

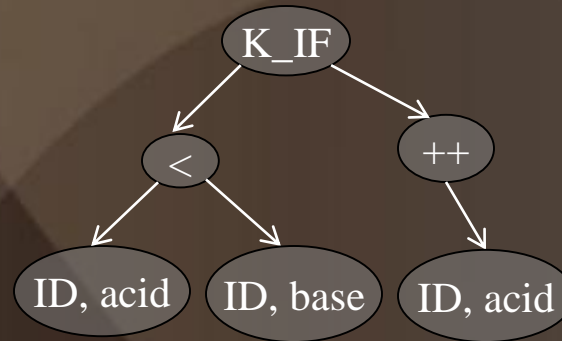
# *Lecture #7*

## Syntax Analysis - I

# Syntax Analysis

1. `if ( acid < base ) ++acid;`

- **Parser Input:** K\_IF, (, <ID, acid>, <, <ID, base>, ), ++, <ID, acid>, DELIM
- **Parser Output:**



- After lexical analysis, we have a series of tokens.
- In **syntax analysis (or parsing)**, we want to interpret what those tokens mean.
- Goal: Obtain the *structure* described by the series of tokens  
Report *errors* if the tokens do not properly encode a structure.

# Formal Grammar

A grammar,  $G$ , is a 4-tuple  $G=\{S,N,T,P\}$ , where:

$S$  is a starting symbol;  $N$  is a set of non-terminal symbols;

$T$  is a set of terminal symbols;  $P$  is a set of production rules.

Example:

LAUGH $\rightarrow$ LAUGH hah	rule 1
/ hah	rule 2

We can use this grammar to create sentences: E.g.:

<u>Rule</u>	<u>Sentential Form</u>
-	LAUGH
1	LAUGH hah
2	hah hah

Such a sequence of rewrites is called a *derivation*

# *Derivations*

$$\alpha A \beta \rightarrow \alpha \gamma \beta \quad \text{if} \quad A \rightarrow \gamma$$

$$\left\{ \begin{array}{l} \alpha \rightarrow \alpha \\ \alpha \xrightarrow{*} \beta \quad \text{and} \quad \beta \rightarrow \gamma \quad \text{then} \quad \alpha \xrightarrow{*} \gamma \end{array} \right.$$

$$S \xrightarrow{*} \alpha \quad \left\{ \begin{array}{l} \alpha \text{ is a sentential form} \\ \alpha \text{ is a sentence if it contains} \\ \text{only terminal symbols} \end{array} \right.$$

*The process of discovering a derivation for some sentence is called parsing!*

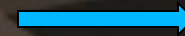
# *Regular Grammars and DFAs*

- We need a *language* for describing valid strings of tokens and a *method* for distinguishing valid from invalid strings of tokens.
- Regular Grammar - Productions have the form:  $A \rightarrow xB$  or  $A \rightarrow \epsilon$
- Regular Grammar to DFA:
  - Terminal symbols – *Input alphabet* for the DFA
  - Non-terminal symbols – Represent *States*
  - If a production has the form:  $A \rightarrow \epsilon$ , then ‘A’ is an accepting state
  - A production of the form  $A \rightarrow xB$  denotes a transition from state ‘A’ to state ‘B’ on input symbol ‘x’.

# Syntax Analysis using DFAs?

- Can we do syntax analysis with DFAs?

We could design a grammar for expressions of the form '*num + num*' using a DFA



$E \rightarrow \text{num } A$   
 $A \rightarrow + B$   
 $B \rightarrow \text{num } C$   
 $C \rightarrow \varepsilon$

- Programming languages have recursive structures
  - For eg:  $EXPR \rightarrow \text{if } EXPR \text{ then } EXPR \text{ else } EXPR / OTHER$
- DFAs cannot count – cannot handle such programming structures
  - Example: Regular grammars cannot define: *Expressions with properly balanced parentheses*  $\{( ^i )^i \mid i \geq 0\}$   $[S \rightarrow ( S ) \mid \varepsilon]$
  - Similarly: *Functions with properly nested block structure*

# *Push Down Automata*

- The situation can be handled if the DFA is augmented with *memory*
  - Memory implemented as a *stack*
  - Such an automata is called *Push Down Automata (PDA)*
- State table for a PDA that recognizes nested parentheses:

		Input Symbol		
		(	)	EOF
State	0	Push 1	Error	Accept
	1	Push 1	Pop	Error

- Determine the stack values and actions at each step as the following string is parsed:  $(( ) ( ( ) ))$

# *Context Free Grammars*

- PDAs have the power to accept Context Free Grammars (CFGs)
- Formally a CFG  $G = (T, N, S, P)$ , where:
  - $T$  is the set of terminal symbols in the grammar (i.e., the set of tokens returned by the lexical analyzer)
  - $N$ , the non-terminals, are variables that denote sets of (sub)strings occurring in the language. These impose a structure on the grammar.
  - $S$  is the goal symbol, a distinguished non-terminal in  $N$  denoting the entire set of strings in  $L(G)$ .
  - $P$  is a finite set of productions specifying how terminals and non-terminals can be combined to form strings in the language. Each production must have a single non-terminal on its left hand side.



# Context Free Grammars

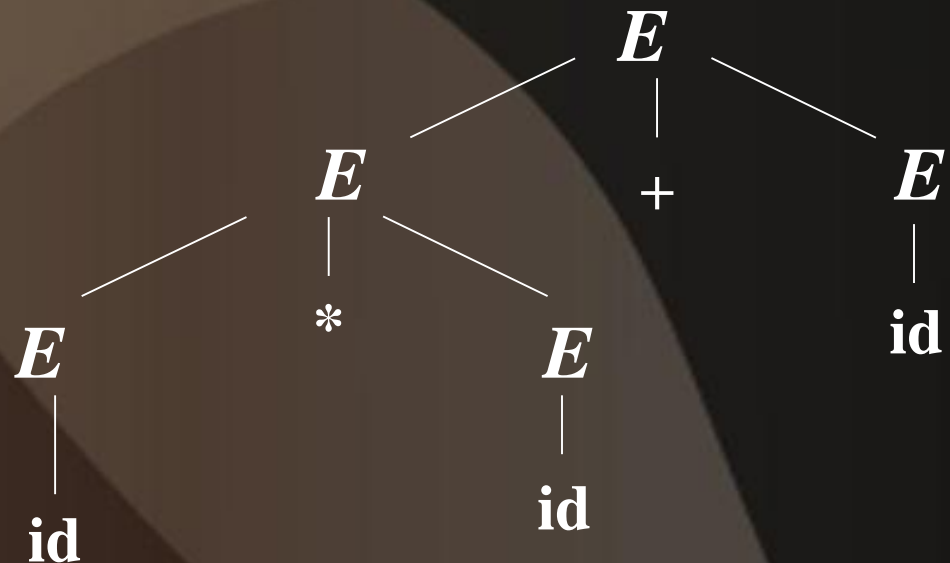
- A possible CFG for arithmetic operations:

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

- Input string:  $\text{id} * \text{id} + \text{id}$

$E \rightarrow E + E$   
 $\rightarrow E * E + E$   
 $\rightarrow \text{id} * E + E$   
 $\rightarrow \text{id} * \text{id} + E$   
 $\rightarrow \text{id} * \text{id} + \text{id}$

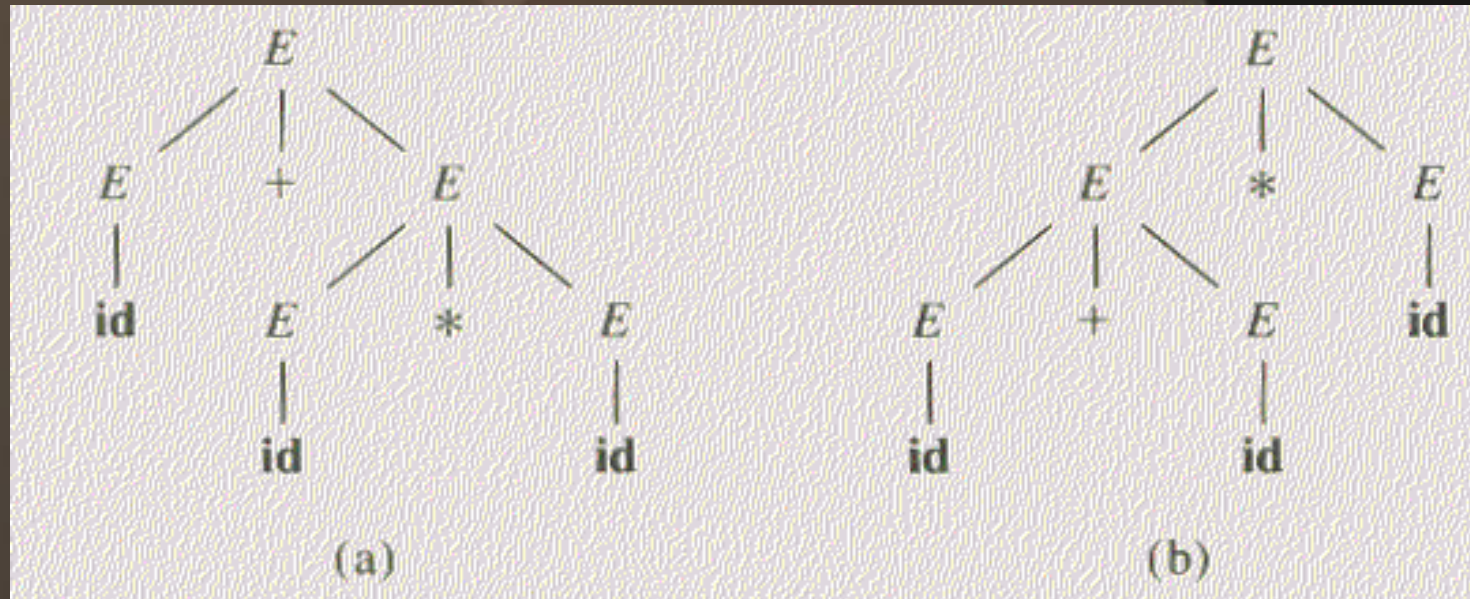
*Left-Most Derivation*



- Inorder traversal of the leaves give the original input string
- All derivations of a string should yield the same parse tree

# Ambiguity

- Derivation defines a parse tree
- A grammar is *ambiguous* if more than one right-most or left-most derivations may be obtained for some string / sentence
- Equivalently, a grammar is ambiguous if more than one parse tree may be obtained for some string



# *Eliminating Ambiguity*

- Ambiguity is bad
  - Leaves meaning of some programs ill-defined
  - Leaves up to the compiler which parse tree to accept
- Several ways to handle ambiguity
- Layering (most direct)

$$\begin{array}{lcl} E \rightarrow E + E \mid E * E \mid (E) \mid \text{id} & \longrightarrow & \begin{array}{l} E \rightarrow E' + E \mid E' \\ E' \rightarrow \text{id} * E' \mid \text{id} \mid (E) * E' \mid (E) \end{array} \end{array}$$

- Enforces precedence of ‘\*’ over ‘+’
- E controls ‘+’ and E’ controls ‘\*’
- All the ‘+’s will be handled before any of the ‘\*’s
- ‘\*’s will always be nested more deeply inside the parse tree than the ‘+’s

# *Eliminating Ambiguity*

- Another Expression

**$E \rightarrow \text{if } E \text{ then } E \mid \text{if } E \text{ then } E \text{ else } E \mid \text{OTHER}$**

- The expression:  **$\text{if } E_1 \text{ then if } E_2 \text{ then } E_3 \text{ else } E_4$**  has two separate parse trees
- The property that we want is: *else matches the closest unmatched then*
- Can be resolved as:

**$E \rightarrow \text{MIF} \mid \text{UIF}$**

**$\text{MIF} \rightarrow \text{if } E \text{ then MIF else MIF} \mid \text{OTHER}$**

**$\text{UIF} \rightarrow \text{if } E \text{ then } E \mid \text{if } E \text{ then MIF else UIF}$**

# *Eliminating Left-Recursion*

- Direct Left-Recursion

$$A \rightarrow A\alpha \mid \beta$$



$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' \mid \varepsilon$$

$$A \rightarrow A\alpha_1 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \dots \mid \beta_n$$



$$A \rightarrow \beta_1 A' \mid \dots \mid \beta_n A'$$

$$A' \rightarrow \alpha_1 A' \mid \dots \mid \alpha_n A' \mid \varepsilon$$

# *Eliminating Left-Recursion*

- Indirect Left-Recursion

$$\begin{aligned} S &\rightarrow Aa \mid b \\ A &\rightarrow Ac \mid Sd \mid \varepsilon \end{aligned}$$

- Algorithm

Arrange the non-terminals in some order  $A_1, \dots, A_n$ .

for (i in 1..n) {

  for (j in 1..i-1) {

    replace each production of the form  $A_i \rightarrow A_j \gamma$  by the  
    productions  $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$  where

$$A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$$

  }

  eliminate the immediate left recursion among  $A_i$  productions

}

# *Left Factoring*

$$A \rightarrow \alpha\beta_1 \mid \dots \mid \alpha\beta_n \mid \gamma$$



$$A \rightarrow \alpha A' \mid \gamma$$

$$A' \rightarrow \beta_1 \mid \dots \mid \beta_n$$

# *Next Lecture*

## Top-Down Parsing