

Programming Languages Recitation Grammars, Ada walk through

Deepti Verma

Computer Science Department

New York University

dv697@cs.nyu.edu

Introduction

- Regular Expressions
- Context Free Grammar (CFG)
- Derivations and Parse trees
- Ada installation

Phases of compiler

- Lexer
- Parser
- Semantic Analyzer
- Intermediate code generator
- Optimization(based on architecture of the system)
- Target code generation

Regular Expressions

- Tokens are the basic building blocks of a program. They are the shortest strings of characters with individual meaning.
- Examples include keywords, identifiers, symbols, constants and numbers.
- In order to specify tokens we use the notation of regular expressions
- Used in the Phase 1 -lexical analysis (Scanner) of the compiler.



Regular Expressions

- Given regular expressions R1 and R2 the following operations can be performed on them:
 - Concatenation: Two regular expressions next to each other. Eg. R1 R2
 - Alternation: Two regular expressions separated by a vertical bar, meaning any string generated by the first one or any string generated by second one. Eg. R1 | R2 (OR operation)
 - Kleene Star (*)

Examples

- `a` - matches the character 'a'
- `ε` - matches a null string
- `a|b|c` - matches 'a' or 'b' or 'c'
- `abc` - matches a concatenated with b concatenated with c
- `(a-z)` - matches any character between 'a' through 'z' (Shorthand)
- Alphabet (Uppercase and Lowercase) \longrightarrow `(A-Za-z)`
- Digit \longrightarrow `0|1|2|3|4|5|6|7|8|9`
- integer \longrightarrow `digit digit*`
- number \longrightarrow `integer|real`
- Identifier \longrightarrow `Alphabet (Digit|Alphabet)*`

Regular Expressions

- Drawbacks:
 - Nesting cannot be expressed in regular expressions which is central to programming languages.
 - For example: Nested parenthesis, palindromes

Context Free Grammar

- More powerful than regular languages/expressions.
- Any set of strings can be defined if we add RECURSION.
- Recognized by parsers.
- Every regular grammar is context free but not every context free grammar is regular.
- Used in the Phase 2 - Syntactic analysis (Parser) of the compiler.

Context Free Grammar

- Consists of
 - Productions (Substitution Rules) : Rules in a CFG of the form
$$A \longrightarrow B$$
 - Nonterminals: Symbols on the left side of the production (A).
 - Terminals: Symbols that make up the strings derived from grammar. They cannot appear on the left hand side of any production. They represent language's tokens. In the production shown above B is a set of terminals and nonterminals.
 - Start symbol: One of the nonterminals is designated as the one on the left side of the first production.
- The notation of CFG is sometimes called Backus-Naur form

Example 1:

- $A \longrightarrow 0A1 \mid \epsilon$
 - Language where it produces equal number of 0's followed by equal number of 1's.
 - Why cannot this be represented by regular expression?

Example 2:

- Give CFG for strings containing only 0's and 1's and contain equal number of 0's and 1's but in any order. Eq: 01, 0011, 0110, 1010, 1100, 11100010..
- Solution: $S \longrightarrow 0S1S \mid 1S0S \mid \epsilon$

Example 3:

- Give context free grammar for all the strings ending with character 'a' and containing the set of characters {a,b}. For example: aa, ba, aaba, abbaa..
- Solution:
 - $S \longrightarrow aS \mid bS \mid a$

Example 4:

- Give context free grammar for all the strings containing the set of characters $\{a,b\}$ and has even number of a's. For example: aa, aba, abba, abab, ababaa, bb
- Solution:
 - $S \longrightarrow aSaS \mid bS \mid \epsilon$

Example 5:

- Give context free grammar for the following code:

read X

read Y

prod=X*Y

write prod

- Solution:

Program \longrightarrow Stmt_List

Stmt_List \longrightarrow Stmt Stmt_List $\mid \epsilon$

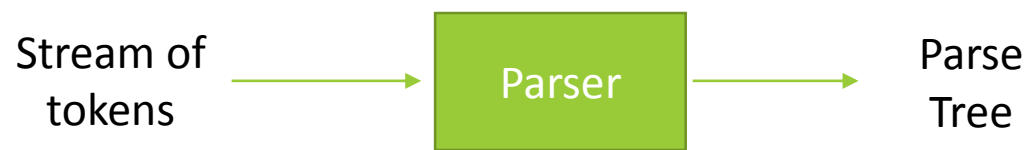
Stmt \longrightarrow id=Expr \mid read id \mid write id

Expr \longrightarrow Expr OP Expr \mid -Expr \mid (Expr) \mid id \mid num

OP \longrightarrow + \mid - \mid * \mid /

Parse Tree

- A CFG shows how to generate syntactically valid string of terminals.
 - Begin with start symbol.
 - Choose a production with the symbol on the left hand side; replace the nonterminal with the right hand side of that production. Start with the production with start symbol on the left.
 - Repeat this process until no nonterminals remain.
- Derivation: A series of replacement operations that shows how to derive a string of terminals from the start symbol.
- Derivation can be represented graphically as a *parse tree*.
- The root of the parse tree is the start symbol of the grammar.
- Based on CFG the parser in compiler produces a parse tree out of the stream of tokens.



Example 1:

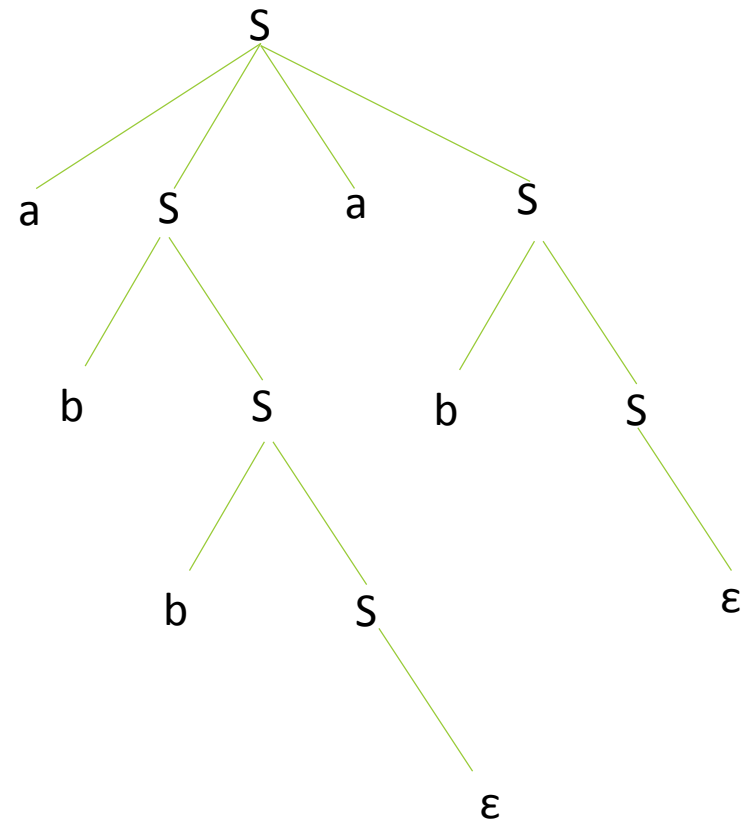
Given CFG for all the strings containing the set of characters $\{a,b\}$ and has even number of a's.

□ $S \longrightarrow aSaS \mid bS \mid \epsilon$

Derivation:

$S \longrightarrow aSaS \longrightarrow abSaS \longrightarrow abbSaS$
 $\longrightarrow abbaS \longrightarrow abbabS \longrightarrow abbab$

Parse Tree for string “abbab”:



Example 2:

Assume that Expression E has CFG:

$E \longrightarrow E + E \mid E * E \mid (E) \mid \text{num} \mid \text{id}$

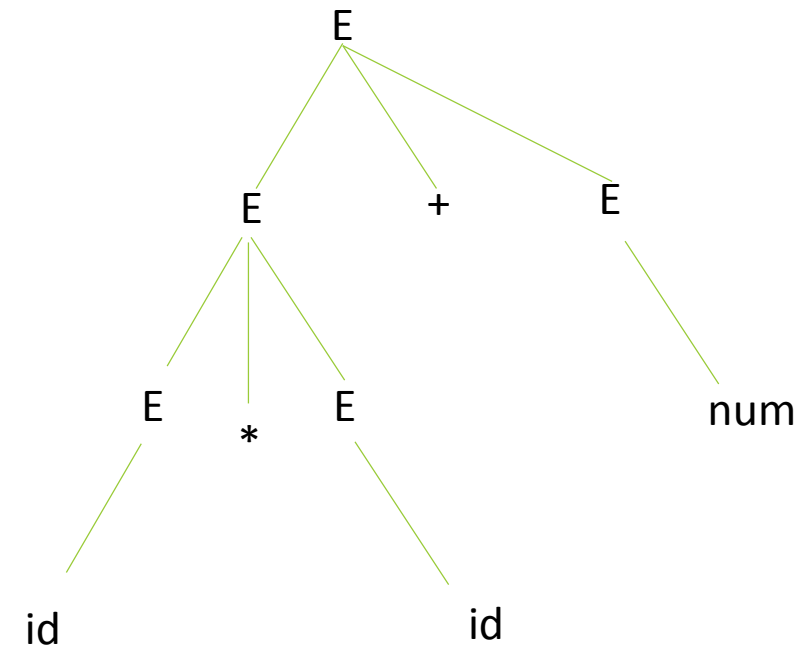
$\text{Id} \longrightarrow (a-b)(a-b \mid 0-9)^*$

$\text{num} \longrightarrow (0-9)^+$

Derivation:

$E \longrightarrow E + E \longrightarrow E + \text{num} \longrightarrow E * E + \text{num}$
 $\longrightarrow \text{id} * E + \text{num} \longrightarrow \text{id} * \text{id} + \text{num}$

Parse Tree for string "id*id+num":



Ada Installation and Example

- Follow link on NYU Classes for installation (Resources -> Ada Resources)
- Compile the examples