# Parsing CFGs and PCFGs with a Chomsky-Schützenberger Representation<sup>\*</sup>

Mans Hulden

University of Helsinki Language Technology P.O. Box 24 FIN-00014 University of Helsinki, Finland mans.hulden@helsinki.fi

**Abstract.** We present a parsing algorithm for arbitrary context-free and probabilistic context-free grammars based on a representation of such grammars as a combination of a regular grammar and a grammar of balanced parentheses, similar to the representation used in the Chomsky-Schützenberger theorem. The basic algorithm has the same worst-case complexity as the popular CKY and Earley parsing algorithms frequently employed in natural language processing tasks.

# 1 Introduction

In this paper we will present a parsing method for probabilistic (PCFG) and non-probabilistic context-free grammars (CFG) based on a finite-state representation derived from the Chomsky-Schützenberger (C-S) Theorem (Chomsky and Schützenberger, 1963). In the representation, a context-free grammar G is constructed as a combination of a regular grammar R and a simple context-free grammar that consists only of balanced parentheses D. The resulting parser has the same cubic asymptotic complexity and bears many similarities to the classical Cocke-Kasami-Younger (CKY) algorithm (Younger, 1967) and the Earley parsing algorithm (Earley, 1968). It is also very simple to implement, assuming one has access to algorithms that perform a basic construction of finite automata. It may also yield more efficient parsing methods for other (P)CFG-based systems.

# 2 Preliminaries

#### 2.1 Notation

We will employ the standard notation of context-free grammars: a context freegrammar (CFG) G is a 4-tuple ( $\Sigma_{NT}, \Sigma_T, P, S$ ), where  $\Sigma_{NT}$  is the set of nonterminal symbols,  $\Sigma_T$  a set of terminal symbols, used in a set of production

<sup>\*</sup> This research has received funding from the European Commission's 7th Framework Program under grant agreement no. 238405 (CLARA).

Z. Vetulani (Ed.): LTC 2009, LNAI 6562, pp. 151–160, 2011.

<sup>©</sup> Springer-Verlag Berlin Heidelberg 2011

rules P of the form  $A \to \alpha$  where  $A \in \Sigma_{NT}$  and  $\alpha$  is a sequence drawn from  $(\Sigma_{NT} \cup \Sigma_T)^+$ . S is a symbol from  $\Sigma_{NT}$  designated as the start symbol. A CFG is in Chomsky Normal Form (CNF) if all productions are of the form  $A \to a$  or  $A \to BC$ , where  $a \in \Sigma_T$  and  $\{B, C\} \in \Sigma_{NT}$ .

Additionally, we will use short extended regular expressions to portray regular languages (and finite automata) and assume familiarity with the notational devices of  $L^*$ ,  $(L \cap L')$ ,  $\neg L$ ,  $(L \cup L')$ , LL', denoting the Kleene closure of a language, intersection, complement, union and concatenation of languages, respectively.

## 2.2 The Chomsky-Schützenberger Theorem

The essential difference between regular grammars and context-free grammars that of the possibility of self-embedding—is captured by the Chomsky-Schützenberger theorem, which essentially says every context-free grammar is a combination of local constraints on well-formedness (expressible as a regular grammar), and global constraints, which reduces to the idea that some elements must be properly nested.

More formally, the theorem states that for every CFL G, there exists regular language R and two homomorphisms g and h, such that

$$L(G) = h(g^{-1}(D) \cap R) \tag{1}$$

where D is a language consisting of different types of balanced parentheses (a Dyck language).

Although the original proof of the theorem did not do so, the language D can be connected directly to the context-free grammar in question and characterized in terms of parenthesis symbols that encode a derivation of strings in the contextfree language. The role of R is to restrict the occurrence of these parentheses locally, while D enforces proper nesting. Under such an interpretation, h is a homomorphism that deletes the parentheses, while g is a homomorphism that deletes actual terminal symbols  $\Sigma_T$ , and hence  $g^{-1}$  'inserts' terminal symbols in a language arbitrarily.

## 2.3 An Encoding for CFGs

A way of performing a C-S encoding of a CFL—one that also yields a constructive proof of the theorem—is to declare a parenthesis language D over an alphabet  $\Sigma_{()}$  that contains 2n symbols for every context-free rule in a grammar, where n is the number of symbols on the right hand side, and each such parenthesis represents a stage of derivation of rule. That is, for each rule of the form:

$$A \to \alpha$$
 (2)

we include symbols  $\begin{pmatrix} 1 \\ A \to \alpha \end{pmatrix}$ ,  $\begin{pmatrix} n \\ A \to \alpha \end{pmatrix}$ ,  $\begin{pmatrix}$ 

For example, if a grammar includes the rules

$$A \to BC, \ B \to x, \ C \to y$$
 (3)

we include symbols  $\binom{1}{A \to BC}$ ,  $\binom{1}{A \to BC}$ ,  $\binom{2}{A \to BC}$ , and  $\binom{2}{A \to BC}$  for the first rule (and likewise for the subsequent ones) and encode a derivation:

$$\begin{array}{c}
A \\
B C \\
| | \\
x y
\end{array}$$
(4)

as the string

$$\binom{1}{A \to BC} \binom{1}{B \to x} \mathbf{x} \binom{1}{B \to x} \binom{1}{A \to BC} \binom{2}{A \to BC} \binom{1}{C \to y} \mathbf{y} \binom{1}{C \to y} \binom{2}{A \to BC}$$

That is, we represent the parse of a string as the bracketed preorder traversal of the derivation tree.

The connection to the above theorem is that if we declare a homomorphism h to be  $\Sigma_{()} \to \epsilon$  and apply it to the above string, the result is naturally a string over only the terminal symbols, namely xy.

Given this, a CFG G with terminal symbols  $\Sigma_T$  and rules of the form  $A \to \alpha$ , where  $\alpha$  is an arbitrary string of terminals and nonterminals, can be represented as  $h(g^{-1}D_{()} \cap R)$ , where  $D_{()}$  is the language of balanced parentheses over the set of parenthesis symbols  $\Sigma_{()}$ , and R an intersection of five regular languages over  $\Sigma = (\Sigma_T \cup \Sigma_{()})$ , specified as follows:

(a) 
$$\begin{pmatrix} 1 \\ S \to \alpha} \Sigma^*$$
 where  $S$  is the initial symbol in  $G$   
(b)  $\begin{pmatrix} i \\ A \to B_1 \dots B_n} \rightsquigarrow \begin{pmatrix} i \\ B_i \to C_1 \dots C_n \end{pmatrix}$  for  $B_i \in \Sigma_{NT}$   
(c)  $\begin{pmatrix} i \\ A \to B_1 \dots B_n \end{pmatrix} \rightsquigarrow \begin{pmatrix} i+1 \\ A \to B_1 \dots B_n \end{pmatrix}$  for  $i < n$   
(d)  $\begin{pmatrix} n \\ A \to B_1 \dots B_n \end{pmatrix} \rightarrow \begin{pmatrix} i \\ A \to B_1 \dots B_n \end{pmatrix} \vee \begin{pmatrix} i \\ A \to B_1 \dots B_n \end{pmatrix}$  for  $i = n$   
(e)  $\begin{pmatrix} i \\ A \to B_1 \dots B_n \end{pmatrix} \sim B_i \begin{pmatrix} i \\ B_i \end{pmatrix} = B_i \dots B_n$  for  $B_i \in \Sigma_T$ 

Here, the notation  $x \rightsquigarrow y$  denotes the idea that any instance of a symbol x must be immediately followed by y. Such constraints are clearly expressible as regular languages. In fact  $x \rightsquigarrow y$  can be considered shorthand for the extended regular expression:

$$\neg(\Sigma^* x \neg (y \Sigma^*))$$

Also, the abstract symbol # denotes the end-of-string. The homomorphism g is simply  $\Sigma_T \to \epsilon$  (delete terminal symbols), and h is  $\Sigma_{()} \to \epsilon$  (delete parentheses). The different regular languages which are intersected encodes strictly local requirements on the ordering of the parenthesis symbols: (b) requires that an open parenthesis be immediately followed by another open parenthesis representing the first nonterminal on the right hand side; (c) enforces that a closing parenthesis representing a nonfinal constituent in a rule be followed immediately by the opening parenthesis for the same rule; (d) says that any parenthesis representing any final constituent in a rule be followed either by any other closing parenthesis or the end-of-word; (e) says that a parenthesis representing a rule constituent yielding a terminal symbol be followed by that terminal symbol and a closing parenthesis. Additionally constraint (a) enforces that any string begin with the start symbol.<sup>1</sup>

The constraints as given above would then apply to the example tree (4) and its string encoding as follows:



## **3** Parsing with C-S Representations

Having in this way made the Chomsky-Schützenberger theorem more concrete by encoding the parse of a sentence as a string, we can see that this immediately yields a parsing algorithm for context-free languages.

Suppose that we have a context-free grammar G and construct from it the regular grammar R as described in steps (a)–(e) above, and wish to parse a sentence  $w_1w_2 \ldots w_n$ . Now we can easily construct a finite-state automaton  $R^w$  encoding  $h^{-1}(w_1w_2 \ldots w_n)$ , that is, an automaton that accepts only the sentence to be parsed, with arbitrary sequences of parentheses interspersed (see figure 1, step 1). This is clearly a matter of first constructing an automaton that accepts only the sentence  $w_1w_2 \ldots w_n$ , and subsequently enhancing the automaton with self-loops for each symbol in  $\Sigma_{()}$ . Now, we can calculate the new finite automaton  $R^{local} = R^w \cap R$  (figure 1 step 2), that accepts all the locally correct parses of the sentence at hand. Obviously  $R^{local}$  overgenerates in the sense that it may contain invalid parses where parenthesis symbols are out of alignment, i.e. not properly nested. However, if we from the automaton representing  $R^{local}$  extract only the set of words where the parentheses are properly nested, this set equals precisely the correct parses of the sentence  $w_1w_2 \ldots w_n$  in question with respect to the grammar G.

To sum up, the steps in figure 1 are:

- (1) Calculate  $R^w = h^{-1}(w)$
- (2) Calculate  $R^{local} = R \cap h^{-1}(w)$
- (3) Extract from (2) the set of words where parentheses are balanced

Step (3), extracting from  $R^{local}$  only those words where parentheses are properly aligned is addressed by algorithm 1.

<sup>&</sup>lt;sup>1</sup> The proof that the constraints are both necessary and sufficient for, together with D, encoding the CFG is a fairly simple induction, see e.g. Salomaa (1973) or Kozen (1997) for proofs with similar representations.

#### 3.1 The Algorithm

The objective of algorithm 1 is to extract all the paths that contain only balanced parentheses from the finite-state automaton produced by step (2). To this end, we maintain an agenda A which contains state pairs. Initially, the only pairs in Aare those where there is a transition from one state to the other with a terminal symbol. From the agenda A we choose a state pair and expand it to produce a new state pair if there are transitions on the left and right with balanced parentheses. This is done in lines 14–16. We also need to merge pairs  $P_1$  and  $P_2$  if they represent parts of the same constituent and  $P_1$  forms a word  $\binom{i}{\alpha} \dots \binom{i}{\alpha}$ and  $P_2 (^j_{\alpha} \dots)^j_{\alpha}$  for some rule  $\alpha$  by checking if  $P_1$  and  $P_2$  share a state on the left or right. This is done in lines 8–13. Note that we do not explicitly need to check the contents of the words formed by  $P_1$  and  $P_2$  or that the indices i and j are such that j = i + 1 since this has already been taken care of by the local grammar. Hence, the finite automaton from which we are extracting the paths containing balanced parentheses will never contain adjacent pairs of states in such a configuration without the indices and rule components being compatible. In other words, we need only look at the state numbers to perform the joining of constituents. Finally, whenever we encounter a state pair such that one state is the initial state and the other a final state, we have found a parse. It is assumed that in lines 17-19 we maintain backtraces of the pairs we added to the agenda A so that the string representing the parse can be reconstructed as one is found.

#### 4 Analysis

Before we move on to consider enhancements to the parsing algorithms, let us first briefly analyze the complexity of parsing a sentence.

First, let us represent the size of the local grammar R, which depends on G, by some constant c. It is worth noting that the constraints (a)–(e) which R is constructed from are all strictly local, maximally 3-testable languages. This leads to that the size of R grows additively with each grammar rule.<sup>2</sup>

The task of constructing from a sentence to be parsed the automaton  $R^w$  (step 1) takes time proportional to the length of the sentence, i.e. |w|. When intersecting two finite automata the result grows as the product of the two. In this case, since |R| = c and  $R^w$  has |w| states, we have that steps (1) and (2) are of time complexity c|w|, where c depends on the grammar, i.e. O(|w|).

For the analysis of the algorithm for EXTRACTBALANCED, let us restrict ourselves to cases where the grammar is given in Chomsky Normal Form. This is reflected in the automaton produced by step (2) in such a way that the maximum index i of a bracket  $\binom{i}{R}$  representing a rule is 2. Interestingly, we can thus

<sup>&</sup>lt;sup>2</sup> This additive growth is important since it is indeed possible to construct many different types of local grammars R that work correctly in tandem with the balanced-parenthesis-grammar D. However, not all of them exhibit subexponential or linear complexity as the construction in this paper does.

Algorithm 1. ExtractBalanced
<b>Input</b> : $FSM = (V, E)$
begin
<b>foreach</b> $(p,q) \in V$ where there is a transition $p \xrightarrow{s} q$ and $s \in \Sigma_T$ do
add $(p,q)$ to A as unmarked
end
while there is an unmarked pair in $A$ do
Choose a pair $P = (p, q)$ from A
$\mathrm{Mark}\ P$
if exists a pair $Q = (p', q')$ in A such that $q' = p$ then
add $(p',q)$ to A as unmarked
end
if exists a pair $Q = (p', q')$ in A such that $p' = q$ then
add $(p,q')$ to A as unmarked
end
if exists transitions $(p' \xrightarrow{s} p)$ , $(q \xrightarrow{t} q')$ , where $s = ({}^i_B, t =)^j_B$ and $i =$
then
add $(p',q')$ to A as unmarked
end
if $p \in V_{start}$ and $q \in V_{final}$ then
Return(Backtrace(P))
end
end
end

j

find an isomorphism between the algorithm that locates paths that contain balanced parentheses in the automaton and the CKY algorithm. That is, each pair p, q added at lines 8–13 represents a complete subtree over some span i, j for the sentence to be parsed. Obviously, the *if* checks at lines 8-17 take O(1) time assuming some suitable indexing of the pairs stored. For any subparse of a subword of length n there are n-1 ways of breaking it up into two constituents  $\binom{1}{R}$ ... $\binom{1}{R}$  $\binom{2}{R}$ ... $\binom{2}{R}$ . This means each possible pair (of which there are  $n^2$ , n being the number of states in the automaton) added to A may be added O(n) times. The total complexity is then  $O(n^3)$ , and given that the size of the automaton nis proportional to |w|, also  $O(|w|^3)$ .

In the above, we are not committing ourselves to any particular queuing strategy for the state pairs in A. There are many options available of how to proceed in this respect. For instance, lines 2–4 immediately add all the spans representing the terminals on the agenda after which we subsequently iterate the search until the agenda is empty. This is not strictly necessary as one can proceed left-toright in the parse by only adding the states representing the first two terminals in the strings, and then, once A is empty, add the third, etc. This strategy would correspond somewhat to Earley's algorithm (Earley, 1968), and could avoid some extra work in that it may reduce the construction of subparses (state pairs) that are not actually parts of complete parses.



Fig. 1. Basic workflow for parsing a CFG or PCFG

### 4.1 Weights and PCFGs

In many cases we would like to include treatment of probabilistic grammars in such a way that each production  $A \to \alpha$  has a weight associated with it, and each parse therefore also has a total probability associated with it. This is simply a matter of storing with each A, an associated probability, or cost, and handling the maintenance of the associated subparses accordingly. There are several potential strategies regarding how this could be handled. We have chosen in our implementation to mark certain transitions in the automaton with a weight or probability. More specifically, each transition with a symbol ( $^1_R$  carries the weight/cost of rule R, other transitions carrying cost 0. Obviously, every correct parse will always include the first opening bracket symbol for each constituent, and hence it is sufficient to mark only these.



Fig. 2. An example of R, a simple CFG approximation

Naturally, if we are only interested in a single most likely parse, we can include a Viterbi approximation and for each A on the agenda, store only the one with lowest probability (recall from the above that each state pair represents a combination of a span and constituent).

## 5 An Example

Let us illustrate the parsing method by a simple grammar with three production rules

$$P = \{S \to SS, S \to a, S \to b\}$$

and a simple string to be parsed: "aaa"

Figure 2 shows the automaton representing the regular approximation R produced by intersection rules (a)–(e). The automaton  $R^w$  is simply a four-state machine that accepts the string "aaa" with any number of bracket symbols intervening, drawn from the bracket alphabet:

$$\{\binom{1}{S \to SS}, \binom{2}{S \to SS}, \binom{1}{S \to SS}, \binom{1}{S \to SS}, \binom{1}{S \to a}, \binom{1}{S \to a}, \binom{1}{S \to b}, \binom{1}{S \to b}\}$$

The resulting automaton  $R^{local} = R^w \cap R$  still contains an arbitrary number of strings, from which the algorithm EXTRACTBALANCED draws out the only two legal parse strings:

$$\begin{pmatrix} 1\\S \to SS \end{pmatrix} \begin{pmatrix} 1\\S \to a \end{pmatrix} = \begin{pmatrix} 1\\S \to a \end{pmatrix} \begin{pmatrix} 1\\S \to SS \end{pmatrix} \begin{pmatrix} 2\\S \to SS \end{pmatrix} \begin{pmatrix} 1\\S \to SS \end{pmatrix} \begin{pmatrix} 1\\S \to a \end{pmatrix} \begin{pmatrix} 1\\S \to a \end{pmatrix} \begin{pmatrix} 2\\S \to SS \end{pmatrix} \begin{pmatrix} 2\\S \to SS \end{pmatrix} \begin{pmatrix} 1\\S \to a \end{pmatrix} \begin{pmatrix} 1\\S \to a \end{pmatrix} \begin{pmatrix} 2\\S \to SS \end{pmatrix} \begin{pmatrix} 2\\S \to SS \end{pmatrix} \begin{pmatrix} 1\\S \to a \end{pmatrix} \begin{pmatrix} 2\\S \to SS \end{pmatrix} \begin{pmatrix} 2\\S \to$$

and

 $\binom{1}{S \to SS} \binom{1}{S \to SS} \binom{1}{S \to a} \binom{1}{S \to a} \binom{1}{S \to a} \binom{2}{S \to SS} \binom{1}{S \to a} \binom{1}{S \to a} \binom{1}{S \to a} \binom{2}{S \to SS} \binom{1}{S \to a} \binom{2}{S \to SS} \binom{1}{S \to a} \binom{1}{S \to a}$ 

representing the parse trees in figure 3 (a) and (b), respectively.

Fig. 3. Two parse trees from the example grammar

# 6 Conclusion and Further Work

We have presented an overall strategy for parsing CFGs and PCFGs through a representation of a CFG as a combination of a regular grammar and a grammar of balanced parentheses. Using this encoding, we present a general cubic-time algorithm for parsing of arbitrary context-free grammars, which can also be used for parsing PCFGs. The resulting algorithm is simple to implement assuming one can construct finite-state automata from regular expressions.<sup>3</sup>

Apart from its direct usability, we also expect the basic approach to be applicable to a number of recent approaches that use CFG representations for NL tasks such as lexicalized PCFGs (Charniak, 1997) or bilexical grammar parsing tasks (Eisner, 1997), to name a few.

# References

- Charniak, E.: Statistical parsing with a context-free grammar and word statistics. In: Proceedings of the 14th National Conference on Artificial Intelligence, pp. 598–603 (1997)
- Chomsky, N., Schützenberger, M.-P.: The algebraic theory of context-free languages. In: Braffort, P., Hirschberg, D. (eds.) Computer Programming and Formal Systems, pp. 118–161. North Holland, Amsterdam (1963)
- Earley, J.: An efficient context-free parsing algorithm. PhD thesis, Carnegie-Mellon University, Pittsburgh, Pa (1968)
- Eisner, J.: Bilexical grammars and a cubic-time probabilistic parser. In: Proceedings of the 1997 International Workshop on Parsing Technologies (1997)
- Hulden, M.: Foma: a finite-state compiler and library. In: Proceedings of EACL 2009, pp. 29–32 (2009)
- Kozen, D.C.: Automata and Computability. Springer, Heidelberg (1997)
- Salomaa, A.: Formal Languages. Academic Press, New York (1973)
- Younger, D.H.: Recognition and parsing of context-free languages in time  $n^3$ . Information and Control 10, 189–208 (1967)

 $<sup>^3</sup>$  We have made an implementation of the basic algorithm for CFGs and a Viterbi PCFG-parser using the finite-state toolkit *foma* (Hulden, 2009).