

Week 4: Book

Wednesday, February 9, 2022 4:20 PM

ϵ (epsilon) : lexical component, anything that you can write in a language.
 ϵ stands for "empty string of symbols"

Parsing : is the problem of taking a string of terminals to figure out how to derive it from the start symbol of the grammar or report syntax errors.

→ We begin with a simple source program:

$9-5+2$ } each character is a TERMINAL

the token's 1st component (token name) are terminals processed by the PARSER.

In reality, a program has multichar lexemes that the lexical converts to tokens

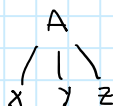
(2.2.1 Derivation)

↑
2.2.3 Parse Trees (page 45)

Parse Tree: Pictorially shows how the START SYMBOL of a grammar derives a string in the language.

→ ex. If NONTERMINAL "A" has a production:
 $A \rightarrow XYZ$

Then a parse tree has an interior NODE A, with 3 children labeled X, Y and Z from left to right:



Given a context-free grammar, a parse tree according to the grammar has the properties:

1. The root is labeled by the START SYMBOL.
 2. Each leaf is labeled by a TERMINAL or by ϵ .
 3. Each interior node is labeled by a NONTERMINAL.
 4. If A is the NONTERMINAL labeling of an interior node and x_1, x_2, \dots, x_n are the labels of the children of that node from $L \rightarrow R$, then there must be a production $A \rightarrow x_1 x_2 \dots x_n$ (x_i is a symbol either T or NONT)
- if $A \rightarrow \epsilon$ is a production, then a node labeled A has a single child ϵ .

Example: The derivation of $9-5+2$ is the tree below

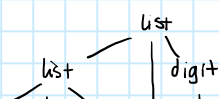
(2.4) Each node in the tree is labeled by a grammar symbol.

- A interior node and its children correspond to a production.

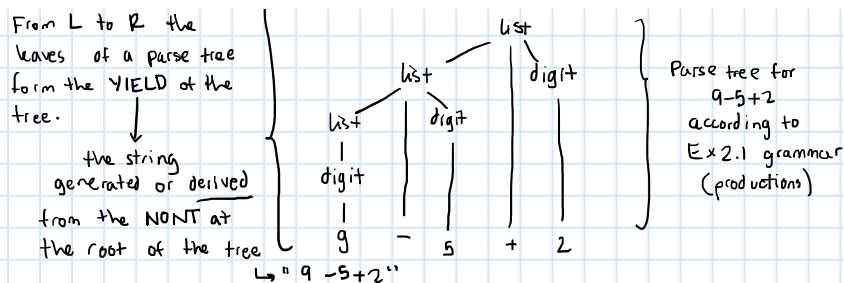
HEAD of the production
→ BODY of the production

- If we got $list \rightarrow list + digit$
 - ✓ the root is labeled list, which is the START SYMBOL of the grammar in Example 2.1
 - ✓ the children of the root are labeled, $L \rightarrow R$:
list, +, digit.

From L to R the leaves of a parse tree form the YIELD of the



Parse tree for $9-5+2$



this way. Any tree imparts a natural left-to-right order to its leaves, based on the idea that if X and Y are two children with the same parent, and X is to the left of Y , then all descendants of X are to the left of descendants of Y .

\rightarrow Another definition for "the language generated by the grammar" is "the set of strings that can be generated by a parse tree"

Parsing a string: the process of finding a parse tree for a given string of terminals.

2.2.4 Ambiguity

- \hookrightarrow the structure of a string according to a grammar
- \hookrightarrow a grammar can have more than 1 parse tree that generates (YIELD) a string of terminals. } This is an ambiguous grammar

- \hookrightarrow to show that a grammar is ambiguous:
 - find a terminal string (YIELD) that is the YIELD of more than 1 parse tree = if a string can be yielded by 2+ trees.

- \hookrightarrow Since a string with more than 1 parse tree can have \neq meanings, it is recommended not to design ambiguous grammars.

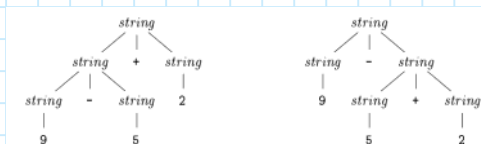
Example 2.5: Suppose we used a single nonterminal *string* and did not distinguish between digits and lists, as in Example 2.1. We could have written the grammar

$string \rightarrow string + string \mid string - string \mid 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

with this grammar (productions), $9+5-2$ has more than 1 tree.

The 2 trees correspond to each:

$$6 \leftarrow (9+5)-2 \qquad 9-(5+2) \rightarrow 2$$



2.2.5 Associativity of Operators

$$9+5+2 = (9+5)+2$$

$$9-5-2 = (9-5)-2$$

When an operand like 5 has operators to its left and right, CONVENTIONS are needed for deciding which operator applies to that operand (5).

- \hookrightarrow We say that + operator ASSOCIATES to the left

\hookrightarrow an operand(5) w/ plus on both sides of it belongs to the operator (+) to its left.

$$(9+5)+2$$

arithmetic operators

$+, -, \times, \div$

: left associative

assignment operator

$$a = (b) = c$$



arithmetic operators

$+, -, \times, \div$: left associative

assignment operator

$=$: right associative, that is, $a=b=c$ is treated as $a=(b=c)$

$a=(b=c)$

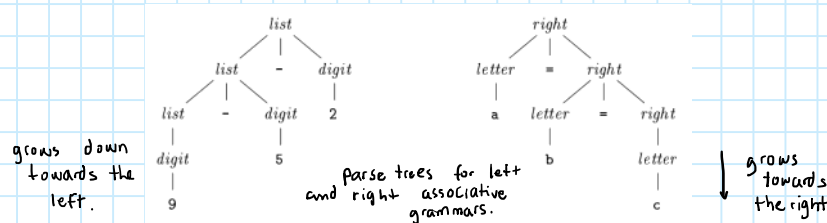


→ Strings like $a=b=c$ with a right associative operator are generated by the grammar

right \rightarrow letter = right letter

letter \rightarrow a | b | ... | z

→ The contrast of the parse tree for a left associative operator and a right associative operator $(=)$ is below.



2.2.6 Precedence of Operators

Consider the expression $9+5*2$

→ Two possible interpretations: $(9+5)*2$ or $9+(5*2)$

→ THE ASSOCIATIVITY RULES FOR $+$ AND $*$ (left association) APPLY TO OCCURRENCES OF THE SAME OPERATOR.



Rules defining the relative precedence of operators are needed when more than one operator.

→ We say that $*$ has HIGHER PRECEDENCE than $+$

- IF $*$ takes its operands $(*)$ before $+$
- Thus, 5 is taken by $*$ in both $9+5*2$ and $9*5+2$

→ which means $9+(5*2)$ and $(9*5)+2$

ex: • a grammar for arithmetic expressions can be constructed (2.6) from a table showing the associativity and precedence of operators

- We start with the four common arithmetic operators and a precedence table, showing the operators in order of increasing precedence.
- Operators on the same line have the same associativity and precedence

left associative: $+$ $-$ → expr
left associative: $*$ $/$ → term

We create two nonterminals *expr* and *term* for the two levels of precedence, and an extra nonterminal *factor* for generating basic units in expressions. The basic units in expressions are presently digits and parenthesized expressions.

$$\text{factor} \rightarrow \text{digit} \mid (\text{expr})$$

Now consider the binary operators, * and /, that have the highest precedence. Since these operators associate to the left, the productions are similar to those for lists that associate to the left.

$$\begin{array}{l} \text{term} \rightarrow \text{term} * \text{factor} \\ \quad \mid \text{term} / \text{factor} \\ \quad \mid \text{factor} \end{array}$$

} highest precedence
= grouped first

∞ ←

Similarly, *expr* generates lists of terms separated by the additive operators.

$$\begin{array}{l} \text{expr} \rightarrow \text{expr} + \text{term} \\ \quad \mid \text{expr} - \text{term} \\ \quad \mid \text{term} \end{array}$$

} lowest precedence
= grouped last

The resulting grammar is therefore

$$\begin{array}{l} \text{expr} \rightarrow \text{expr} + \text{term} \mid \text{expr} - \text{term} \mid \text{term} \\ \text{term} \rightarrow \text{term} * \text{factor} \mid \text{term} / \text{factor} \mid \text{factor} \\ \text{factor} \rightarrow \text{digit} \mid (\text{expr}) \end{array}$$

↗ higher level

← higher level

expr: list of terms separated by +/

↓
any parenthesized expr is factor

=

with parentheses we can develop expressions that have deep nesting (deep trees) arbitrarily.

() protect

A factor is something that cannot be torn apart: digit | (expr)

n=2
3 NONT:
expr,
term,
factor.

A term (that is not also a factor) is an expression that can be torn apart by operators of the highest precedence: * and / but not by the lower-precedence operators. An expression (that is not a term or factor) can be torn apart by any operator.

We can generalize this idea to any number n of precedence levels. We need n + 1 nonterminals. The first, like *factor* in Example 2.6, can never be torn apart. Typically, the production bodies for this nonterminal are only single operands and parenthesized expressions. Then, for each precedence level, there is one nonterminal representing expressions that can be torn apart only by operators at that level or higher. Typically, the productions for this nonterminal have bodies representing uses of the operators at that level, plus one body that is just the nonterminal for the next higher level.