

## Lecture 4

# Top-Down Parsing: Psycholinguistic Adequacy

### 4.1 Two Basic Advantages

#### 4.1.1 Incrementality

The most basic requirement for a psychologically plausible parser is that it works in an incremental fashion, that is to say, parsing can take place as soon as the first word of the sentence has been uttered, rather than delaying parsing until the full sentence has been heard. You showed in an earlier exercise that the top-down parser only needs to keep track of the starting position of an item, which entails that it can be run in an incremental fashion.

#### 4.1.2 Simplicity

Top-down parsers are conceptually simple. The top-down orientation of the parser mirrors the view of CFGs as top-down generators. In addition, the inference rules are extremely simple (compared to, say, *left-corner parsing*, which we will discuss at a later point), seeing how the inference rule of the parser directly mirrors the rewrite rules of the grammar. This is particularly appealing for proponents of a *transparent parser*. A transparent parser must use the grammar in as direct a fashion as possible. That precludes processing-based explanations of syntactic phenomena — e.g. island effects — and prioritizes a maximally simple parser.

### 4.2 Garden Path Effects

#### 4.2.1 Garden Paths as Backtracking

One of the best known phenomena in syntactic processing are *garden path effects*, which arise with sentences that are grammatically well-formed but nonetheless difficult to parse (Frazier 1979; Frazier and Rayner 1982). More precisely, a garden path sentence is a sentence  $w$  that up to some  $w_i$  has a strongly preferred analysis that must be discarded at  $w_{i+1}$ . Here's a list of examples that includes more than the usual suspects:

- (5) *Structural ambiguity*

- a. The horse raced past the barn fell.
  - b. The raft floated down the river sank.
  - c. The player tossed a Frisbee smiled.
  - d. The doctor sent for the patient arrived.
  - e. The cotton clothing is made of grows in Mississippi.
  - f. Fat people eat accumulates.
  - g. I convinced her children are noisy.
- (6) *Lexical ambiguity*
- a. The old train the young.
  - b. The old man the boat.
  - c. Until the police arrest the drug dealers control the street.
  - d. The dog that I had really loved bones.
  - e. The man who hunts ducks out on weekends.

[Frazier \(1979\)](#) proposes to treat garden path effects as a result of reanalysis: the parser is forced to abandon its current set of hypotheses, backtrack to an earlier point, and build a new structure from there. If for some reason this process proves too difficult, the parser gets stuck and assigns no structure at all.

This kind of analysis is compatible with top-down parsing, depending on the data structure we use. Remember that the data structure keeps track of the items in the current parse, but also of alternative parses. The control structure includes a set of instructions for which parses should be prioritized. Let's try to make [Frazier's](#) account precise using a top-down parser with priority queues as data structure and a serial strategy for expanding multiple parses.

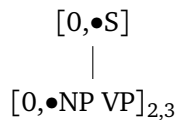
#### 4.2.2 Priority Queues and Prefix Trees

Let's assume that our data structure is a *priority queue*, i.e. a list of parse tables that can be reordered, to which we can add new entries, and from which we can delete old entries (note: even though I call it a list, that doesn't mean queues are best implemented as lists; in Python, for example, you're better off using arrays).

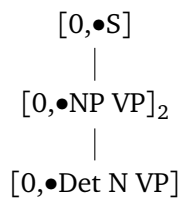
As so often, a more abstract perspective is helpful here. Instead of a list of parse tables, we can use a *prefix tree* as a unified way of representing multiple parses at once. The root of the prefix tree is always our start axiom  $[0, S]$  (since the final position is redundant for top-down parsing, it can be safely omitted). If parse item  $q$  is obtained from item  $p$  via a prediction or scan rule,  $q$  is added to the tree as a daughter of  $p$ . Furthermore, each node that is not a leaf is also subscripted with the list of unused predict rules. We call a node in this tree a *parse leaf* iff it is a leaf or has at least one subscript, and a (strictly descending) path spanning from the root to a leaf or a parse leaf is called a *parse path*. The set of parse tables, then, corresponds exactly to the set of parse paths. Putting all of this together, we can implement a priority queue as a ranking of parse paths through the prefix tree.

**Example 4.1**

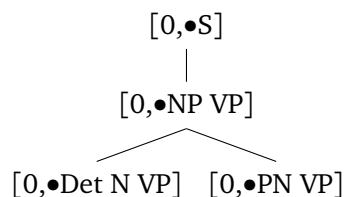
Parse Tables as Trees Suppose we are using the grammar from the beginning of Lecture 3 to parse *the anvil hit Daffy* with a recursive descent parser. Our parse table starts with  $[0, \bullet S]$  as usual, from which we can only predict  $[0, \bullet NP VP]$ . This can be represented as the tree below.



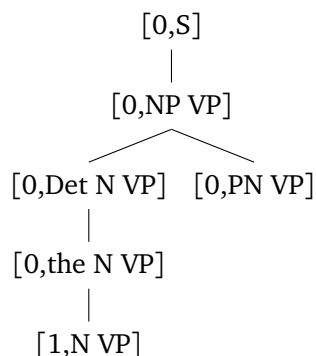
Note that  $[0, \bullet NP VP]$  has subscripts 2 and 3 to indicate that we can use rules 2 and 3 of our CFG to predict new parse items. Suppose the parser uses  $\text{predict}(3)$  to create the item  $[0, \bullet \text{Det } N VP]$ . Then this item is added as a daughter of  $[0, \bullet NP VP]$  to the previous tree, and the subscript 3 is also removed from  $[0, \bullet NP VP]$ . We do not need to add a subscript to  $[0, \bullet \text{Det } N VP]$  since it is a leaf.



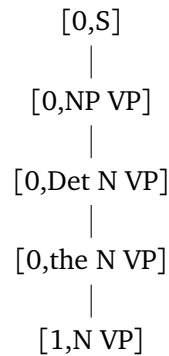
The next step of the parser now depends on how the control structure prioritizes distinct parses. We can either expand the table that ends in  $[0, \bullet \text{Det } N VP]$ , or go back to the one ending in  $[0, \bullet NP VP]$  and apply  $\text{predict}(2)$ . Suppose we do the latter, so that  $[0, \bullet PN VP]$  is added as the second daughter of  $[0, \bullet NP VP]$ , which thus loses its last subscript.



Now let's expand  $[0, \bullet \text{Det } N VP]$  again to obtain  $[0, \bullet \text{the } N VP]$  and then apply a scan step, yielding  $[1, \bullet N VP]$ .



If our parser is really smart, it will be able to infer at this point that the scanned word *the* can never be obtained from  $[0, \text{PN VP}]$  (technically this is achieved by associating every parse item  $p$  with a regular expression that describes the possible left edges of the strings that can be derived from  $p$ ). So the parser can remove  $[0, \text{PN VP}]$  from the tree, which is tantamount to discarding the parse table where NP was rewritten as PN.



What makes the tree representation of parse tables appealing is that the construction and prioritization of parse tables can be reduced to strategies for tree building. In particular, the familiar notions of depth-first and breadth-first carry over in a natural fashion.

**depth first/serial** expand a parse item  $p$  that was introduced during the previous parse step; if the parse item  $p$  cannot be expanded, expand the lowest parse leaf  $l$  that dominates  $p$

**breadth first/parallel** before expanding a parse item introduced during parse step  $j$ , all parse items that were introduced at parsing step  $i$  must have been rewritten, for every  $i < j$

The depth first strategy leads to a parser that always builds a single complete parse history rather than multiple partial ones. If the parse history cannot be expanded anymore, the parser either stops (successful parse) or backtracks to the last choice point in the parse history and tries a different choice instead. This corresponds exactly to the notion of *serial parsing* in the psycholinguistic literature.

The intuitive counterpart of serial parsing is *fully parallel parsing*, where all parse tables are built up at the same time. This is exactly the behavior of the breadth-first strategy. However, even proponents of parallel parsing usually do not assume that the human parser computes and stores all options all the time. Instead, only a subset of very likely parses is claimed to be worked on in parallel, with all others either discarded or at least not expanded on. On a technical level we may formalize this in terms of a probabilistic procedure for which parse items may be expanded, and in which order. The details are not important here, the basic insights is simply that just as we can add probabilities to the control schema to regulate which inference rules the parser may apply, we can also use probabilities to prioritize the expansion of certain parse tables over others.

### 4.2.3 Formal Account of Garden Paths

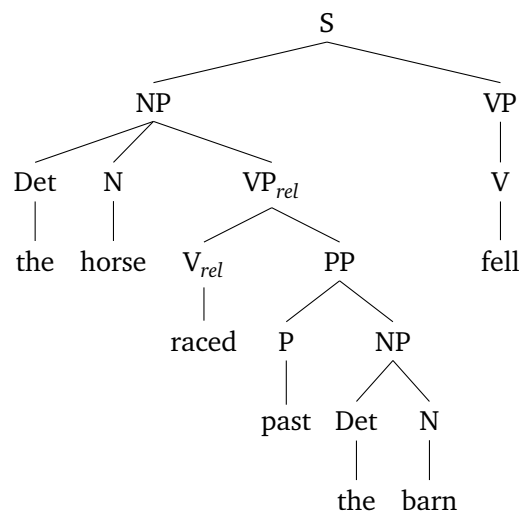
A recursive descent parser with a depth-first expansion strategy for prefix trees, coupled with a preference ranking over prediction steps, can easily account for garden path effects. In the case of *the horse raced past the barn fell*, for example, the parser builds a single parse table, and due to how certain rewrite rules are preferred over others, this parse table encodes the structure for *the horse raced past the barn*. When encountering *fell*, the parser has to backtrack. A successful parse requires backtracking to the point where *raced* is analyzed as part of the VP rather than the NP — that's quite a distance.

#### Example 4.2 Backtracking in *the horse raced past the barn fell*

Assume we have the following (massively simplified) grammar:

- |  |                              |
|--|------------------------------|
| 1) S → NP VP                               | 8) Det → the                 |
| 2) NP → Det N                              | 9) N → barn                  |
| 3) NP → Det N VP <sub>rel</sub>            | 10) N → horse                |
| 4) VP → V                                  | 11) P → past                 |
| 5) VP → V PP                               | 12) V → fell                 |
| 6) VP <sub>rel</sub> → V <sub>rel</sub> PP | 13) V → raced                |
| 7) PP → P NP                               | 14) V <sub>rel</sub> → raced |

Here's the resulting tree for our garden path sentence.



If the parser prefers NP → Det N over NP → Det N CP and operates in a recursive descent fashion, the construction of the first parse table results in the prefix tree below.

```

[0,•S]
|
[0,•NP VP]3
|
[0,•Det N VP]
|
[0,•the N VP]
|
[1,•N VP]9
|
[1,•horse VP]
|
[2,•VP]4
|
[2,•V PP]12
|
[2,•raced PP]
|
[3,•PP]
|
[3,•P NP]
|
[3,•past NP]
|
[4,•NP]3
|
[4,•Det N]
|
[4,•the N]
|
[5,•N]10
|
[5,•barn]
|
[6,•]

```

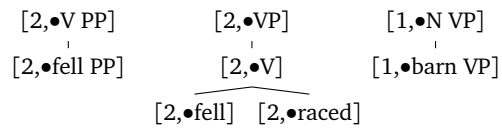
Since this parse does not succeed, the parser needs to backtrack. The closest choice point is  $[5, \bullet N]_{10}$ , which obviously does not fix the problem of integrating *fell* into the structure, as can be verified after a single scan step. The next choice point is  $[4, \bullet NP]_3$ . Here the parser still has the option of replacing NP by Det N CP, which won't help much either, but it takes quite a while to realize this because the tree for the parse table obtained by following this option involves a choice point, too.

```

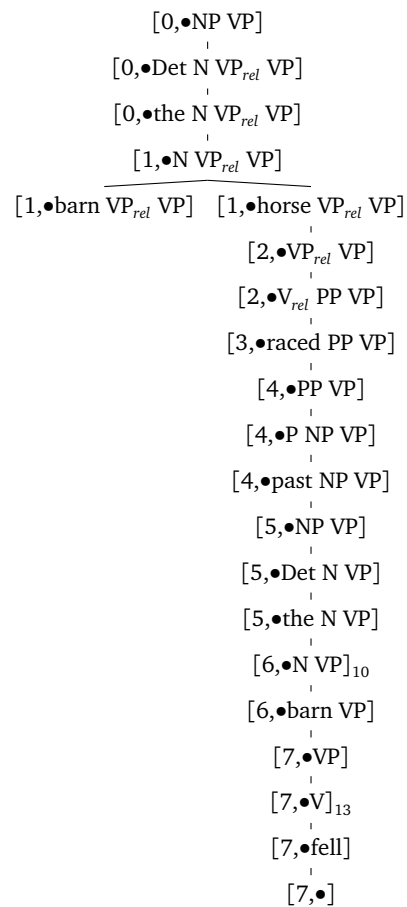
[4,•NP]
|
[4,•Det N VPrel]
|
[4,•the N VPrel]
|
[5,•N VPrel]
|
[5,•barn VPrel] [5,•horse VPrel]
|
[6,•VPrel]
|
[6,•Vrel PP]
|
[6,•fell PP]

```

Since this venue didn't yield a successful parse either, the parse backtracks to  $[2, \bullet V PP]_{12}$ , and after this fails, to  $[2, \bullet VP]_4$ . Once again it is not successful, and the same holds once it expands  $[1, \bullet N VP]$ .



Only if the parser backtracks all the way to  $[0, \bullet \text{NP VP}]_3$ , essentially undoing all its work so far, can it find a working parse.



The prefix tree for all the parse histories built by the parser before it encounters a successful parse is much bigger than the simple phrase structure tree for the sentence.





needed to close the containing CP, which can subsequently be removed from memory. Center embedding, on the other hand, requires that both CPs be stored in memory, so it is natural limited to one level of embedding.

Gibson's (1998) *Syntactic Prediction Locality Theory* (SPLT) contends that memory burden increases the more dependencies need to be stored at the same time, and since center embedding of the type NP-S-VP necessarily involves starting a new dependency before the old one between NP and VP has been finished, center embedding is difficult.

This argument can be ported into our more formal setting via a *linking hypothesis* between control structure and processing difficulty. Remember that we can annotate phrase structure trees to indicate the order in which nodes are introduced and rewritten by the parser. Obviously any item that is not rewritten immediately after its introduction — i.e. any item whose prefix and suffix differ by more than 1 — needs to be stored in memory. Let's put these items in boxes so that they are easy to pick out at a glance. Parsing difficulty, then, can be measured in a variety of ways.

**Tenure** the *tenure of node  $n$*  is the difference between its subscript and its superscript

**Payload** the *payload of tree  $t$*  is the number of nodes whose tenure is strictly greater than 1

These two measures can be combined to create a variety of difficulty metrics such as the two below (cf. Koble et al. 2012 and Graf and Marcinek 2014).

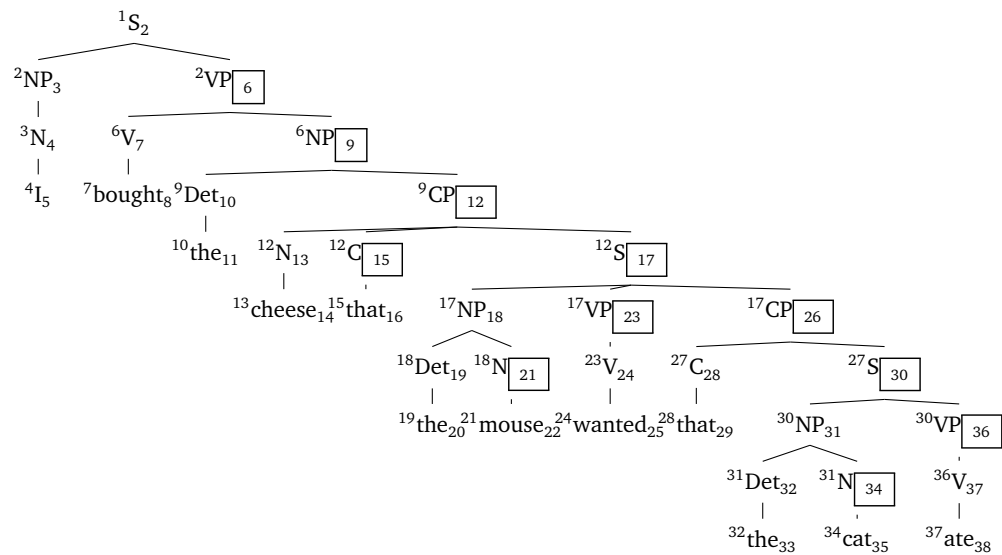
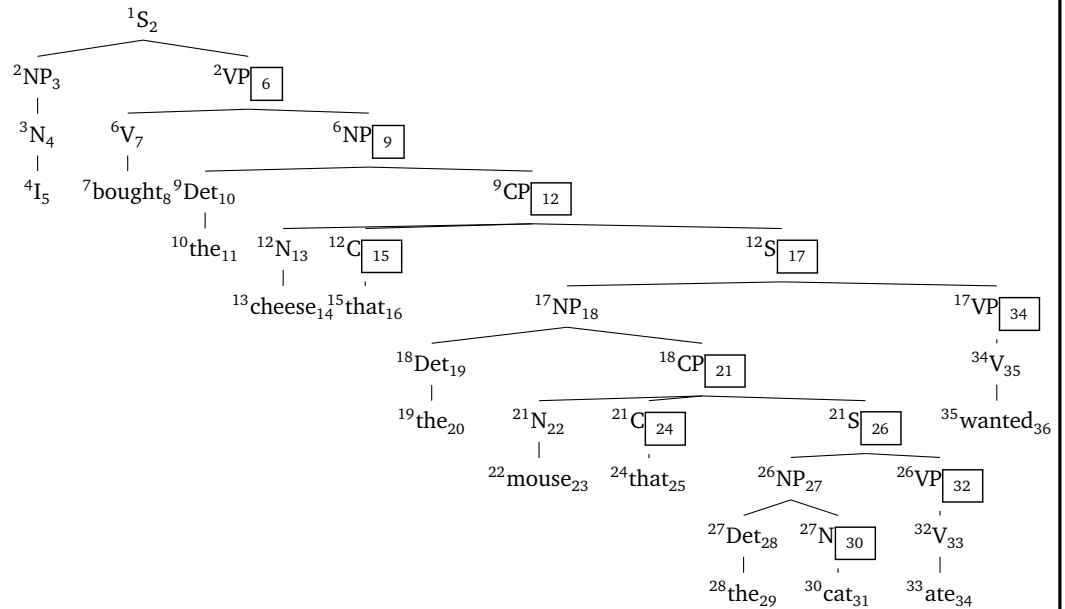
**MaxTen** greatest tenure among all the nodes;  $\max(\{\text{tenure}(n)\})$

**SumTen** total number of steps that items must be memorized;  $\sum(\{\text{tenure}(n) > 1\})$

All of these metrics capture the fact that right embedding isn't harder than center embedding. **MaxTen** and **SumTen** also predict right embedding to be easier.

#### Example 4.3 Center embedding and maximum tenure

Consider the center embedding and right embedding trees below, which have been annotated according to the behavior of a recursive descent parser. For center embedding we are using a promotion-style analysis of relative clauses (Vergnaud 1974; Kayne 1994), which posits that the relative clause is an argument of the determiner, with the head noun residing in a CP specifier. For right embedding the extraposed relative clause is considered a daughter of S.



The difficult center embedding sentence has a payload of 11, a maximum tenure of 17, and a sum tenure of 55. The easier right embedding sentence has a payload of 11, a maximum tenure of 9, and a sum tenure of 48. Irrespective of how we weigh or rank these three metrics, then, right embedding is never predicted to be harder than center embedding. As long as we do not take payload as the only difficulty metric, right embedding is correctly predicted to be easier than center embedding.

Notice that our explanation for the difficulty of center embedding is very different from the one we used for garden path effects. Garden path effects were explained in terms of the difficulty of finding a right parse among the many possible ones. We used prefix trees as a way of representing the parser's route through this search space, and since a prefix tree encodes many parsing tables at ones, our account is ultimately

based on the processing challenges of coordinating multiple parse tables — if all parse tables could be easily stored in memory and expanded in a breadth-first manner, garden path effects would not arise.

Center embedding, on the other hand, is explained purely in terms of how the parser constructs the correct parse. How the parser actually finds this parse does not factor into the explanation. The claim is that even if the parser had a perfect *oracle*, a machine that could tell it at every step which inference step must be taken to get the correct structure, the difference between center embedding and right embedding would not disappear because the former still puts a higher load on working memory than the latter.

## 4.3 Problems

### 4.3.1 Memory Usage in Left Embedding

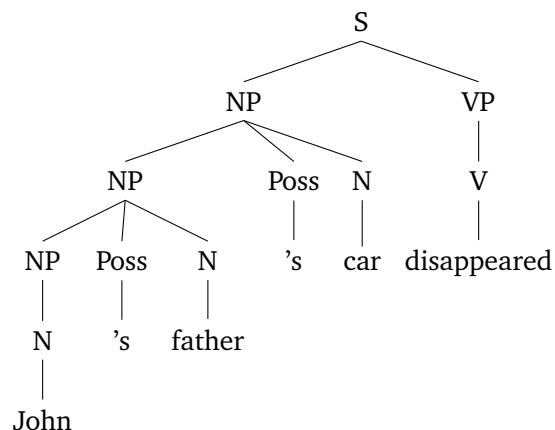
A careful examination of the top-down parsers behavior over center embedding constructions reveals that these configurations aren't the only ones that should cause an increased memory burden. Coupled with our metrics of processing difficulty, a recursive descent parser predicts that left embedding should be hard to parse, too — which is not borne out empirically.

- (8) a. The exhaust pipe of the car of the father of the mechanic is broken.  
 b. The mechanic's father's car's exhaust pipe is broken.

The cause for the predicted difficulty spike with left embedding is straight-forward. Suppose the parser conjectures at step  $i$  that node  $n$  in the tree has daughters  $d_1$  and  $d_2$ . Steps  $i + 1$  to  $j$  are spent by the parser expanding the subtree rooted in  $d_1$ . The bigger the subtree, the higher the value of  $j + 1$ , the point at which the parser can move on to expanding  $d_2$ . Since  $d_2$  must be kept in memory from step  $i$  until step  $j + 1$ , and the value of  $j$  increases with the size of the subtree rooted in  $d_1$ , every construction that increases the size of a left sibling — including left embedding — should increase parsing difficulty.

### 4.3.2 Looping in Left Recursion

Left embedding constructions are problematic for top-down parsers in more respects than just psycholinguistic adequacy. Consider a grammar that licenses possessive structures such as the one below.



In order to arrive at this structure, the parser must first infer  $[0, \bullet \text{NP VP}]$ , and then expand the NP to get  $[0, \bullet \text{NP Poss N VP}]$ . But now the parser once again has to expand the NP, possibly creating  $[0, \bullet \text{NP Poss N Poss N VP}]$ . It is easy to see that the parser can keep rewriting NP *ad infinitum*, essentially looping the NP rewriting step and producing longer and longer parse items. Without a smart control structure that detects this kind of looping, the parser will soon run out of working memory, or if there is no limited on the amount of memory, keep looping forever without ever constructing a parse.

A standard solution is to use a probabilistic control structure that prunes away all parse tables that fall below a certain probability threshold. As the parser keeps conjecturing bigger and bigger structures, the probability of these parses decreases until they are eventually pruned away. This also called a *beam parser*. Beam parsers require a method for automatically inferring the necessary probabilities from a corpus, and their overall behavior can be difficult to figure out. So the appealing simplicity of top-down parsers is lost to some extent.

### 4.3.3 Merely Local Syntactic Coherence Effects

Merely local syntactic coherence effects is a term coined by Tabor et al. (2004) to refer to cases where an analysis that is locally well-formed but incompatible with the structure built so far nonetheless induces a temporary increase in parsing difficulty (see also Konieczny 2005, Konieczny et al. 2009, and Bicknell et al. 2009). This is exemplified by the contrast in (9). During self-paced reading experiments, subjects' reading speed of (9a) decreases at *tossed*, indicating an increase in processing difficulty. The effect disappears, however, when *tossed* is replaced by *thrown*.

- (9) a. The coach smiled at the player tossed a frisbee.
- b. The coach smiled at the player thrown a frisbee.

It seems as if the fact that *the player* can be locally analyzed as the subject of *tossed* confuses the parser. Since this interpretation is not available with *thrown* due to the unambiguous past participle morphology, no slowdown occurs.

Merely local syntactic coherence effects are completely unexpected with a top-down parser because there is no way the parser could infer an item that treats *tossed* as a finite verb with *the player* as its subject. The preceding parts of the sentence have already disambiguated the structure of the sentence to an extent where this is no longer a viable hypothesis. These effects are expected, however, if the parser proceeds bottom-up instead of top-down.

## References and Further Reading

- Bicknell, Klinton, Roger Levy, and Vera Demberg. 2009. Correcting the incorrect: Local coherence effects modeled with prior belief update. In *Proceedings of the 35th Annual Meeting of the Berkeley Linguistics Society*.
- Frazier, Lyn. 1979. *On comprehending sentences: Syntactic parsing strategies*. Doctoral Dissertation, University of Connecticut.
- Frazier, Lyn, and Keith Rayner. 1982. Making and correcting errors during sentence comprehension: Eye movements in the analysis of structurally ambiguous sentences. *Cognitive Psychology* 14:178–210.

- Gibson, Edward. 1998. Linguistic complexity: Locality of syntactic dependencies. *Cognition* 68:1–76.
- Graf, Thomas, and Bradley Marcinek. 2014. Evaluating evaluation metrics for minimalist parsing. In *Proceedings of the 2014 ACL Workshop on Cognitive Modeling and Computational Linguistics*, 28–36.
- Kayne, Richard S. 1994. *The antisymmetry of syntax*. Cambridge, Mass.: MIT Press.
- Kimball, John. 1973. Seven principles of surface structure parsing in natural language. *Cognition* 2:15–47.
- Kobele, Gregory M., Sabrina Gerth, and John T. Hale. 2012. Memory resource allocation in top-down minimalist parsing. In *Proceedings of Formal Grammar 2012*.
- Konieczny, Lars. 2005. The psychological reality of local coherences in sentence processing. In *Proceedings of the 27th Annual Conference of the Cognitive Science Society*.
- Konieczny, Lars, Daniel Müller, Wibke Hachmann, Sarah Schwarzkopf, and Sascha Wolfer. 2009. Local syntactic coherence interpretation. Evidence from a visual world study. In *Proceedings of the 31st Annual Conference of the Cognitive Science Society*, 1133–1138.
- Tabor, Whitney, Bruno Galantucci, and Daniel Richardson. 2004. Effects of merely local syntactic coherence on sentence processing. *Journal of Memory and Language* 50:355–370.
- Vergnaud, Jean-Roger. 1974. *French relative clauses*. Doctoral Dissertation, MIT.