

Incremental Parsing in Bounded Memory

William Schuler

Department of Linguistics
The Ohio State University
schuler@ling.osu.edu

Abstract

This tutorial will describe the use of a factored probabilistic sequence model for parsing speech and text using a bounded store of three to four incomplete constituents over time, in line with recent estimates of human short-term working memory capacity. This formulation uses a grammar transform to minimize memory usage during parsing. Incremental operations on incomplete constituents in this transformed representation then define an extended domain of locality similar to those defined in mildly context-sensitive grammar formalisms, which can similarly be used to process long-distance and crossed-and-nested dependencies.

1 Introduction

This paper will describe the derivation of a factored probabilistic sequence model from a probabilistic context-free grammar (PCFG). The resulting sequence model can incrementally parse input sentences approximately as accurately as a bottom-up CKY-style parser, incrementally estimating the contents of a bounded memory store of intended constituents, consisting of only three to four working memory elements, in line with recent estimates of human short-term working memory capacity (Cowan, 2001). The detailed derivation of this model is intended to illustrate how probabilistic dependencies from an original PCFG (or other types of syntax-derived dependencies) can be preserved in a processing model with human-like working memory constraints.

1.1 Notation

This paper will associate variables for syntactic categories c , trees or subtrees τ , and string yields \bar{x} with constituents in phrase structure trees, identified using subscripts that describe the path from the root of the tree containing this constituent to the constituent itself. These paths may consist of left branches (indicated by ‘0’s in the path) and right branches (indicated by ‘1’s), concatenated into sequences η (or ι or κ). Thus, if a path η identifies a constituent, that constituent’s left child would be identified by $\eta 0$, and that constituent’s right child would be identified by $\eta 1$. The empty path ϵ will be used to identify the root of a tree.

The probabilistic parsers defined here will also use an indicator function $\llbracket \cdot \rrbracket$ to denote deterministic probabilities: $\llbracket \phi \rrbracket = 1$ if ϕ is true, 0 otherwise.

2 Parsing

For a phrase structure subtree rooted at a constituent of category c_η with yield \bar{x}_η (a sequence of individual words x), the task of parsing will require the calculation of the inside likelihood probability or Viterbi (best subtree) probabilities. When the domain X of words x is a subset of the domain C of category labels c , these can be calculated using rule probabilities in a probabilistic context-free grammar (PCFG) model θ_G , notated:

$$P_{\theta_G}(c_\eta \rightarrow c_{\eta 0} c_{\eta 1}) = P_{\theta_G}(c_{\eta 0} c_{\eta 1} | c_\eta) \quad (1)$$

Any yield \bar{x}_η can be decomposed into prefix $\bar{x}_{\eta 0}$ and suffix $\bar{x}_{\eta 1}$ yields:

$$\bar{x}_\eta = \bar{x}_{\eta 0} \bar{x}_{\eta 1} \quad (2)$$

Therefore, inside likelihood probabilities can be defined by marginalizing over all such decompositions:

$$P_{\theta_{\text{Ins(G)}}}(\bar{x}_\eta | c_\eta) = \begin{cases} \text{if } |\bar{x}_\eta| = 1 : \llbracket \bar{x}_\eta = c_\eta \rrbracket \\ \text{otherwise : } \sum_{\bar{x}_{\eta 0} c_{\eta 0}, \bar{x}_{\eta 1} c_{\eta 1}} P_{\theta_G}(c_\eta \rightarrow c_{\eta 0} c_{\eta 1}) \cdot P_{\theta_{\text{Ins(G)}}}(\bar{x}_{\eta 0} | c_{\eta 0}) \cdot P_{\theta_{\text{Ins(G)}}}(\bar{x}_{\eta 1} | c_{\eta 1}) \end{cases} \quad (3)$$

and Viterbi scores (the probability of the best tree) can be defined by maximizing over all such decompositions:

$$P_{\theta_{\text{Vit(G)}}}(\bar{x}_\eta | c_\eta) = \begin{cases} \text{if } |\bar{x}_\eta| = 1 : \llbracket \bar{x}_\eta = c_\eta \rrbracket \\ \text{otherwise : } \max_{\bar{x}_{\eta 0} c_{\eta 0}, \bar{x}_{\eta 1} c_{\eta 1}} P_{\theta_G}(c_\eta \rightarrow c_{\eta 0} c_{\eta 1}) \cdot P_{\theta_{\text{Vit(G)}}}(\bar{x}_{\eta 0} | c_{\eta 0}) \cdot P_{\theta_{\text{Vit(G)}}}(\bar{x}_{\eta 1} | c_{\eta 1}) \end{cases} \quad (4)$$

From these, it is possible to obtain the probability of a sentence \bar{x}_ϵ :

$$P(\bar{x}_\epsilon) = \sum_{c_\epsilon} P_{\theta_{\text{Ins(G)}}}(\bar{x}_\epsilon | c_\epsilon) \cdot P(c_\epsilon) \quad (5)$$

and the most likely tree:

$$\hat{\tau}_\epsilon = \underset{\tau_\epsilon}{\text{argmax}} P_{\theta_{\text{Vit(G)}}}(\bar{x}_\epsilon | c_{\tau_\epsilon}) \cdot P(c_{\tau_\epsilon}) \quad (6)$$

3 Incremental Parsing using Incomplete Constituents

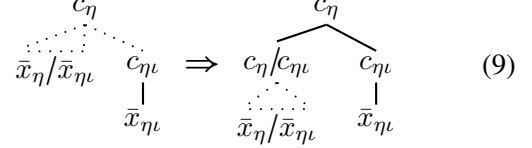
Note that any prefix of a yield \bar{x}_η can be rewritten as an ‘incomplete yield’ consisting of the complete yield lacking some suffix of that yield $\bar{x}_{\eta\iota}$:

$$\bar{x}_\eta = (\bar{x}_\eta / \bar{x}_{\eta\iota}) \bar{x}_{\eta\iota} \quad (7)$$

It is therefore possible to decompose any inside likelihood or Viterbi probability into two parts: first, an ‘incomplete constituent’ probability of generating this incomplete yield $\bar{x}_\eta / \bar{x}_{\eta\iota}$ along with an *awaited* constituent of category $c_{\eta\iota}$, given an *active* constituent of category c_η ; and second, an ordinary inside or Viterbi probability (or ‘complete constituent’ probability) of generating $\bar{x}_{\eta\iota}$ given $c_{\eta\iota}$:

$$P_{\theta_{\text{Ins(G)}}}(\bar{x}_\eta | c_\eta) = \sum_{\iota \in 1^+, \bar{x}_{\eta\iota}, c_{\eta\iota}} P_{\theta_{\text{IC(G)}}}(\bar{x}_\eta / \bar{x}_{\eta\iota}, c_{\eta\iota} | c_\eta) \cdot P_{\theta_{\text{Ins(G)}}}(\bar{x}_{\eta\iota} | c_{\eta\iota}) \quad (8)$$

This decomposition can be represented graphically as a transformation of the structure of a recurrence from the standard PCFG recurrence above, corresponding to PCFG dependencies, to one involving incomplete constituents (in particular, this will transform the end of a right-expanding sequence into the beginning of a left-expanding sequence of incomplete constituents):

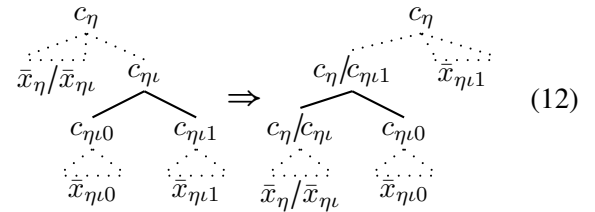


This decomposition gives incomplete constituent probabilities that can then be decomposed into other incomplete constituent probabilities (this will transform the middle of a right-expanding sequence into the middle of a left-expanding sequence of incomplete constituents):

$$P_{\theta_{\text{IC(G)}}}(\bar{x}_\eta / \bar{x}_{\eta\iota 1}, c_{\eta\iota 1} | c_\eta) = \sum_{\bar{x}_{\eta\iota}, c_{\eta\iota}} P_{\theta_{\text{IC(G)}}}(\bar{x}_\eta / \bar{x}_{\eta\iota}, c_{\eta\iota} | c_\eta) \cdot P_{\theta_{\text{IC(G)}}}(\bar{x}_{\eta\iota} / \bar{x}_{\eta\iota 1}, c_{\eta\iota 1} | c_{\eta\iota}) \quad (10)$$

$$= \sum_{\bar{x}_{\eta\iota}, c_{\eta\iota}} P_{\theta_{\text{IC(G)}}}(\bar{x}_\eta / \bar{x}_{\eta\iota}, c_{\eta\iota} | c_\eta) \cdot \sum_{\bar{x}_{\eta\iota 0}, c_{\eta\iota 0}} P_{\theta_G}(c_{\eta\iota} \rightarrow c_{\eta\iota 0} c_{\eta\iota 1}) \cdot P_{\theta_{\text{Ins(G)}}}(\bar{x}_{\eta\iota 0} | c_{\eta\iota 0}) \quad (11)$$

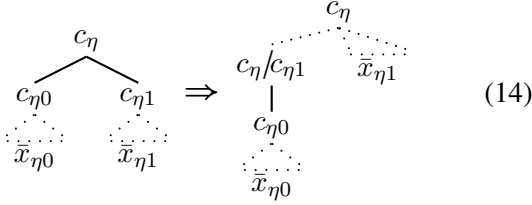
represented graphically:



or into the product of a grammar rule probability and an inside probability at the end of a sequence of such decompositions (this will transform the beginning of a right-expanding sequence into the end of a left-expanding sequence of incomplete constituents):

$$P_{\theta_{\text{IC(G)}}}(\bar{x}_\eta / \bar{x}_{\eta 1}, c_{\eta 1} | c_\eta) = \sum_{\bar{x}_{\eta 0}, c_{\eta 0}} P_{\theta_G}(c_\eta \rightarrow c_{\eta 0} c_{\eta 1}) \cdot P_{\theta_{\text{Ins(G)}}}(\bar{x}_{\eta 0} | c_{\eta 0}) \quad (13)$$

represented graphically:



Essentially this transformation turns right-expanding sequences of constituents (with subscript addresses ending in ‘1’) into left-expanding sequences of incomplete constituents, which can be composed together as they are recognized incrementally from left to right. The decompositions of Equations 8–13, taken together, are a model-based right-corner transform (Schuler, 2009).

This transformation can then be defined on depth-specific rules, allowing the sequence model to keep track of a bounded amount of center-embedded phrase structure.

$$P_{\theta_{\text{Ins}(G),d}}(\bar{x}_\eta | c_\eta) = \sum_{\nu \in 1^+, \bar{x}_{\eta\nu}, c_{\eta\nu}} P_{\theta_{\text{IC}(G),d}}(\bar{x}_\eta / \bar{x}_{\eta\nu}, c_{\eta\nu} | c_\eta) \cdot P_{\theta_{\text{Ins}(G),d}}(\bar{x}_{\eta\nu} | c_{\eta\nu}) \quad (15)$$

$$P_{\theta_{\text{IC}(G),d}}(\bar{x}_\eta / \bar{x}_{\eta\nu 1}, c_{\eta\nu 1} | c_\eta) = \sum_{\bar{x}_{\eta\nu}, c_{\eta\nu}} P_{\theta_{\text{IC}(G),d}}(\bar{x}_\eta / \bar{x}_{\eta\nu}, c_{\eta\nu} | c_\eta) \cdot P_{\theta_{\text{IC}(G),d}}(\bar{x}_{\eta\nu} / \bar{x}_{\eta\nu 1}, c_{\eta\nu 1} | c_{\eta\nu}) \quad (16)$$

$$= \sum_{\bar{x}_{\eta\nu}, c_{\eta\nu}} P_{\theta_{\text{IC}(G),d}}(\bar{x}_\eta / \bar{x}_{\eta\nu}, c_{\eta\nu} | c_\eta) \cdot \sum_{\bar{x}_{\eta\nu 0}, c_{\eta\nu 0}} P_{\theta_{G-R,d}}(c_{\eta\nu} \rightarrow c_{\eta\nu 0} c_{\eta\nu 1}) \cdot P_{\theta_{\text{Ins}(G),d+1}}(\bar{x}_{\eta\nu 0} | c_{\eta\nu 0}) \quad (17)$$

$$P_{\theta_{\text{IC}(G),d}}(\bar{x}_\eta / \bar{x}_{\eta 1}, c_{\eta 1} | c_\eta) = \sum_{\bar{x}_{\eta 0}, c_{\eta 0}} P_{\theta_{G-L,d}}(c_\eta \rightarrow c_{\eta 0} c_{\eta 1}) \cdot P_{\theta_{\text{Ins}(G),d}}(\bar{x}_{\eta 0} | c_{\eta 0}) \quad (18)$$

Estimating probabilities for the beginning and ending of these transformed left-expanding sequences will require the estimation of expected counts for repeated recursive decompositions of yet-unrecognized incomplete constituents according to Equation 8, marginalizing over values of \bar{x} . Since left children and right children are decomposed differently, the expected counts $\theta_{G-RL^*,d}$ will use PCFG probabilities $\theta_{G-R,d}$ and $\theta_{G-L,d}$ that are conditioned

on whether the expanding constituent is a right or left child, and on the center-embedding depth d of the expanding constituent. The expected counts $\theta_{G-RL^*,d}$ are of constituent categories $c_{\eta\nu}$ anywhere in the left progeny of a right child of category c_η :

$$E_{\theta_{G-RL^*,d}}(c_\eta \xrightarrow{0} c_{\eta 0} \dots) = \sum_{c_{\eta 1}} P_{\theta_{G-R,d}}(c_\eta \rightarrow c_{\eta 0} c_{\eta 1}) \quad (19)$$

$$E_{\theta_{G-RL^*,d}}(c_\eta \xrightarrow{k} c_{\eta 0^k 0} \dots) = \sum_{c_{\eta 0^k}} E_{\theta_{G-RL^*,d}}(c_\eta \xrightarrow{k-1} c_{\eta 0^k} \dots) \cdot \sum_{c_{\eta 0^{k-1}}} P_{\theta_{G-L,d}}(c_{\eta 0^k} \rightarrow c_{\eta 0^{k-1}} c_{\eta 0^k 1}) \quad (20)$$

$$E_{\theta_{G-RL^*,d}}(c_\eta \xrightarrow{*} c_{\eta\nu} \dots) = \sum_{k=0}^{\infty} E_{\theta_{G-RL^*,d}}(c_\eta \xrightarrow{k} c_{\eta\nu} \dots) \quad (21)$$

$$E_{\theta_{G-RL^*,d}}(c_\eta \xrightarrow{+} c_{\eta\nu} \dots) = E_{\theta_{G-RL^*,d}}(c_\eta \xrightarrow{*} c_{\eta\nu} \dots) - E_{\theta_{G-RL^*,d}}(c_\eta \xrightarrow{0} c_{\eta\nu} \dots) \quad (22)$$

In practice the infinite sum is estimated to some constant K using value iteration (Bellman, 1957).

These expected counts can then be used to calculate left progeny probabilities:

$$P_{\theta_{G-RL^*,d}}(c_\eta \xrightarrow{*} c_{\eta\nu} \dots) = \frac{E_{\theta_{G-RL^*,d}}(c_\eta \xrightarrow{*} c_{\eta\nu} \dots)}{\sum_{c_{\eta\nu}} E_{\theta_{G-RL^*,d}}(c_\eta \xrightarrow{*} c_{\eta\nu} \dots)} \quad (23)$$

which can be used to calculate forward or Viterbi probabilities for all incomplete constituents in the memory store, along with their yields:

$$P_{\theta_{\text{Fwd}}}((\bar{x}_{\eta_d} / \bar{x}_{\eta_d \nu_d 1})_{d=1}^D, (c_{\eta_d})_{d=1}^D, (c_{\eta_d \nu_d 1})_{d=1}^D) = \prod_{d=1}^D P_{\theta_{G-RL^*,d}}(c_{\eta_{d-1} \nu_{d-1} 1} \xrightarrow{*} c_{\eta_d} \dots) \cdot P_{\theta_{\text{IC}(G)}}(\bar{x}_{\eta_d} / \bar{x}_{\eta_d \nu_d 1}, c_{\eta_d \nu_d 1} | c_{\eta_d}) \quad (24)$$

Putting all the pieces together, probabilities for stores of incomplete constituents can now be defined in terms of transitions from possible previous stores of incomplete constituents (set apart by parentheses in the equation below) and reductions of incomplete constituents and terminal symbols into complete constituents (set apart by square brackets). These transitions either:

- add a layer of structure below an awaited constituent at some depth level (the first case below),
- add a layer of structure to the top of an active constituent at some depth level (the second case below), or
- carry forward the probability of an incomplete constituent at a depth at which no transition takes place (the third case below):

$$\begin{aligned}
& P_{\theta_{\text{Fwd}}}((\bar{x}_{\eta_d}/\bar{x}_{\eta_d \iota_d 1})_{d=1}^D, (c_{\eta_d})_{d=1}^D, (c_{\eta_d \iota_d 1})_{d=1}^D) = \\
& \prod_{d=1}^D \left\{ \begin{array}{l}
\text{if } c_{\eta_{d+1}} = \bar{x}_{\eta_{d+1}}, \iota_d \neq \epsilon : \\
\sum_{\bar{x}_{\eta_d \iota_d}, c_{\eta_d \iota_d}} \left(\begin{array}{l}
P_{\theta_{\text{G-RL}^*, d}}(c_{\eta_{d-1} \iota_{d-1} 1} \xrightarrow{*} c_{\eta_d} \dots) \\
\cdot P_{\theta_{\text{IC(G), d}}}(\bar{x}_{\eta_d}/\bar{x}_{\eta_d \iota_d}, c_{\eta_d \iota_d} | c_{\eta_d})
\end{array} \right) \\
\cdot \sum_{\bar{x}_{\eta_d \iota_d 0}, c_{\eta_d \iota_d 0}} \sum_{\kappa} s.t. \bar{x}_{\eta_d \iota_d 0 \kappa} = c_{\eta_d \iota_d 0 \kappa} \\
\left(\begin{array}{l}
P_{\theta_{\text{G-RL}^*, d+1}}(c_{\eta_d \iota_d} \xrightarrow{*} c_{\eta_d \iota_d 0} \dots) \\
\cdot P_{\theta_{\text{IC(G), d+1}}}(\bar{x}_{\eta_d \iota_d 0}/\bar{x}_{\eta_d \iota_d 0 \kappa}, c_{\eta_d \iota_d 0 \kappa} | c_{\eta_d \iota_d 0})
\end{array} \right) \\
\left[\begin{array}{l}
P_{\theta_{\text{G-RL}^*, d+1}}(c_{\eta_d \iota_d} \xrightarrow{0} c_{\eta_d \iota_d 0} \dots) \\
P_{\theta_{\text{G-RL}^*, d+1}}(c_{\eta_d \iota_d} \xrightarrow{*} c_{\eta_d \iota_d 0} \dots) \\
\cdot \llbracket \bar{x}_{\eta_d \iota_d 0 \kappa} = c_{\eta_d \iota_d 0 \kappa} \rrbracket
\end{array} \right] \\
P_{\theta_{\text{G-R}, d}}(c_{\eta_d \iota_d} \rightarrow c_{\eta_d \iota_d 0} c_{\eta_d \iota_d 1}) \\
E_{\theta_{\text{G-RL}^*, d+1}}(c_{\eta_d \iota_d} \xrightarrow{0} c_{\eta_d \iota_d 0} \dots) \\
\text{if } c_{\eta_{d+1}} = \bar{x}_{\eta_{d+1}}, \iota_d = \epsilon : \\
\sum_{\bar{x}_{\eta_d 0}, c_{\eta_d 0}} \sum_{\kappa} s.t. \bar{x}_{\eta_d 0 \kappa} = c_{\eta_d 0 \kappa} \\
\left(\begin{array}{l}
P_{\theta_{\text{G-RL}^*, d}}(c_{\eta_{d-1} \iota_{d-1} 1} \xrightarrow{*} c_{\eta_d 0} \dots) \\
\cdot P_{\theta_{\text{IC(G), d}}}(\bar{x}_{\eta_d 0}/\bar{x}_{\eta_d 0 \kappa}, c_{\eta_d 0 \kappa} | c_{\eta_d 0})
\end{array} \right) \\
\left[\begin{array}{l}
P_{\theta_{\text{G-RL}^*, d}}(c_{\eta_{d-1} \iota_{d-1} 1} \xrightarrow{+} c_{\eta_d 0} \dots) \\
P_{\theta_{\text{G-RL}^*, d}}(c_{\eta_{d-1} \iota_{d-1} 1} \xrightarrow{*} c_{\eta_d 0} \dots) \\
\cdot \llbracket \bar{x}_{\eta_d 0 \kappa} = c_{\eta_d 0 \kappa} \rrbracket
\end{array} \right] \\
P_{\theta_{\text{G-RL}^*, d}}(c_{\eta_{d-1} \iota_{d-1} 1} \xrightarrow{*} c_{\eta_d} \dots) \\
P_{\theta_{\text{G-RL}^*, d}}(c_{\eta_{d-1} \iota_{d-1} 1} \xrightarrow{+} c_{\eta_d 0} \dots) \\
\cdot P_{\theta_{\text{G-L}, d}}(c_{\eta_d} \rightarrow c_{\eta_d 0} c_{\eta_d 1}) \\
\text{if } c_{\eta_{d+1}} \neq \bar{x}_{\eta_{d+1}} : \\
\left(\begin{array}{l}
P_{\theta_{\text{G-RL}^*, d}}(c_{\eta_{d-1} \iota_{d-1} 1} \xrightarrow{*} c_{\eta_d} \dots) \\
\cdot P_{\theta_{\text{IC(G), d}}}(\bar{x}_{\eta_d}/\bar{x}_{\eta_d \iota_d}, c_{\eta_d \iota_d} | c_{\eta_d})
\end{array} \right)
\end{array} \right\} \quad (25)
\end{aligned}$$

Note that the left progeny probabilities above cancel out over time, leaving only the relevant original PCFG probabilities at the end of each sentence.

4 Parsing in a Factored Sequence Model

Store elements can now be abstracted away from (i.e. marginalized over) individual constituent structure addresses. Store elements are therefore defined to contain the only the active and awaited ($c_{s_t^d}^A$ and $c_{s_t^d}^W$) constituent categories necessary to compute an incomplete constituent probability:

$$s_t^d \stackrel{\text{def}}{=} \langle c_{s_t^d}^A, c_{s_t^d}^W \rangle \quad (26)$$

$$\stackrel{\text{def}}{=} \langle c_{\eta_d}, c_{\eta_d \iota_d} \rangle s.t. \iota_d \in 1^+, \eta_d \in \eta_{d-1} \iota_{d-1} 0^+ \quad (27)$$

Reduction states are defined to contain only the complete constituent category $c_{r_t^d}$ necessary to compute an inside likelihood probability, as well as a flag $f_{r_t^d}$ indicating whether a reduction has taken place (to end a sequence of incomplete constituents):

$$r_t^d \stackrel{\text{def}}{=} \langle c_{r_t^d}, f_{r_t^d} \rangle \quad (28)$$

$$\stackrel{\text{def}}{=} \langle c_{\eta_d}, \llbracket \bar{x}_{\eta_d} = c_{\eta_d} \rrbracket \rangle \quad (29)$$

Since $\iota_d \in 1^+$, it follows that:

$$x_{1..t} = (\bar{x}_{\eta_d}/\bar{x}_{\eta_d \iota_d})_{d=1}^D \quad (30)$$

This allows store elements to be abstracted away from (marginalized over) tree addresses in the calculation of forward or Viterbi probabilities:

$$\begin{aligned}
& P_{\theta_{\text{Fwd}}}(x_{1..t}, s_t^{1..D}) = \\
& \sum_{\eta_{1..D}, \iota_{1..D}} P_{\theta_{\text{Fwd}}}((\bar{x}_{\eta_d}/\bar{x}_{\eta_d \iota_d 1})_{d=1}^D, (c_{\eta_d})_{d=1}^D, (c_{\eta_d \iota_d 1})_{d=1}^D) \quad (31)
\end{aligned}$$

Forward probabilities can then be factored into contributions from previous store states (θ_{Fwd} at $s_{t-1}^{1..D}$ below, parenthesized in Equation 25), reductions of terminal symbols (θ_{B} below, bracketed in Equation 25), and transition operations (θ_{A} below, not set apart in Equation 25):

$$\begin{aligned}
& P_{\theta_{\text{Fwd}}}(x_{1..t}, s_t^{1..D}) = \\
& \sum_{s_{t-1}^{1..D}} P_{\theta_{\text{Fwd}}}(x_{1..t-1}, s_{t-1}^{1..D}) \cdot P_{\theta_{\text{A}}}(s_t^{1..D} | s_{t-1}^{1..D}) \\
& \cdot P_{\theta_{\text{B}}}(x_t | s_t^{1..D}) \quad (32)
\end{aligned}$$

Probabilities for recognized sequences of incomplete constituents (Equations 8 through 13), and expected left progeny counts for unrecognized sequences of incomplete constituents (Equations 19

through 22) can be combined in a probabilistic push-down automaton with a bounded pushdown store (to simulate the bounded working memory store of a human comprehender). This is essentially an extension of a Hierarchical Hidden Markov Model (HHMM) (Murphy and Paskin, 2001), which obtains a most likely sequence of hidden store states $\hat{s}_{1..T}^{1..D}$ of some length T and some maximum depth D , given a sequence of observations (e.g. words) $x_{1..T}$:

$$\hat{s}_{1..T}^{1..D} \stackrel{\text{def}}{=} \operatorname{argmax}_{s_{1..T}^{1..D}} \prod_{t=1}^T P_{\theta_A}(s_t^{1..D} | s_{t-1}^{1..D}) \cdot P_{\theta_B}(x_t | s_t^{1..D}) \quad (33)$$

The model generates each successive store only after considering whether each nested sequence of incomplete constituents has completed and reduced:

$$P_{\theta_A}(s_t^{1..D} | s_{t-1}^{1..D}) \stackrel{\text{def}}{=} \sum_{r_t^1..r_t^D} \prod_{d=1}^D P_{\theta_R}(r_t^d | r_t^{d+1} s_{t-1}^d s_{t-1}^{d-1}) \cdot P_{\theta_S}(s_t^d | r_t^{d+1} r_t^d s_{t-1}^d s_{t-1}^{d-1}) \quad (34)$$

The model probabilities for these store elements and reduction states can then be defined (from Murphy and Paskin 2001) to expand a new incomplete constituent after a reduction has taken place ($f_{r_t^d} = 1$), transition along a sequence of store elements if no reduction has taken place ($f_{r_t^d} = 0$):

$$P_{\theta_S}(s_t^d | r_t^{d+1} r_t^d s_{t-1}^d s_{t-1}^{d-1}) \stackrel{\text{def}}{=} \begin{cases} \text{if } f_{r_t^{d+1}} = 1, f_{r_t^d} = 1 : P_{\theta_{S-E,d}}(s_t^d | s_{t-1}^{d-1}) \\ \text{if } f_{r_t^{d+1}} = 1, f_{r_t^d} = 0 : P_{\theta_{S-T,d}}(s_t^d | r_t^{d+1} r_t^d s_{t-1}^d s_{t-1}^{d-1}) \\ \text{if } f_{r_t^{d+1}} = 0, f_{r_t^d} = 0 : \llbracket s_t^d = s_{t-1}^d \rrbracket \end{cases} \quad (35)$$

and possibly reduce a store element (terminate a sequence) if the store state below it has reduced ($f_{r_t^{d+1}} = 1$):

$$P_{\theta_R}(r_t^d | r_t^{d+1} s_{t-1}^d s_{t-1}^{d-1}) \stackrel{\text{def}}{=} \begin{cases} \text{if } f_{r_t^{d+1}} = 0 : \llbracket r_t^d = \mathbf{r}_\perp \rrbracket \\ \text{if } f_{r_t^{d+1}} = 1 : P_{\theta_{R-R,d}}(r_t^d | r_t^{d+1} s_{t-1}^d s_{t-1}^{d-1}) \end{cases} \quad (36)$$

where $s_t^0 = \mathbf{s}_\top$ and $r_t^{D+1} = \mathbf{r}_\top$ for constants \mathbf{s}_\top (an incomplete root constituent), \mathbf{r}_\top (a complete lexical constituent) and \mathbf{r}_\perp (a null state resulting from the

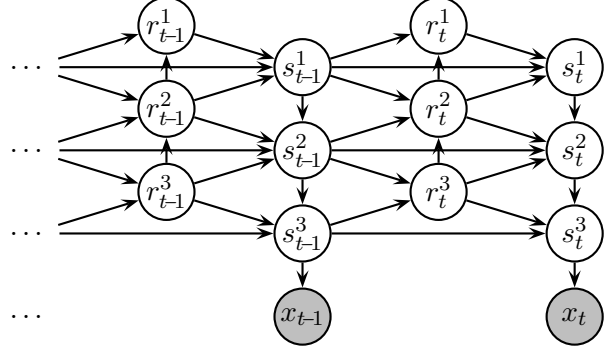


Figure 1: Graphical representation of the dependency structure in a standard Hierarchic Hidden Markov Model with $D = 3$ hidden levels that can be used to parse syntax. Circles denote random variables, and edges denote conditional dependencies. Shaded circles denote variables with observed values.

failure of an incomplete constituent to complete). The model is shown graphically in Figure 1.

These pushdown automaton operations are then refined for right-corner parsing (from Schuler 2009), distinguishing active transitions $\theta_{S-T-A,d}$ (in which an incomplete constituent is completed, but not reduced, and then immediately expanded to a new incomplete constituent in the same store element) from awaited transitions $\theta_{S-T-W,d}$ (which involve no completion):

$$P_{\theta_{S-T,d}}(s_t^d | r_t^{d+1} r_t^d s_{t-1}^d s_{t-1}^{d-1}) \stackrel{\text{def}}{=} \begin{cases} \text{if } r_t^d \neq \mathbf{r}_\perp : P_{\theta_{S-T-A,d}}(s_t^d | s_{t-1}^{d-1} r_t^d) \\ \text{if } r_t^d = \mathbf{r}_\perp : P_{\theta_{S-T-W,d}}(s_t^d | s_{t-1}^d r_t^{d+1}) \end{cases} \quad (37)$$

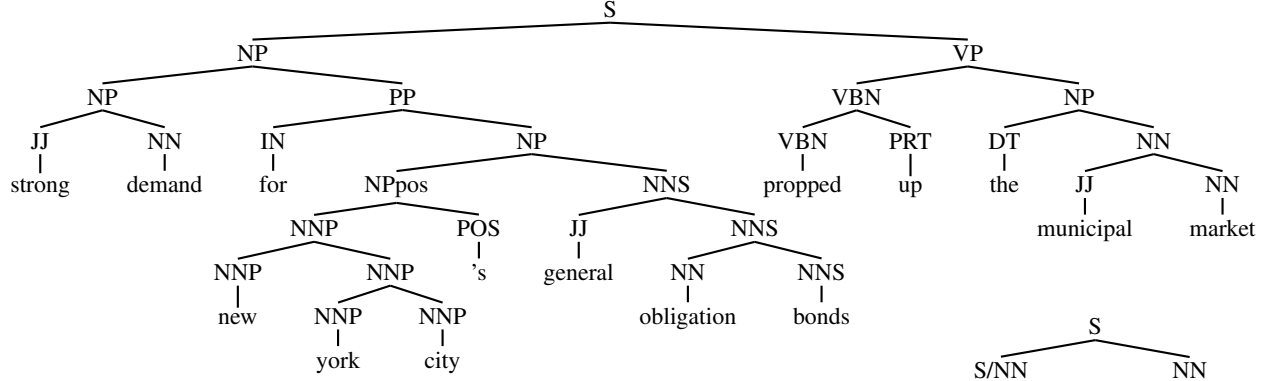
$$P_{\theta_{R-R,d}}(r_t^d | r_t^{d+1} s_{t-1}^d s_{t-1}^{d-1}) \stackrel{\text{def}}{=} \begin{cases} \text{if } c_{r_t^{d+1}} \neq x_t : \llbracket r_t^d = \mathbf{r}_\perp \rrbracket \\ \text{if } c_{r_t^{d+1}} = x_t : P_{\theta_{R-R,d}}(r_t^d | s_{t-1}^d s_{t-1}^{d-1}) \end{cases} \quad (38)$$

These right-corner parsing operations then construct a full inside probability decompositions, using Equations 8 through 13 and Equations 19 through 22, marginalizing out all constituents that are not required in each term:

- for expansions:

$$P_{\theta_{S-E,d}}(\langle c_{\eta_\nu}, c'_{\eta_\nu} \rangle | \langle -, c_{\eta_\nu} \rangle) \stackrel{\text{def}}{=} E_{\theta_{G-RL^*,d}}(c_{\eta_\nu} \xrightarrow{*} c_{\eta_\nu} \dots) \cdot \llbracket x_{\eta_\nu} = c'_{\eta_\nu} = c_{\eta_\nu} \rrbracket \quad (39)$$

a) binarized phrase structure tree:



b) result of right-corner transform:

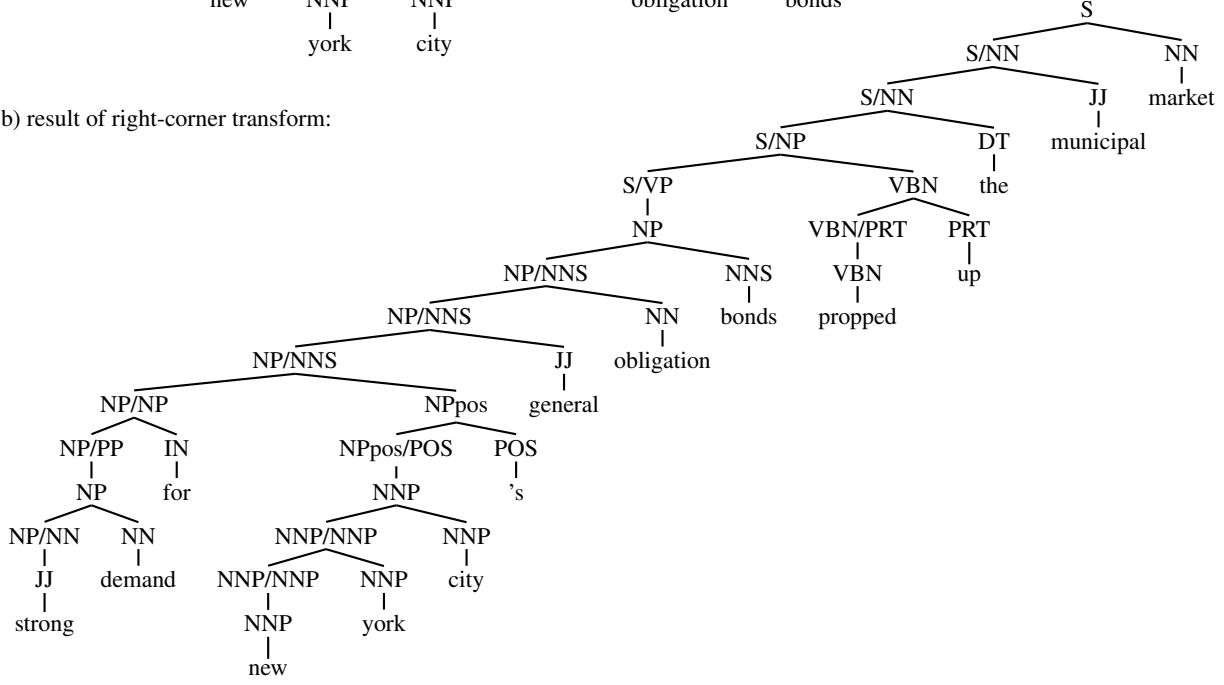


Figure 2: (a) Sample binarized phrase structure tree for the sentence *Strong demand for New York City's general obligations bonds propped up the municipal market*, and (b) a right-corner transform of this binarized tree.

- for awaited transitions, from Equation 11:

$$P_{\theta_{S-T-W,d}}(\langle c_\eta, c_{\eta l 1} \rangle | \langle c'_\eta, c_{\eta l} \rangle c_{\eta l 0}) \stackrel{\text{def}}{=} \llbracket c_\eta = c'_\eta \rrbracket \cdot \frac{P_{\theta_{G-R,d}}(c_{\eta l} \rightarrow c_{\eta l 0} c_{\eta l 1})}{E_{\theta_{G-RL^*,d}}(c_{\eta l} \xrightarrow{0} c_{\eta l 0} \dots)} \quad (40)$$

- for active transitions, from Equations 8 and 13:

$$\frac{P_{\theta_{S-T-A,d}}(\langle c_{\eta l}, c_{\eta l 1} \rangle | \langle -, c_\eta \rangle c_{\eta l 0}) \stackrel{\text{def}}{=} E_{\theta_{G-RL^*,d}}(c_\eta \xrightarrow{*} c_{\eta l} \dots) \cdot P_{\theta_{G-L,d}}(c_{\eta l} \rightarrow c_{\eta l 0} c_{\eta l 1})}{E_{\theta_{G-RL^*,d}}(c_\eta \xrightarrow{+} c_{\eta l 0} \dots)} \quad (41)$$

- for cross-element reductions:

$$P_{\theta_{R-R,d}}(c_{\eta l}, \mathbf{1} | \langle -, c_\eta \rangle \langle c'_\eta, - \rangle) \stackrel{\text{def}}{=} \llbracket c_{\eta l} = c'_\eta \rrbracket \cdot \frac{E_{\theta_{G-RL^*,d}}(c_\eta \xrightarrow{0} c_{\eta l} \dots)}{E_{\theta_{G-RL^*,d}}(c_\eta \xrightarrow{*} c_{\eta l} \dots)} \quad (42)$$

- for in-element reductions:

$$P_{\theta_{R-R,d}}(c_{\eta l}, \mathbf{0} | \langle -, c_\eta \rangle \langle c'_\eta, - \rangle) \stackrel{\text{def}}{=} \llbracket c_{\eta l} = c'_\eta \rrbracket \cdot \frac{E_{\theta_{G-RL^*,d}}(c_\eta \xrightarrow{+} c_{\eta l} \dots)}{E_{\theta_{G-RL^*,d}}(c_\eta \xrightarrow{*} c_{\eta l} \dots)} \quad (43)$$

A sample phrase structure tree is shown as a right-corner transformed recursive structure in Figure 2,

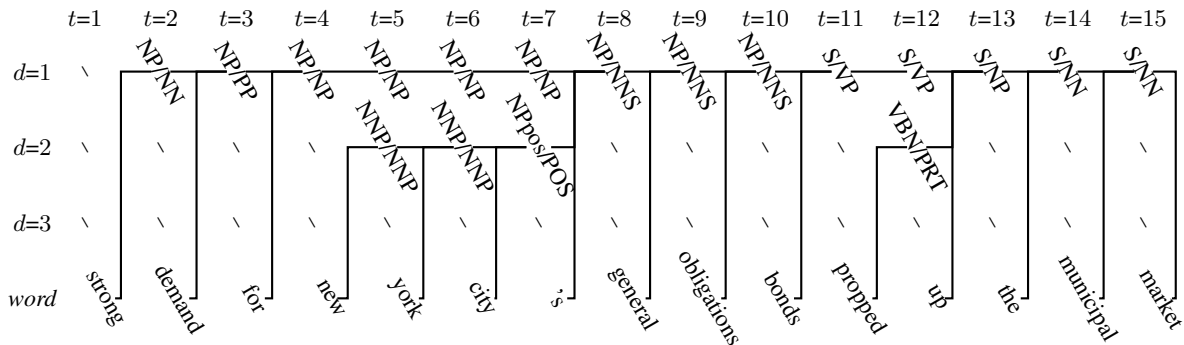


Figure 3: Sample tree from Figure 2 mapped to s_t^d variable positions of an HHMM at each depth level d (vertical) and time step t (horizontal). This tree uses at most only two memory store elements. Values for reduction states r_t^d are not shown.

and as a sequence of store states in Figure 3, corresponding to the output of Viterbi most likely sequence estimator. This estimation runs in linear time on the length of the input, so the parser can be run continuously on unsegmented or unpunctuated input. An ordinary phrase structure tree can be obtained by applying the transforms from Section 3, in reverse, to the right-corner recursive phrase structure tree represented in the sequence of store states.

5 Discussion

This paper has presented a derivation of a factored probabilistic sequence model from a probabilistic context-free grammar, using a right-corner transform to minimize memory usage in incremental processing. This sequence model characterization is attractive as a cognitive model because it does not posit any internal representation of complex phrasal structure beyond the pair of categories in each incomplete constituent resulting from the application of a right-corner transform; and because these incomplete constituents represent contiguous connected chunks of phrase structure, in line with characterizations of chunking in working memory (Miller, 1956). Experiments on large phrase-structure annotated corpora (Marcus et al., 1993) show this model could process the vast majority of sentences in a typical newspaper using only three or four store elements (Schuler et al., 2008; Schuler et al., 2010), in line with recent estimates of human short-term working memory capacity (Cowan, 2001).

This derivation can be applied to efficiently in-

crementalize any PCFG model, preserving the probabilistic dependencies in the original model. But, since the model is ultimately defined through transitions on entire working memory stores, it is also possible to relax the independence assumptions from the original PCFG model, and introduce additional dependencies across store elements that do not correspond to context-free dependencies. These additional dependencies might be used approximate dependencies of mildly context-sensitive grammar formalisms like tree adjoining grammars (Joshi, 1985), e.g. to model long-distance dependencies in filler-gap constructions (Kroch and Joshi, 1986), or crossed and nested dependencies in languages like Dutch (Shieber, 1985).

Figure 4 shows a sample store sequence corresponding to a parse of a noun phrase modified by an object relative clause. Here, the affix ‘-xNP’ is used to identify a phrase containing an extracted NP constituent, and the category label ‘GCnpS’ is used to identify the maximal projection of a gapped clause. If the GC constituent (and only this constituent) is associated with the referent or dependency information of the filler constituent (the bike in this example), this information can be made available during processing from the immediately superior incomplete constituent at the gap position in the relative clause, without passing dependency information down the tree as any kind of feature (Pollard and Sag, 1994), even though this position is not adjacent to the filler in the phrase structure tree (they are separated by the referent or dependency information of the bridge verb ‘said’). This is important to

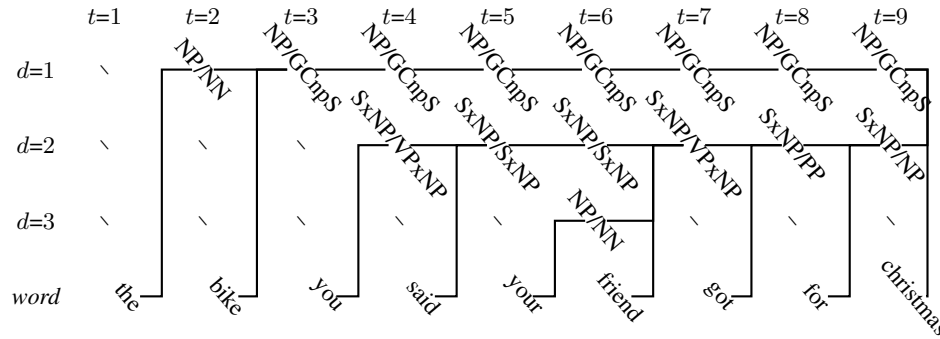


Figure 4: Sample store sequence containing long-distance dependency in a filler-gap construction.

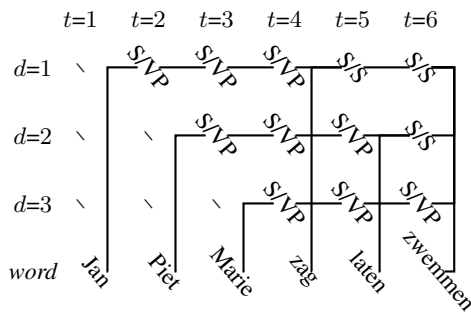


Figure 5: Sample store sequence containing crossed and nested dependencies.

satisfy claims that the model constructs chunks only for contiguous dependency structures (which is not true of propagated ‘slash’ features).

Figure 5 shows a sample store sequence corresponding to a parse of a Dutch sentence containing crossed and nested dependencies, featuring reductions across non-adjacent depth levels. This requires a more severe relaxation of PCFG independence assumptions, and is beyond the capacity of the HHMM as defined above, but this does preserve the notion of incomplete constituents and general composition established above, and is not beyond the capacity of factored sequence models with human-like memory bounds in general (note that the example sentence can still be parsed with only three store elements). This suggests a promising avenue of generalizing this model to learn parsing transitions that may be broader than those of a strict pushdown automaton.

References

Richard Bellman. 1957. *Dynamic Programming*. Princeton University Press, Princeton, NJ.

Nelson Cowan. 2001. The magical number 4 in short-term memory: A reconsideration of mental storage capacity. *Behavioral and Brain Sciences*, 24:87–185.

Aravind K. Joshi. 1985. How much context sensitivity is necessary for characterizing structural descriptions: Tree adjoining grammars. In L. Karttunen D. Dowty and A. Zwicky, editors, *Natural language parsing: Psychological, computational and theoretical perspectives*, pages 206–250. Cambridge University Press, Cambridge, U.K.

Anthony S. Kroch and Aravind K. Joshi. 1986. The linguistic relevance of tree adjoining grammars. *Linguistics and Philosophy*.

Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. 1993. Building a large annotated corpus of English: the Penn Treebank. *Computational Linguistics*, 19(2):313–330.

George A. Miller. 1956. The magical number seven, plus or minus two: Some limits on our capacity for processing information. *Psychological Review*, 63:81–97.

Kevin P. Murphy and Mark A. Paskin. 2001. Linear time inference in hierarchical HMMs. In *Proc. NIPS*, pages 833–840, Vancouver, BC, Canada.

Carl Pollard and Ivan Sag. 1994. *Head-driven Phrase Structure Grammar*. University of Chicago Press, Chicago.

William Schuler, Samir AbdelRahman, Tim Miller, and Lane Schwartz. 2008. Toward a psycholinguistically-motivated model of language. In *Proceedings of COLING*, pages 785–792, Manchester, UK, August.

William Schuler, Samir AbdelRahman, Tim Miller, and Lane Schwartz. 2010. Broad-coverage incremental parsing using human-like memory constraints. *Computational Linguistics*, 36(1).

William Schuler. 2009. Parsing with a bounded stack using a model-based right-corner transform. In *Proceedings of NAACL*, pages 344–352, Boulder, Colorado.

Stuart Shieber. 1985. Evidence against the context-freeness of natural language. *Linguistics and Philosophy*, 8:333–343.