositions), and *if* introduces condition clauses. However, the fact that this type of

information is available to the language system does not explain how it is that the learner knows to pursue this particular strategy (or knows where to draw the line between types of container nodes). It may be possible to capture part of this behavior with *innate, domain-general* preferences for certain types of hypotheses (either more specific hypotheses, such as subcategorize all container nodes, or more general hypotheses, such as subcategorize no container nodes) coupled with a *domain-specific* proposal about the types of information in the *learning mechanism* that could be used to correct mistaken hypotheses. But this simply pushes back the question to one about how the system knows which evidence to look for to correct mistaken hypotheses (i.e., is it innate or derived?). In short, much like the container node bias, the empirical necessity of subcategorizing CPs raises difficult questions for any theory of the acquisition of syntactic islands.

5.3 Tracking the frequency of container node trigrams

The proposed algorithm decomposes the container node sequence into trigrams (a moving window of three container nodes). Once again, this is an empirical necessity: The corpus analysis in section 2 suggests that the learning algorithm must decompose the container node sequences into smaller units, otherwise three of the (grammatical) MATRIX | ISLAND conditions would be erroneously characterized as ungrammatical. Similar to the previous biases, it is an open question how this bias arises. Learning models based on sequences of three units have been proposed and are consistent with children's observable behavior for other linguistic knowledge (e.g., the comparison of three sequential transitional probabilities for word segmentation: Saffran et al. 1996, Aslin et al. 1998, Graf Estes et al. 2007, Pelucchi et al. 2009a, Pelucchi et al. 2009b; frequent frames consisting of three

sequential units for grammatical categorization: Mintz 2006, Wang & Mintz 2008); additionally, these learning models are consistent with human behavior for non-linguistic phenomena (Saffran et al. 1996) and also with learning behavior in non-human primates (Saffran et al. 2008). Given this, such a bias is likely *domain-general* (and clearly about the *learning mechanism*); however, the fact that trigrams are an available option does not explain how it is that the learning algorithm knows to leverage trigrams (as opposed to other n-grams) for syntactic islands.

A more easily solved issue concerns the potential issue of data sparseness that could occur with a trigram model, such that the learner could not possibly hope to have enough input to observe examples of all legal trigrams.⁸ However, that is not likely to be a problem for the learner we propose, since we are constructing trigrams over units much more abstract than individual vocabulary items. If we have fewer than 15 (as we might if we only use IP, VP, NP, PP, AdjP, and CP subtypes as the relevant phrasal constituents), then the number of trigrams children must track is less than 15^3 (3375). This is likely less than the number of vocabulary items children know by the time they would be learning grammaticality preferences about dependency structures, and so doesn't seem particularly taxing for children to track.

5.4 Calculating the probability of a container node sequence based on trigrams

Another basic component of the proposed algorithm is that the learner has the ability to track the frequency of units in the input, and then calculate the probabilities of those units. This is a relatively uncontroversial assumption, as many learning theories, both in language

⁸ Additionally, tracking a huge number of trigrams may strain a learner's memory.

and other cognitive domains, assume that the learner can track frequencies and calculate probabilities. The bias to track frequencies and calculate probabilities is likely an *innate, domain-general* bias about the *learning mechanism*. Still, the interesting question about the ability to track frequencies and calculate probabilities is not so much the existence of the ability itself, but rather the units that are tracked, which we discussed above.

5.5 Learning bias summary

Table 3 summarizes the learning biases required for the proposed acquisition process along the relevant dimensions for the UG debate: domain-specific vs. domain-general, and innate vs. derived. Note that none of the learning biases (or their components) are definitively both innate and domain-specific simultaneously (though some very well could be). If these biases (and their components) turn out not to be both innate and domain-specific, they would then not be part of a nativist/UG-based approach to the acquisition of island constraints. In other words, the learning model that we have constructed here would not be based on any Universal Grammar assumptions.

Table 3. Classification of the learning biases required by the proposed acquisition process. The critical bias types (domain-specific and innate) are shaded to help illustrate the fact that no process in this learning model requires a bias that is clearly both domain-specific and innate simultaneously, though questions still remain about how some of these biases arise in the learner.

Description of process	Domain-specific	Domain-general	Innate	Derived
Parse utterance into a phrase structure tree	*		?	?
Characterize dependency as container node sequence	*		?	?
Identify trigrams & update probability		*	*	
Calculate probability of utterance's dependency		*	*	

6. Discussion & conclusion

In this chapter, we have proposed a statistical model for the acquisition of syntactic constraints on *wh*-dependencies that does not rely on innate, domain-specific knowledge of island constraints. Instead, our psychologically plausible learning model is able to implicitly derive knowledge of islands from the input using a series of relatively uncontroversial assumptions, such as the ability to parse sentences into phrase structure trees, the ability to track the nodes that contain the gap location of a *wh*-dependency, the ability to track the frequency of trigrams of container nodes, and the ability to construct a grammaticality preference for a dependency based on its trigrams. This suggests that children (and adults) do not need innate, domain-specific knowledge about islands, which in turn suggests that explicit constraints against island structures do not have to be part of Universal Grammar. In addition, we find that the learning strategy capable of doing this doesn't even need to involve sophisticated probabilistic inference abilities, such as Bayesian updating (e.g., Feldman et al. 2009, Foraker et al. 2009, Frank et al. 2009, Goldwater et al. 2009, Pearl et al. 2011, Perfors et al. 2011). Instead, the probabilistic learning component is fairly simple and

involves tracking frequencies of particular linguistic representations that are small in size (trigrams of container nodes).

However, these results do raise interesting questions about how feasible this learner would be for the full range of constraints on *wh*-dependencies. Though this statistical model demonstrates that syntactic islands can in principle be learned from child-directed input, this particular model cannot capture certain exceptions to syntactic island constraints, such as *parasitic gap* constructions (Engdahl, 1983). Parasitic gap constructions are *wh*-questions in which the *wh*-word is associated with two gap positions: one gap position occurs in a licit gap location (i.e., not inside a syntactic island) while the other gap position occurs inside a syntactic island. Whereas a single gap within an island structure results in unacceptability (15a and 16a), the addition of another gap outside of the island seems to eliminate the unacceptability (15b and 16b) (see Phillips 2006 for experimentally collected acceptability judgments):

- (15) a. *Which book did you laugh [before reading __]?
b. Which book did you judge ___{true} [before reading ___{parasitic}]?
- (16) a. *What did [the attempt to repair __] ultimately damage the car?
b. What did [the attempt to repair ___{parasitic}] ultimately damage ___{true}?

The two gaps in a parasitic gap construction are often described as the *true gap*, which occurs outside of the island, and the *parasitic gap*, which occurs inside of the island. The name is a metaphorical reference to the fact that the *parasitic gap* could not exist without

the *true gap*, much like a parasite cannot exist without a host. Though there are several structural restrictions on parasitic gap constructions (e.g., the true gap cannot c-command the parasitic gap), there is no constraint on the linear order of the two gaps, as illustrated by (15-16).

We believe the grammaticality of parasitic gap constructions pose a problem for our statistical learner. This is because the probability of the trigram sequence for the dependency between the *wh*-word and the parasitic gap will be the same as the probability of the trigram sequence for the relevant syntactic island violation. In other words, our learner would infer that parasitic gap constructions are ungrammatical. For example, the container node sequences for (15) would be as in (17). The sequence for both the ungrammatical gap in (15a) and the grammatical (parasitic) gap in (15b) are identical, and in fact would be as (un)acceptable as other adjunct islands, such as those using the complementizer *if*.

(17)

- a. *Which book did [IP you [VP laugh [CP without [IP [VP reading ___]]]]]?
 Ungrammatical gap sequence: IP-VP-CP_{without}-IP-VP
- b. Which book did [IP you [VP judge ____{true} [CP without [IP [VP reading ____{parasitic}]]]]]?
 Parasitic gap sequence: IP-VP-CP_{without}-IP-VP

Given that this is not the desired target state, the learning algorithm proposed here is unlikely to be the one children use in practice. However, it may be possible to modify the

learning model to account for these constructions. For example, recent studies demonstrate that the human parser continues to actively search for a second gap even after encountering a licit first gap (Wagers & Phillips, 2009). It could be that the learning algorithm assembles a grammaticality preference based on some kind of aggregation of all container node sequences for gaps in a given utterance. However, unless there is an innate, domain-specific bias to aggregate gap information (which would then make this a UG bias), this would need to be derived from linguistic experience somehow. One way is for children to have experience with multiple gaps associated with the same *wh*-element. In order for this to be true, child-directed input (or adult-directed, if acquisition is relatively late) must contain examples of *wh*-elements associated with multiple gaps, such as examples of parasitic gaps. We are currently examining additional syntactically-annotated child-directed corpora to answer this (and other) questions.

The implications of these findings for the grammar versus reductionism debate are substantial. Many of the reductionist proposals for capturing island effects without grammatical constraints have at their heart the notion that fewer grammatical constraints will lead to “simpler” grammars, and thus less motivation for innate, domain-specific learning biases (i.e., Universal Grammar). However, as we have just seen, syntactic constraints on *wh*-dependencies can be learned in principle from input available to children without the need for innate, domain-specific biases ((i.e., Universal Grammar). Therefore there appears to be little psychological motivation to “simplify” grammatical theories above and beyond the quest for truth in science, which in this case would be the desire to accurately characterizing the grammatical system itself. We believe that this changes the nature of this debate significantly, as the question about the right

characterization of island effects is no longer tied to assumptions about the nature of language acquisition, but is instead simply one question among many that must be answered to arrive at a complete understanding of the human language faculty.

Acknowledgements

This research was supported in part by National Science Foundation grant BCS-0843896 to Lisa Pearl and Jon Sprouse. We have benefited greatly from discussions with audiences at Input & Syntactic Acquisition 2009 and Input & Syntactic Acquisition 2012. All errors remain our own.

References

- Aslin, R., Saffran, J., & Newport, E. 1998. Computation of Conditional Probability Statistics by 8-Month-Old Infants. *Psychological Science* 9(4): 321-324.
- Brown, R. 1973. *A first language: The early stages*. Cambridge, MA: Harvard University Press.
- Chomsky, N. 1973. Conditions on transformations. In *A Festschrift for Morris Halle*, ed. S. Anderson and P. Kiparsky, pp. 237-286. New York: Holt, Rinehart and Winston.
- Crain, S., and J. Fodor. 1985. How can grammars help parsers? In *Natural language parsing: psycholinguistic, computational, and theoretical approaches*, ed. D. Dowty, L. Karttunen, and A. Zwicky, pp. 94–128. Cambridge University Press.

- Dresher, E. 1999. Charting the learning path: Cues to parameter setting. *Linguistic Inquiry* 30: 27-67.
- Engdahl, E. 1980. Wh-constructions in Swedish and the relevance of subadjacency. In J. T. Jensen (Ed.), *Cahiers Linguistiques D'Ottawa: Proceedings of the Tenth Meeting of the North East Linguistic Society*, (pp. 89-108). Ottawa, ONT: University of Ottawa Department of Linguistics.
- Feldman, N., Griffiths, T., & Morgan, J. 2009. The influence of categories on perception: Explaining the perceptual magnet effect as optimal statistical inference. *Psychological Review* 116: 752-782.
- Fodor, J. D. 1998a. Unambiguous Triggers. *Linguistic Inquiry* 29: 1-36.
- Fodor, J. D. 1998b. Parsing to learn. *Journal of Psycholinguistic Research* 27(3): 339-374.
- Fodor, J. D. 2009. Syntax Acquisition: An Evaluation Measure After All? In Piatelli Palmarini, M., Uriagereka, J., and Salaburu, P. (eds.), *Of Minds and Language: The Basque Country Encounter with Noam Chomsky*, pp.256-277. New York: Oxford University Press.
- Foraker, S., Regier, T., Khetarpal, A., Perfors, A., & Tenenbaum, J. 2009. Indirect evidence and the poverty of the stimulus: The case of anaphoric *one*. *Cognitive Science* 33: 287-300.
- Frank, M.C., Goodman, S., & Tenenbaum, J. 2009. Using Speakers' Referential Intentions to Model Early Cross-Situational Word Learning. *Psychological Science* 20(5): 578- 585.
- Gerken, L. 2006. Decision, decisions: infant language learning when multiple generalizations are possible. *Cognition* 98: B67-B74.
- Gibson, E. & Wexler, K. 1994. Triggers, *Linguistic Inquiry* 25: 355-407.

- Goldwater, S., T. Griffiths, & M. Johnson. 2009. A Bayesian Framework for Word Segmentation: Exploring the Effects of Context. *Cognition* 112(1): 21-54.
- Graf Estes, K., Evans, J., Alibali, M., & Saffran, J. 2007. Can Infants Map Meaning to Newly Segmented Words? *Psychological Science* 18(3): 254-260.
- Griffiths, T. & Tenenbaum, J. 2005. Structure and strength in causal induction. *Cognitive Psychology* 51: 334-384.
- Hale, J. (2001). A probabilistic Earley parser as a psycholinguistic model. In *Proceedings of the Second Meeting of the North American Chapter of the Association for Computational Linguistics*, pp. 159–166.
- Jaeger, T.F & Snider, N. 2008. Implicit learning and syntactic persistence: Surprisal and Cumulativity. In *Proceedings of the 30th Annual Meeting of the Cognitive Science Society*, pp. 1061-1066.
- Klein, D. & Manning, C. 2002. A Generative Constituent-Context Model for Improved Grammar Induction. In *Proceedings of the 40th Annual Meeting for the Association for Computational Linguistics*, pp. 128-135. Association for Computational Linguistics: Stroudsburg, PA.
- Levy, R. 2008. Expectation-based syntactic comprehension. *Cognition* 106: 1126–1177.
- Levy, R. 2011a. Integrating surprisal and uncertain-input models in online sentence comprehension: formal techniques and empirical results. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics*.
- Lightfoot, D. 1991. *How to Set Parameters: arguments from language change*, Cambridge, MA: MIT Press.

- Lightfoot, D. 2010. Language acquisition and language change. *Wiley Interdisciplinary Reviews: Cognitive Science* 1: 677-684. doi: 10.1002/wcs.39.
- MacWhinney, B. 2000. *The CHILDES project: Tools for analyzing talk*. Mahwah, NJ: Lawrence Erlbaum Associates.
- Manning, C. & Schütze, H. 1999. *Foundations of Statistical Natural Language Processing*, Cambridge, MA: MIT Press.
- Marcus, M., Santorini, B., Marcinkiewicz, M., & Taylor, A. 1999. *Treebank-3*. Philadelphia: Linguistic Data Consortium.
- Mintz, T. 2003. Frequent frames as a cue for grammatical categories in child directed speech. *Cognition* 90: 91-117.
- Mintz, T. 2006. Finding the verbs: distributional cues to categories available to young learners. In Hirsh-Pasek, K. & Golinkoff, R.M. (eds.), *Action Meets Word: How Children Learn Verbs*, pp. 31-63. New York: Oxford University Press.
- Niyogi, P. & Berwick, R. 1996. A language learning model for finite parameter spaces. *Cognition* 61: 161-193.
- Pearl, L., Goldwater, S., & Steyvers, M. 2011. Online Learning Mechanisms for Bayesian Models of Word Segmentation, *Research on Language and Computation*, special issue on computational models of language acquisition. doi 10.1007/s11168-011-9074-5.
- Pearl, L., & Lidz, J. 2009. When domain-general learning fails and when it succeeds: Identifying the contribution of domain-specificity. *Language Learning and Development* 5(4): 235-265.

- Pearl, L. & Mis, B. 2011. How Far Can Indirect Evidence Take Us? Anaphoric One Revisited. In *Proceedings of the 33rd Annual Conference of the Cognitive Science Society*. Boston, MA: Cognitive Science Society.
- Pearl, L. & Mis, B. (submitted). What Indirect Evidence Can Tell Us About Universal Grammar: Anaphoric One Revisited. Ms., University of California, Irvine.
- Perfors, A., Tenenbaum, J., & Regier, T. 2011. The learnability of abstract syntactic principles. *Cognition* 118: 306–338.
- Pelucchi, B., Hay, J., & Saffran, J. 2009a. Statistical Learning in Natural Language by 8-Month-Old Infants. *Child Development* 80(3): 674-685.
- Pelucchi, B., Hay, J., & Saffran, J. 2009b. Learning in reverse: Eight-month-old infants track backward transitional probabilities. *Cognition* 113(2): 244-247.
- Phillips, C. 2006. The real-time status of island constraints. *Language* 82: 795-823.
- Saffran, J., Aslin, R., & Newport, E. 1996. Statistical Learning by 8-Month-Old Infants. *Science* 274: 1926-1928.
- Saffran, J. R., Hauser, M., Seibel, R. L., Kapfhamer, J., Tsao, F., & Cushman, F. 2008. Grammatical pattern learning by infants and cotton-top tamarin monkeys. *Cognition* 107: 479-500.
- Saffran, J. R., Johnson, E. K., Aslin, R. N., & Newport, E. L. 1999. Statistical learning of tone sequences by human infants and adults. *Cognition* 70(1): 27-52.
- Sagae, K., Davis, E., Lavie, A., MacWhinney, B., & Wintner, S. 2010. Morphosyntactic annotation of CHILDES transcript. *Journal of Child Language*, 37(3), 705-729.

- Sakas, W.G. & Fodor, J.D. 2001. The structural triggers learner. In Bertolo, S. (ed.) *Language Acquisition and Learnability*, pp. 172-233. Cambridge, UK: Cambridge University Press.
- Sprouse, J. *this volume*. Defining the terms of the debate: Reductionist theories and the superadditive nature of island effects. In Sprouse, J. and Hornstein, N. (eds.), *Experimental Syntax and Island Effects*. Cambridge University Press.
- Sprouse, J., M. Wagers, and C. Phillips. 2012. A test of the relation between working memory capacity and syntactic island effects. *Language* 88.
- Stowe, L. 1986. Parsing wh-constructions: Evidence for on-line gap location. *Language and Cognitive Processes* 1: 227-245.
- Suppes, P. 1974. The semantics of children's language. *American Psychologist* 29: 103- 114.
- Tenenbaum, J. & Griffiths, T. 2001. Generalization, similarity and Bayesian inference. *Behavioral and Brain Sciences* 24: 629-640.
- Tribus, M. 1961. *Thermodynamics and Thermostatistics: An Introduction to Energy, Information and States of Matter, with Engineering Applications*. New York, NY: D. Van Nostrand Company Inc.
- Valian, V. 1991. Syntactic subjects in the early speech of American and Italian children. *Cognition* 40: 21-81.
- Wagers, M., & Phillips, C. 2009. Multiple dependencies and the role of the grammar in real-time comprehension. *Journal of Linguistics*, 45, 395-433.
- Wang, H. & Mintz, T. 2008. A Dynamic Learning Model for Categorizing Words Using Frames. In Chan, H., Jacob, H., & Kapia, E. (eds.), *BUCLD 32 Proceedings*, pp.525-536. Somerville, MA: Cascadilla Press.

Xu, F., & Tenenbaum, J. 2007. Word learning as Bayesian inference. *Psychological Review* 114: 245-272.

Yang, C. 2002. *Knowledge and Learning in Natural Language*. Oxford: Oxford University Press.