

Elizabeth Scott Explained

Parsing from Earley Recognisers

Zoe Wheeler

University of Texas at Austin
zoe.donnellon.wheeler@gmail.com

Walter Xia

University of Texas at Austin
swilery@utexas.edu

Abstract

Earley's Algorithm is able to recognize general context-free grammars in $O(n^3)$, where n is the size of the string to be recognized. However, there are times in which we want more than just a yes or no answer. There are times in which we want an actual parse tree, and for ambiguous grammars, there are times in which we want all possible parse trees. Fortunately, there is a paper by Dr. Elizabeth Scott, [2], that presents a technique to produce a data structure known as a Shared Packed Parse Forest (SPPF), able to represent even an infinite number of parse trees. Unfortunately this paper is poorly written, making it very difficult to understand. Our paper is a re-explanation of Scott's techniques. It is agreed by many that Earley's Algorithm is also difficult to understand. Fortunately, there exists a data structure due to Dr. Gianfranco Bilardi and Dr. Keshav Pingali, [1], known as Grammar Flow Graphs (GFGs) that significantly ease the understanding of the algorithm by reformulating parsing problems as path problems in a graph. Our technique will use GFGs.

Categories and Subject Descriptors F.7.2 [Semantics and Reasoning]: Program Reasoning–Parsing

General Terms Context-Free Languages, Cubic Generalized Parsing, Earley Parsing

Keywords Earley Sets, Grammar Flow Graphs, Non-Deterministic Finite Automaton, Shared Packed Parse Forest

1. Introduction

It is important here for us to distinguish between recognisers and parsers for a grammar. Recognizers determine whether or not a string is part of a language defined by a grammar whereas parsers construct parse trees that reveal *how* a string satisfies the syntax dictated by a grammar. For about the past five decades, there already exist general recognizers like Cocke-Younger-Kasami (CYK) and Earley's Algorithms that run cubic relative to the size of the string to be recognized. Alternatively, Generalized LR (GLR) is an algorithm that produces parsers but has the very undesirable property that it is unbounded. Dr. Elizabeth Scott extended the Earley Recogniser into a parser that is able run in cubic space and time,

[2]. The challenge was to successfully apply the parser to ambiguous grammars that produces multiple, perhaps infinite, parse trees for a string in the grammar. Note that simply disallowing ambiguous grammars is not a solution since there exists grammars that are intrinsically ambiguous. The solution she used was a representation known as a Shared Packed Parse Forest (SPPF), which is in essence a Directed Acyclic Graph (DAG).

Earley's Algorithm is a highly complex algorithm. To dramatically simplify its understanding, we view it from the perspective of Grammar Flow Graphs (GFGs) that restructure parsing as finding certain paths within the graph, [1]. For those of you familiar with automata theory, GFGs play the same role for context-free grammars as finite-state automata play for regular grammars. The rest of the paper is organized as follows:

- Section 2 will introduce GFGs
- Section 3 will introduce Earley's Algorithm using GFGs
- Section 4 will introduce SPPFs
- Section 5 will introduce Dr. Scott's Algorithm for producing SPPFs
- Section 6 will discuss our implementation
- Section 7 will discuss our results
- Section 8 will conclude

2. Grammar Flow Graphs

Let us begin with the standard definition of a context-free grammar.
Definition: A context-free grammar, CFG, is a tuple (N, T, P, S) , where, [1]:

- ▷ N is a finite set of elements called *nonterminals*,
- ▷ T is a finite set of elements called *terminals*,
- ▷ $P \subseteq N \times (N \cup T)^*$ is the set of *productions* that map nonterminals to a sequence of nonterminals or terminals, and
- ▷ $S \in N$ is the unique *start symbol* that appears once on the left-hand side of a single production.

An example of a grammar is the following, where $|$ signifies or:

$$S \longrightarrow N t \mid t N$$

$$N \longrightarrow t t$$

Now we are in a position to introduce the GFG.

Definition: Let $CFG = (N, T, P, S)$ be a context-free grammar and let ϵ denote the empty string. The *grammar flow graph* (GFG) of CFG, $GFG(CFG) = (V(CFG), G(CFG))$, is the smallest directed graph that has the following properties, [1]:

- ▷ For each nonterminal $M \in N$, there exist $(\bullet M), (M \bullet) \in V(CFG)$ called *start nodes* and *end nodes* respectively,

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

CONF 'yy, Month d-d, 20yy, City, ST, Country.

Copyright © 20yy ACM 978-1-nnnn-nnnn-n/yy/mm...\$15.00.

<http://dx.doi.org/10.1145/nnnnnnn.nnnnnnn>



- Figure 1 depicts the GFG associated with the preceding grammar. The following definition comes naturally.
- Definition:* A path in a GFG *generates* the word w by concatenating the labels along its sequence of edges.

Proof: Please see [1].

$$\begin{array}{l} S \longrightarrow N t \mid t N \\ N \longrightarrow t t \end{array}$$

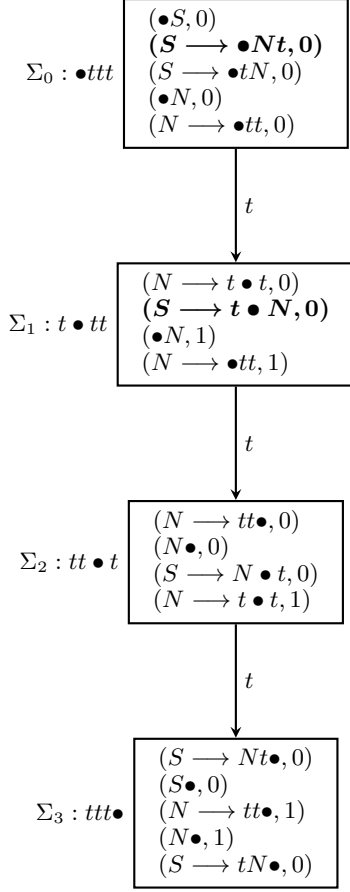


Figure 2. Example of Earley Sets for the preceding grammar on string "ttt". Call nodes are in bold.

4. Shared Packed Parse Forest (SPPF)

Even though a Shared Packed Parse Forest was the purpose of her algorithm, Dr. Elizabeth Scott only gave it a brief description. Here, we will try to provide a more intuitive understanding of this data structure. First, consider the following ambiguous grammar:

$$S \rightarrow S S \mid u$$

Figure 3 displays the two parse trees of the string "uuu" that is part of the above grammar. The meaning of the numbers at each node are as follows. Giving the string "uuu", number the string as shown:

$$0 \quad u \quad 1 \quad u \quad 2 \quad u \quad 3$$

So $u(0, 1)$ indicates that the u of interest is between the numbers 0 and 1, $u(1, 2)$ indicates that the u of interest is between the numbers 1 and 2, and so on. The numbers of the interior nodes are just natural extensions of this pattern.

Let us ignore these numbers for the moment and just consider the most natural way that allows us to share the common subtrees that appear in both parse trees in Figure 3. Notice that both subtrees consist of a root node S whose children consist of two other nodes, let's call them S_l and S_r for the moment. The difference between the two parse trees only begins to appear at level 3 in which we have two alternatives:

- Either S_l contains two children, S_m, S_s and S_r goes directly to a leaf u

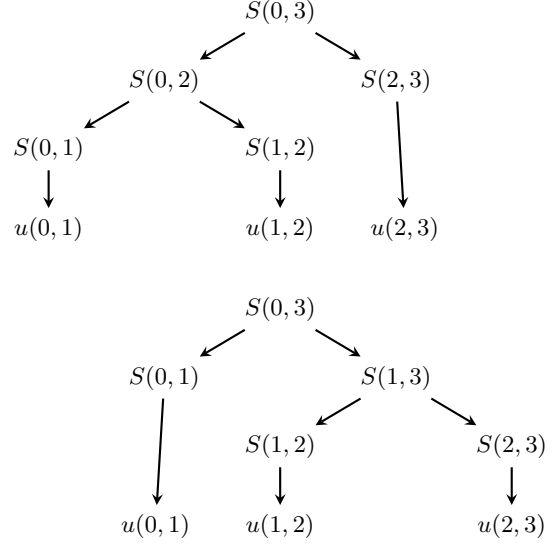


Figure 3. The two parse trees of "uuu" for the preceding grammar

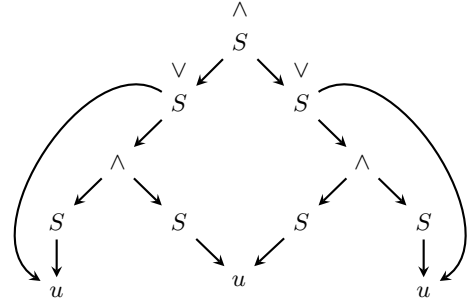


Figure 4. The and-or tree representation of the two parse trees in Figure 3.

- or S_l goes directly to a leaf u and S_r contains two children, S_m, S_s .

We can capture this intuitive notion by the *and-or* tree in Figure 4, where the subscripts are ignored since they are immaterial. An *and* node represents the action that you consider both children of the node while an *or* node represents the action that you consider one child of the node. However, observe that this tree represents strings other than "uuu". To see this, consider starting at the root node and following its edges to both of its children. Here, you are at two *or* nodes both labeled S . Since we are at *or* nodes, we choose a single child for each node. Suppose we choose the leaf children, resulting in the string "uu". To correct this error, we reintroduce the numbers associated with each node and require that only nodes with identical names and numbers can share structure from the tree. Figure 5 shows our correction.

The data structure we have constructed is what Dr. Elizabeth Scott calls a Shared Packed Parse Forest (SPPF). For sake of completeness, we present her definition. Note that in her definition, her *packed* nodes are our nameless *and* nodes.

Definition: A *Shared Packed Parse Forest (SPPF)* is a representation that reduces the space used to represent multiple parse trees of an ambiguous string. Nodes are named (x, j, i) , when node x matches substring $a_{j+1} \dots a_i$. Nodes which contain the same sub-

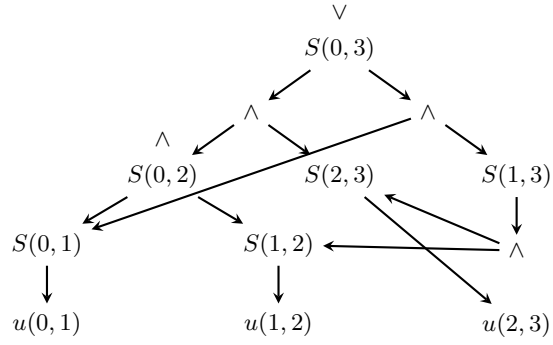


Figure 5. The SPPF representation of the two parse trees in Figure 3.

tree are shared and nodes which represent the same substring for nonterminals with different parse trees are packed.

A. Appendix Title

This is the text of the appendix, if you need one.

Acknowledgments

We would like to thank Dr. Keshav Pingali for his guidance on this project. We would also like to thank Sepideh Maleki for helping us understand Scott's algorithm.

References

- [1] Gianfranco Bilardi, and Keshav Pingali. Parsing with Pictures. UTCS Tech Reports, 2012. This is a full TECHREPORT entry.
- [2] Elizabeth Scott. SPPF-Style Parsing From Earley Recognisers. *Electronic Notes in Theoretical Computer Science*, 203(53-67), 2008. This is a full ARTICLE entry.