

An Introduction To Parsing

Wednesday, November 14, 2007

Reading: Stoughton 4.3

CS235 Languages and Automata

Department of Computer Science
Wellesley College

The Parsing Problem

Given a context-free grammar G and a string s :

- if $s \in \text{Lang}(G)$, return a parse tree in G that yields s ;
- if $s \notin \text{Lang}(G)$, reject it.

Efficient solutions to this problem are important for analyzing both programming languages and natural languages.

The Parsing Problem: Examples

Example Grammar:

$S \rightarrow T \mid OS1$
$T \rightarrow \% \mid 10T$

Example Strings:

000111 011011

101010 001011

010101 011001

Parsing Intro 30-3

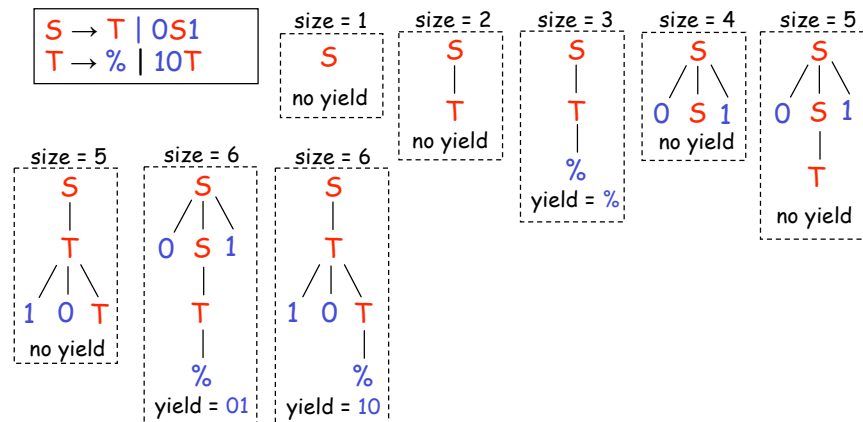
Goals For the Next Few Lectures

- Begin with some simple, but inefficient, parsing algorithms that work for all grammars.
- Study some more efficient approaches that work for various kinds of restricted grammars
- Build our way up to grammars suitable for a parser-generator program (ML-Yacc)

Parsing Intro 30-4

Naive Algorithm : Top-Down Generate and Test

Idea: generate all parse trees top down from the start variable in order of size (# of nodes or height) until find one that yields s .



How to know when to reject? Helps to have grammar in Chomsky Normal Form or Greibach Normal form. Why?

Parsing Intro 30-5

A Relevant Joke

An **engineer**, a **physicist** and a **mathematician** are staying in a hotel.

The **engineer** wakes up and smells smoke. He goes out into the hallway and sees a fire, so he fills a trash can from his room with water and douses the fire. He goes back to bed.

Later, the **physicist** wakes up and smells smoke. He opens his door and sees a fire in the hallway. He walks down the hall to a fire hose and after calculating the flame velocity, distance, water pressure, trajectory, etc. extinguishes the fire with the minimum amount of water and energy needed.

Later, the **mathematician** wakes up and smells smoke. He goes to the hall, sees the fire and then the fire hose. He thinks for a moment and then exclaims, "Ah, a solution exists!" and then goes back to bed.

(From <http://www.math.utah.edu/~cherk/mathjokes.html>)

Parsing Intro 30-6

Problems with Top-Down Generate and Test

Although it works for any grammar, top-down generate and test has big problems:

- Extremely inefficient approach to parsing.
- Repeat much of the same work for any string.
- Never take advantage of the structure of the string!

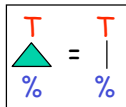
Parsing Intro 30-7

Working Bottom-Up: The Enzyme* Algorithm

Idea: Build up all parse trees that can match substrings in the given string s . There are only a finite number of these. Accept if there is a parse tree for s rooted at the start variable. (This is the algorithm presented in Stoughton 4.3.)

grammar G
$$\begin{array}{l} S \rightarrow T \mid OS1 \\ T \rightarrow \% \mid 10T \end{array}$$
 string s 001011

Step 1: Find all productions that yield $\%$ or some symbol in s .



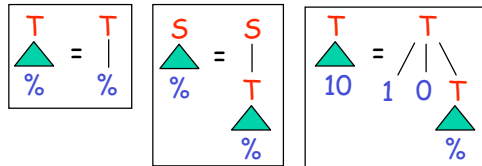
Step 2: Based on productions in G , add parse trees that yield some substring of s .

* "Enzyme" algorithm is my term. It is not standard.

Parsing Intro 30-8

Enzyme Example: The Next Round

grammar G $\begin{matrix} S \rightarrow T \mid OS1 \\ T \rightarrow \% \mid 10T \end{matrix}$ string s 001011

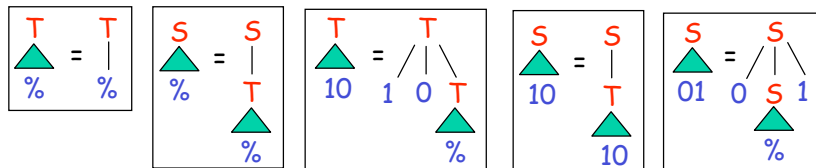


* "Enzyme" algorithm is my term. It is not standard.

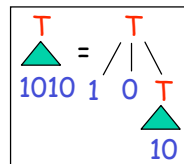
Parsing Intro 30-9

Enzyme Example Continued

grammar G $\begin{matrix} S \rightarrow T \mid OS1 \\ T \rightarrow \% \mid 10T \end{matrix}$ string s 001011



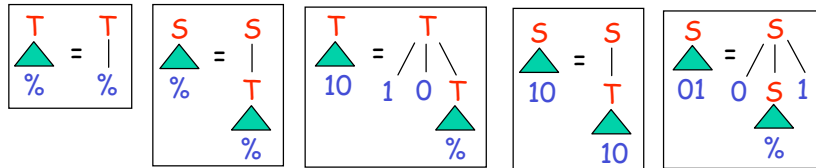
Note: we do not include this tree because 1010 isn't a substring of s .



Parsing Intro 30-10

Finishing The Enzyme Example

grammar G $\begin{matrix} S \rightarrow T \mid OS1 \\ T \rightarrow \% \mid 10T \end{matrix}$ string s 001011



Flesh out the remaining trees in this example:

Parsing Intro 30-11

Problems with Enzyme Algorithm

Although it works for any grammar G and takes input string s into account, enzyme algorithm still can do much unnecessary work generating parse trees for substrings that can never appear in the final parse tree.

Parsing Intro 30-12

