



**Compiler Design**

## Syntax Analysis

Amey Karkare  
 Department of Computer Science and Engineering  
 IIT Kanpur  
 karkare@iitk.ac.in

### Drawbacks of Regular Expressions:

\* nested expressions can't be matched

\* whenever we have to count things.

\* dynamic counting isn't possible.

### Limitations of regular languages

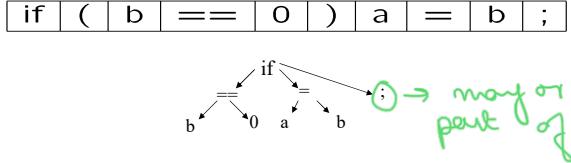
- How to describe language syntax precisely and conveniently. Can regular expressions be used?
- Many languages are not regular, for example, string of balanced parentheses
  - $((\dots)))$
  - $\{(\cdot)^i \mid i \geq 0\}$
  - There is no regular expression for this language
- A finite automata may repeat states, however, it cannot remember the number of times it has been to a particular state
- A more powerful language is needed to describe a valid string of tokens

3

## Syntax Analysis

- Check syntax and construct abstract syntax tree

if | ( | b | == | 0 | ) | a | = | b | ;



- ✓ Error reporting and recovery
- Model using context free grammars
- Recognize using Push down automata/Table Driven Parsers (for compilers)

2

## Syntax definition

- Context free grammars  $\langle T, N, P, S \rangle$ 
  - T: a set of **tokens** (terminal symbols)
  - N: a set of **non terminal** symbols
  - P: a set of **productions** of the form **nonterminal  $\rightarrow$  String of terminals & non terminals**
  - S: a **start symbol**  $\in N$
- A grammar derives strings by **beginning with a start symbol** and repeatedly **replacing a non terminal by the right hand side** of a production for that non terminal.
- The strings that can be derived from the start symbol of a grammar G form the language  $L(G)$  defined by the grammar.

4

## Examples

- String of balanced parentheses  $S \rightarrow (S)S | \epsilon$  *Not regular (pumping lemma)*

- Grammar

$$\begin{aligned} \text{list} &\rightarrow \text{list} + \text{digit} \\ &| \text{list} - \text{digit} \\ &| \text{digit} \\ \text{digit} &\rightarrow 0 | 1 | \dots | 9 \end{aligned}$$

Consists of the language which is a list of digit separated by + or -.

5

## Derivation

$$\begin{aligned} \text{list} &\rightarrow \underline{\text{list}} + \text{digit} \\ &\rightarrow \underline{\text{list}} - \text{digit} + \text{digit} \\ &\rightarrow \underline{\text{digit}} - \text{digit} + \text{digit} \\ &\rightarrow 9 - \underline{\text{digit}} + \text{digit} \\ &\rightarrow 9 - 5 + \underline{\text{digit}} \\ &\rightarrow 9 - 5 + 2 \end{aligned}$$

Does  $9-5+2 \in \text{Grammar}$   
Proof is sufficient

Therefore, the string  $9-5+2$  belongs to the language specified by the grammar

The name context free comes from the fact that use of a production  $X \rightarrow \dots$  does not depend on the context of  $X$  (we can choose any option of production)

\* There are context sensitive grammars as well  
 $\Rightarrow$  Chomsky hierarchy

## Examples ...

- Simplified Grammar for C block *not terminals*  
 $\begin{aligned} \text{block} &\rightarrow \{ \text{decls statements} \} \\ \text{statements} &\rightarrow \text{stmt-list} | \epsilon \\ \text{stmt-list} &\rightarrow \text{stmt-list stmt} ';' \\ &| \text{stmt} ';' \\ \text{decls} &\rightarrow \text{decls declaration} | \epsilon \\ \text{declaration} &\rightarrow \dots \text{int} | \text{float} | \dots \end{aligned}$

*we have to have all declarations.* <sup>7</sup>

## Syntax analyzers

- Testing for membership whether  $w$  belongs to  $L(G)$  is just a "yes" or "no" answer
- However the syntax analyzer
  - Must generate the parse tree
  - Handle errors gracefully if string is not in the language
- Form of the grammar is important
  - Many grammars generate the same language
  - Tools are sensitive to the grammar
    - We may have to keep restrictions bcoz of grammar* <sup>8</sup>

- + If a variable has been declared or not  
→ can't be checked.

$\alpha, \beta$ : string of terminals and non-terminals  
A: Non-terminal

Type compatibility

Declaration

Initialization

### What syntax analysis cannot do!

- To check whether variables are of types on which operations are allowed
- To check whether a variable has been declared before use  $wCw \rightarrow \text{Not CFL}$   
 $wCw\gamma \rightarrow \text{CFL}$
- To check whether a variable has been initialized
- These issues will be handled in semantic analysis

9

### Derivation

- If there is a production  $A \rightarrow \alpha$  then we say that A derives  $\alpha$  and is denoted by  $A \Rightarrow \alpha$
- $\alpha A \beta \Rightarrow \alpha \gamma \beta$  if  $A \rightarrow \gamma$  is a production
- If  $\alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_n$  then  $\alpha_1 \Rightarrow^+ \alpha_n$
- Given a grammar G and a string w of terminals in  $L(G)$ , we can write  $S \Rightarrow^+ w$
- If  $S \Rightarrow^* \alpha$  where  $\alpha$  is a string of terminals and non terminals of G then we say that  $\alpha$  is a sentential form of G

w is part of  
 $L(G)$   
(S derives w)

They are like sentences but not fully formed.

Not part  
of grammar  
intermediate

### Derivation ...

- If in a sentential form only the leftmost non terminal is replaced then it becomes leftmost derivation
- Every leftmost step can be written as  $wAy \Rightarrow^{lm*} w\delta y$  (A is the first non terminal)  
where w is a string of terminals and  $A \rightarrow \delta$  is a production
- Similarly, right most derivation can be defined
- An ambiguous grammar is one that produces more than one leftmost (rightmost) derivation of a sentence

11

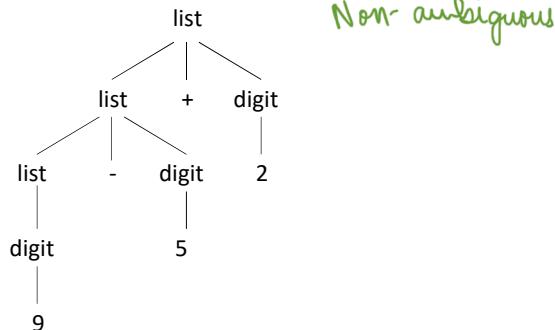
### Parse tree

- shows how the start symbol of a grammar derives a string in the language
- root is labeled by the start symbol
- leaf nodes are labeled by tokens
- Each internal node is labeled by a non terminal
- if A is the label of a node and  $x_1, x_2, \dots, x_n$  are labels of the children of that node then  $A \rightarrow x_1 x_2 \dots x_n$  is a production in the grammar

12

## Example

Parse tree for 9-5+2



13

## Ambiguity

- A Grammar can have more than one parse tree for a string

- Consider grammar

$\text{list} \rightarrow \text{list} + \text{list}$

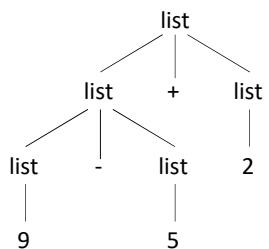
|  $\text{list} - \text{list}$

| 0 | 1 | ... | 9

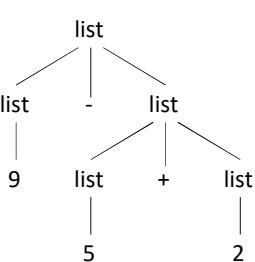
→ Ambiguous

- String 9-5+2 has two parse trees

14



15



## Ambiguity ...

- Ambiguity is problematic because meaning of the programs can be incorrect
- Ambiguity can be handled in several ways
  - Enforce associativity and precedence
  - Rewrite the grammar (cleanest way)
- There is no algorithm to convert automatically any ambiguous grammar to an unambiguous grammar accepting the same language
- Worse; there are inherently ambiguous languages!  $a^nb^nc^nd^n \cup S$

$G_1$

$G_2$

16

\* Python doesn't have dangling else problem bcoz of alignment rules.

Any else matches closest unmatched if.

## Ambiguity in Programming Lang.

- Dangling else problem

$\text{stmt} \rightarrow \text{if expr stmt}$

|  $\text{if expr stmt else stmt}$

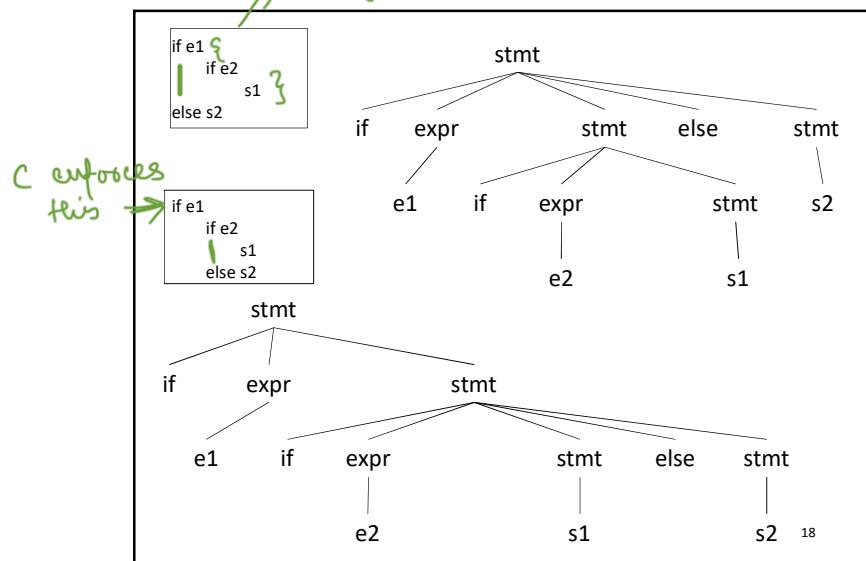
- For this grammar, the string

$\text{if e1 if e2 then s1 else s2}$

has two parse trees

17

C  $\Rightarrow$  Need of brackets



## Resolving dangling else problem

- General rule: match each **else** with the closest previous **unmatched if**. The grammar can be rewritten as

$\text{stmt} \rightarrow \text{matched-stmt}$

|  $\text{unmatched-stmt}$

$\text{matched-stmt} \rightarrow \text{if expr matched-stmt}$   
else matched-stmt

| others ??.

$\text{unmatched-stmt} \rightarrow \text{if expr stmt}$

| if expr matched-stmt  
else unmatched-stmt

} count of  
if  $\neq$  else

Parsing:

There will always be shift/reduce conflict because of if/else.

19

## Associativity

- If an operand has operator on both the sides, the side on which operator takes this operand is the associativity of that operator

In  $a+b+c$   $b$  is taken by left +

$+, -, *, /$  are left associative

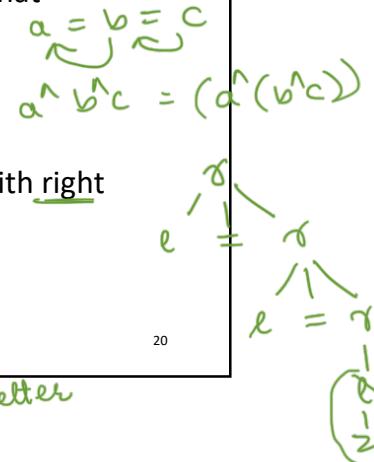
$^, =$  are right associative

- Grammar to generate strings with right associative operators

right  $\rightarrow$  letter = right | letter

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

- left-associative



20

left  $\rightarrow$  left = letter | letter

## Precedence

- String  $a+5*2$  has two possible interpretations because of two different parse trees corresponding to  $(a+5)*2$  and  $a+(5*2)$
- Precedence determines the correct interpretation.
- Next, an example of how precedence rules are encoded in a grammar

21

## Precedence/Associativity in the Grammar for Arithmetic Expressions

Ambiguous

$$\begin{array}{l} E \rightarrow E + E \\ | \\ E * E \\ | \\ (E) \\ | \\ \text{num} \mid \text{id} \end{array}$$

$$\begin{array}{l} 3 + 2 + 5 \\ 3 + 2 * 5 \end{array}$$

- Unambiguous, with precedence and associativity rules honored

$$\begin{array}{l} E \rightarrow E + T \mid T \\ T \rightarrow T * F \mid F \end{array}$$

$$\begin{array}{l} F \rightarrow (E) \mid \text{num} \\ | \\ \text{id} \end{array}$$

22

left recursive

{ highest precedence }

Top-down: Grammars can't be left recursive  
Something like left-factoring needs to be done.

## Parsing

- Process of determination whether a string can be generated by a grammar
- Parsing falls in two categories:

Easier to understand

Top-down parsing: [Brute force]

Construction of the parse tree starts at the root (from the start symbol) and proceeds towards leaves (tokens or terminals)

Bottom-up parsing: → (constructive)

More powerful

Construction of the parse tree starts from the leaf nodes (tokens or terminals of the grammar) and proceeds towards root (start symbol)

23

ANTLR: Top-down parser.

\* For same amount of look-ahead bottom-up parsers are much more powerful.

