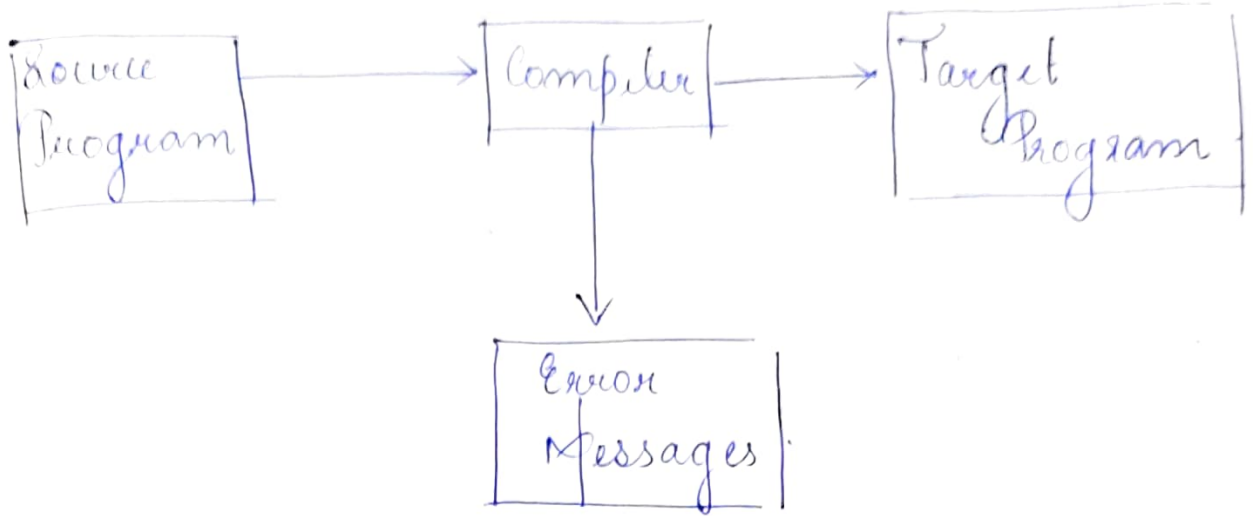


# Introduction to Compilers



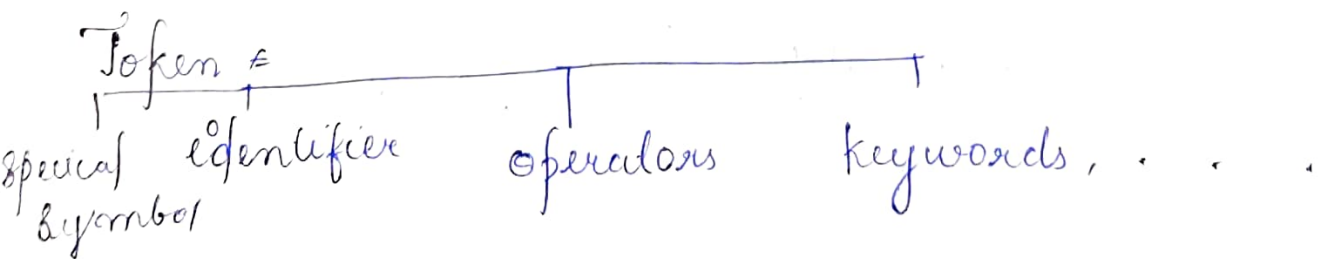
## Phases of a Compiler

Analysis of the source program

- ① Lexical Analysis (or scanning or linear)
- ② Hierarchical Analysis (Syntax/Parsing)
- ③ Semantic Analysis.

## Linear Analysis

$sum = a + b;$  (letter followed by (letter/digit))



The lexical analysis phase the given input source program is scanned from left to right thereby group into tokens.

### Token

→ sequence of characters having collective meaning and moreover it must match with a pattern (regular expressions)

- identifiers
- keyword
- operator
- special symbol
- constant

### Lexeme

→ sequence of characters in the source program is matched by the pattern in order to

identify a specific token eg:  $\text{Position} \equiv \text{initial} + \text{rate} * 60$

eg:  $\text{sum} = a + b;$

$$\text{Id1} = \text{Id2} + \text{Id3} * 60$$

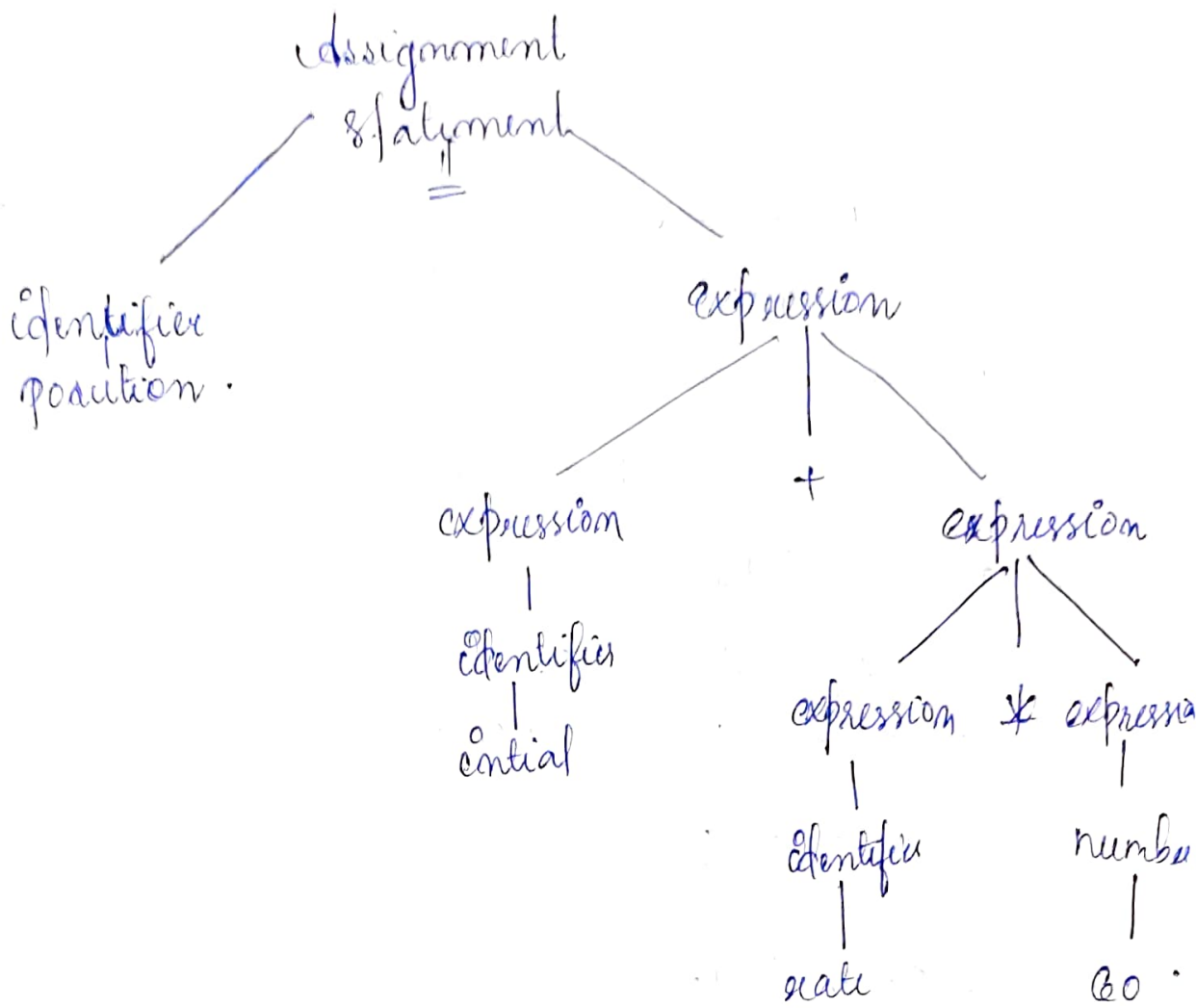
$\text{printf} ("enter the numbers");$

### Hierarchical Analysis (Parsing / Syntax)

#### CFG Rule

$$E \rightarrow E + E \mid E * E \mid (E) \mid id \mid num.$$

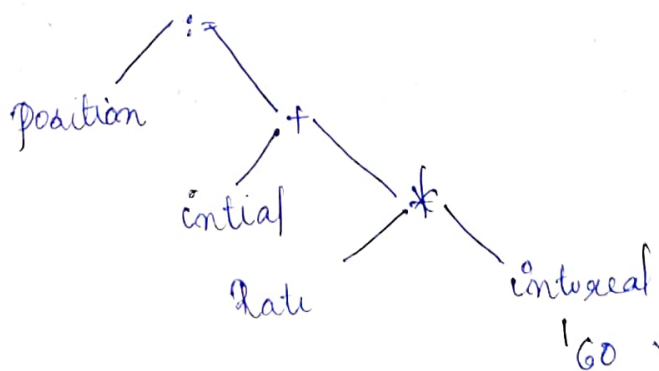
$$\text{Position} = \text{initial} + \text{rate} * 60$$



→ It involves grouping the tokens of source program (received during lexical analysis phase) into grammatical phases that can be used by compiler to synthesize output (parse tree based on CFG (based on the given source program))

### Semantic analysis

eg: position = initial + rate \* 60



## Intermediate Code Generation

- it should be easy to produce
- easy to translate into target program.
- three address code (representation)
- position = initial + rate \* 60

$$id1 = id2 + id3 * 60$$

$$temp1 = \text{int to real}(60)$$

$$temp2 = id3 * temp1$$

$$temp3 = id2 + temp2$$

$$id1 = temp3$$

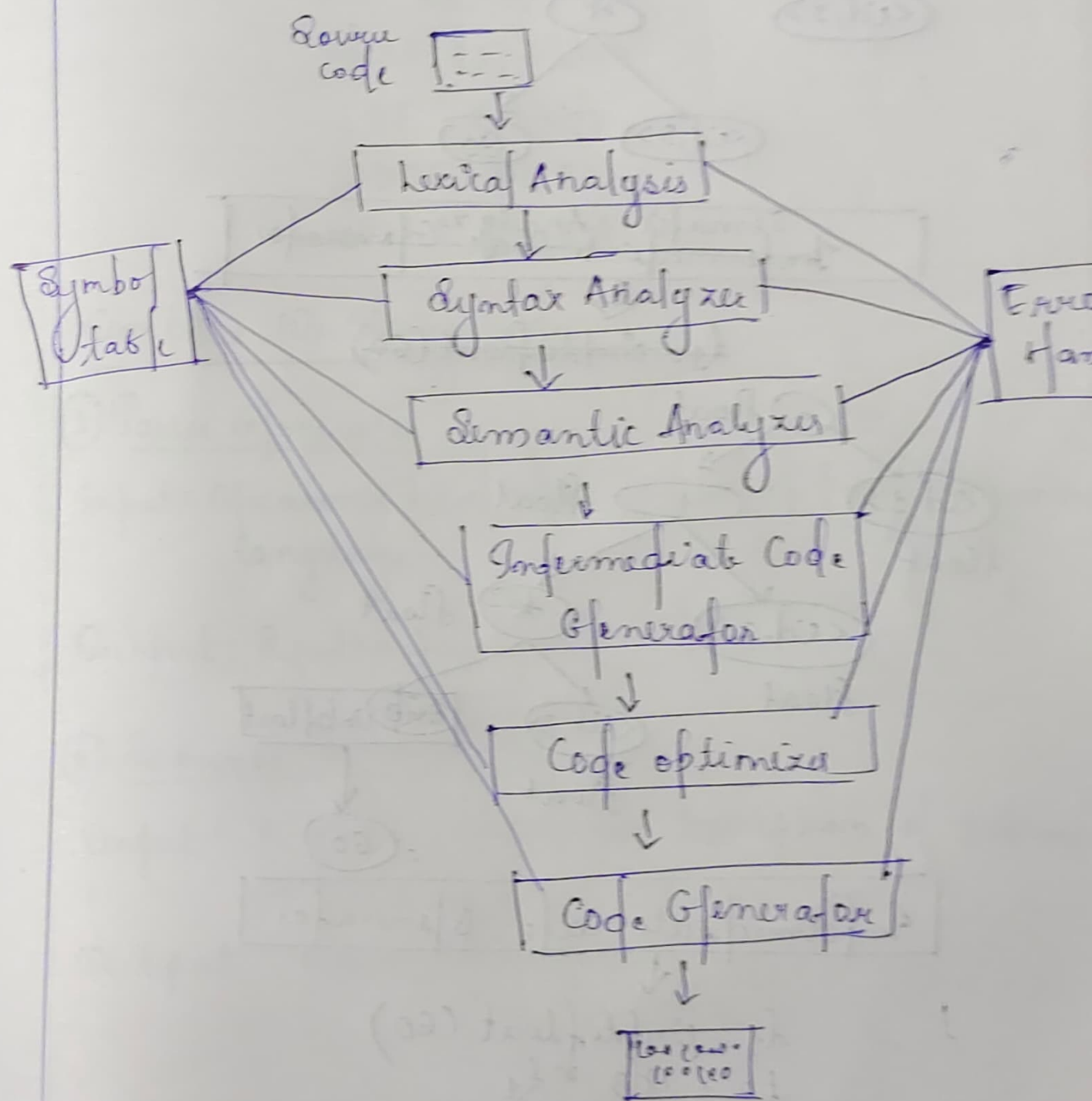
## Code Optimisation

$$Temp1 = id3 * \text{int to real}(60)$$

$$id1 = id2 + temp1$$

## Code Generation

- Final phase of the compiler is the generation of target code consisting of relocatable / machine code.
- Memory locations are selected for each of the variables used by the program.

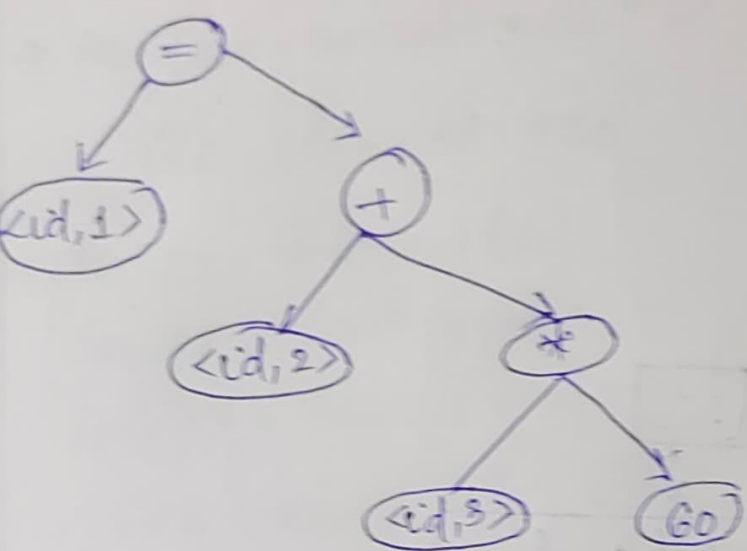




Lexical Analyzer

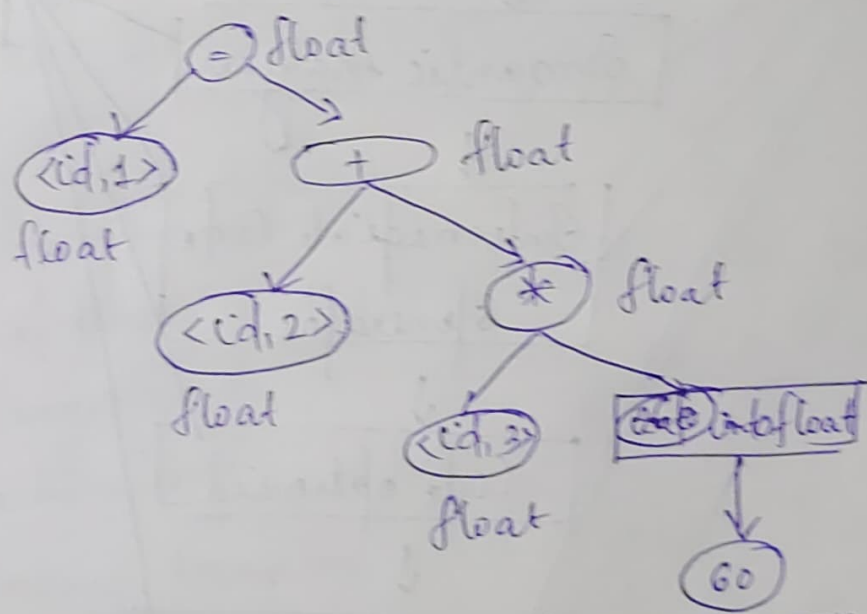
$\langle id, 1 \rangle \langle = \rangle \langle id, 2 \rangle \langle + \rangle \langle id, 3 \rangle \langle * \rangle \langle 60 \rangle$

Syntax Analyzer



~~Semantic Analyzer / Intermediate Code Generator~~

~~$t_1 = \text{inttofloat}(60)$~~



Intermediate Code Generator

$t_1 = \text{inttofloat}(60)$   
 $t_2 = id3 * t_1$   
 $t_3 = id1 + t_2$

$cd1 = t3$

Code Optimizer



$t1 = cd3 * 60.0$

$cd1 = cd2 + t1$

Code Generator



LDF R2, cd3

MULF R2, #60.0

LDF R1, cd2

ADDF R1, R2

STF cd1, R1

## Compiler Writing Tools

### ① Parser Generator

Input: Grammatical description of programming language

Output: Syntax Analyzer

### ② Scanner Generator

Input: Regular expression description of tokens of a language

Output: Lexical Analyzer

## Module 2

### Context Free Grammar

- set of tokens terminals
- non-terminal
- set of productions ( $A \rightarrow B$ )
  - $A \rightarrow \text{L.H.S.}$
  - $B \rightarrow \text{R.H.S.}$
- Start symbol.
- $S \rightarrow AB$ 
  - $A \rightarrow a \mid \epsilon$
  - $AB \rightarrow \text{Non-Terminals}$
  - $a \mid \epsilon \rightarrow \text{non-terminals}$

### Terminals

- $a, b, \dots$  (lowercase letters)
- operator symbol
- punctuation symbols
- id or if

### Non-terminals

- uppercase letters
- lower italic names (expression, statement, operator (op))
- 

→ if  $A \rightarrow \alpha_1 \mid \alpha_2, \dots, A \rightarrow \alpha_k$

we can write

$$A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_k.$$



## Derivation

→ sequence of production rules / sequence of replacements of non-terminals.

eg: derivation of  $(id * id)$  → generate derivation for  $(id * id)$ .

$$\left. \begin{array}{l} E \rightarrow E + E \\ \rightarrow E * E \\ \rightarrow (E) \\ \rightarrow -E \\ \rightarrow id \end{array} \right\} \text{ productions.}$$

∴  $E \rightarrow (E)$   
→  $(E * E)$   
→  $(id * E)$   
→  $(id * id)$

$S \xRightarrow{*} \alpha$  if  $\alpha$  contains non-terminals, it is called sentential form of  $G$ .

if  $\alpha$  <sup>does not</sup> contains non-terminals it is called sentence of  $G$ .

→ if we choose the leftmost non-terminal in each derivation step. This derivation called leftmost derivation. (used in top down)

→ if we always choose the rightmost ~~derivation~~ non-terminal that derivation rightmost derivation. (used in bottom up parse tree)

⇒ leftmost derivation.

generate derivation  $(id * id) + id$

$$\begin{array}{ll} E \rightarrow E + E \\ \rightarrow E * E & \rightarrow -E \\ \rightarrow (E) & \rightarrow id \end{array}$$

$$E \rightarrow E * E$$

$$E \rightarrow E + E$$

$$\rightarrow (E) + E$$

$$\rightarrow (E * E) + E$$

$$\rightarrow \cancel{(id * id) + id} \dots$$

$$\rightarrow (id * E) + E$$

$$\rightarrow (id * id) + E$$

$$\rightarrow (id * id) + id$$

Rightmost derivation

$$(id * id) + id.$$

$$E \rightarrow E$$

$$\rightarrow E + E$$

$$\rightarrow E + id$$

$$\rightarrow (E) + id$$

$$\rightarrow (E * E) + id$$

$$\rightarrow (E * id) + id$$

$$\rightarrow (id * id) + id.$$

Parse tree using <sup>left</sup>rightmost derivation

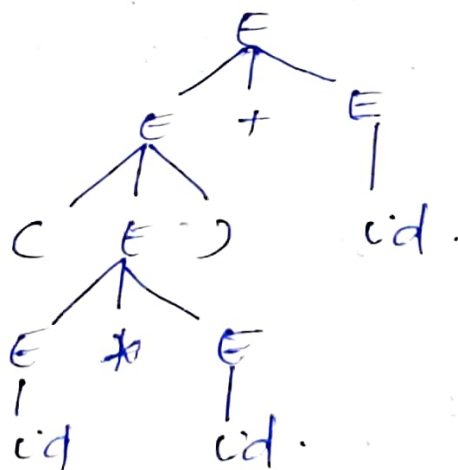
$$(id * id) + id$$

$$E \rightarrow E$$

$$\rightarrow E + E$$

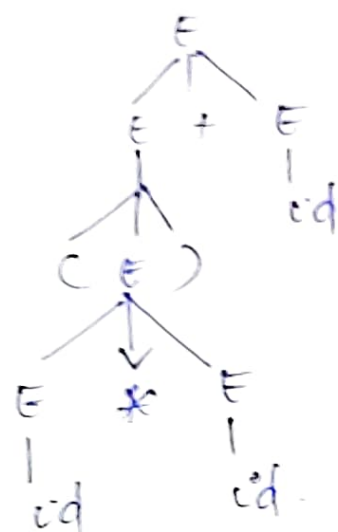
$$\rightarrow E + id$$

$$\rightarrow$$



using rightmost derivation.

$(id+id)*id$



### Ambiguous Grammar.

→ A grammar is ambiguous if there are multiple parse trees for the same sentence. Grammar for programming language is ambiguous.

→  ~~$(id+id)*id$~~

→  $id+id*id$

rightmost derivation

$$\begin{aligned} E &\rightarrow E \\ &\rightarrow E + E \\ &\rightarrow id + \cancel{E * E} \\ &\rightarrow id + E * E \\ &\rightarrow id + id * E \\ &\rightarrow id + id * id. \end{aligned}$$

→ leftmost derivation #2

$$\begin{aligned} E &\rightarrow E \\ &\rightarrow \cancel{E * E} \\ &\rightarrow E + E * E \\ &\rightarrow id + E * E \\ &\rightarrow id + id * E \\ &\rightarrow id + id * id. \end{aligned}$$

## Left Recursion

If we have production of form.

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

it could be replaced by non-recursive productions

$$\boxed{\begin{array}{l} A \rightarrow \beta A' \\ A' \rightarrow \alpha A' \mid \epsilon \end{array}}$$

Q.  $E \rightarrow TE'$

$$E' \rightarrow +TE' \mid \epsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \epsilon$$

$$F \rightarrow (E) \mid id$$

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

$$E \rightarrow E+T \mid T$$

$$A = E$$

$$\alpha = +T$$

$$\beta = T$$

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \epsilon$$

Q.  $E \rightarrow E+T \mid T$

$$T \rightarrow T*F \mid F$$

$$F \rightarrow (E) \mid id$$

$$T \rightarrow T*F \mid F$$

$$A = T$$

$$\alpha = *F$$

$$\beta = F$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \epsilon$$

$$F \rightarrow (E) \mid id$$

## Indirect Left Recursion

$$S \rightarrow Aa/b$$

$$A \rightarrow Ac | Sc | \epsilon$$

Remove Left Recursion.

## Left Factoring

Production of form:  $A \rightarrow \alpha B_1 | \alpha B_2$

Replaced by

$$\boxed{\begin{array}{l} A \rightarrow \alpha A' \\ A' \rightarrow \beta_1 | \beta_2 \end{array}}$$

Productions form:  $A \rightarrow \alpha \beta_1 | \alpha \beta_2 | \dots | \alpha \beta_n | \gamma$

Replaced by:

$$\boxed{\begin{array}{l} A \rightarrow \alpha A' | \gamma \\ A' \rightarrow \beta_1 | \beta_2 | \dots | \beta_n \end{array}}$$

Q.  $S \rightarrow cE \mid S | cE \mid ScS | a$

$$E \rightarrow b$$

General:  $A \rightarrow \alpha \beta_1 | \alpha \beta_2 | \dots | \gamma :$

$$A \rightarrow \alpha A' | \alpha$$

$$A' \rightarrow \beta_1 | \beta_2$$

Q  $A = S$

$$S \rightarrow cE \mid S$$

$$S \rightarrow \gamma$$



$S \rightarrow \text{if } ts S' / a$

$S' \rightarrow a \in \mid es$

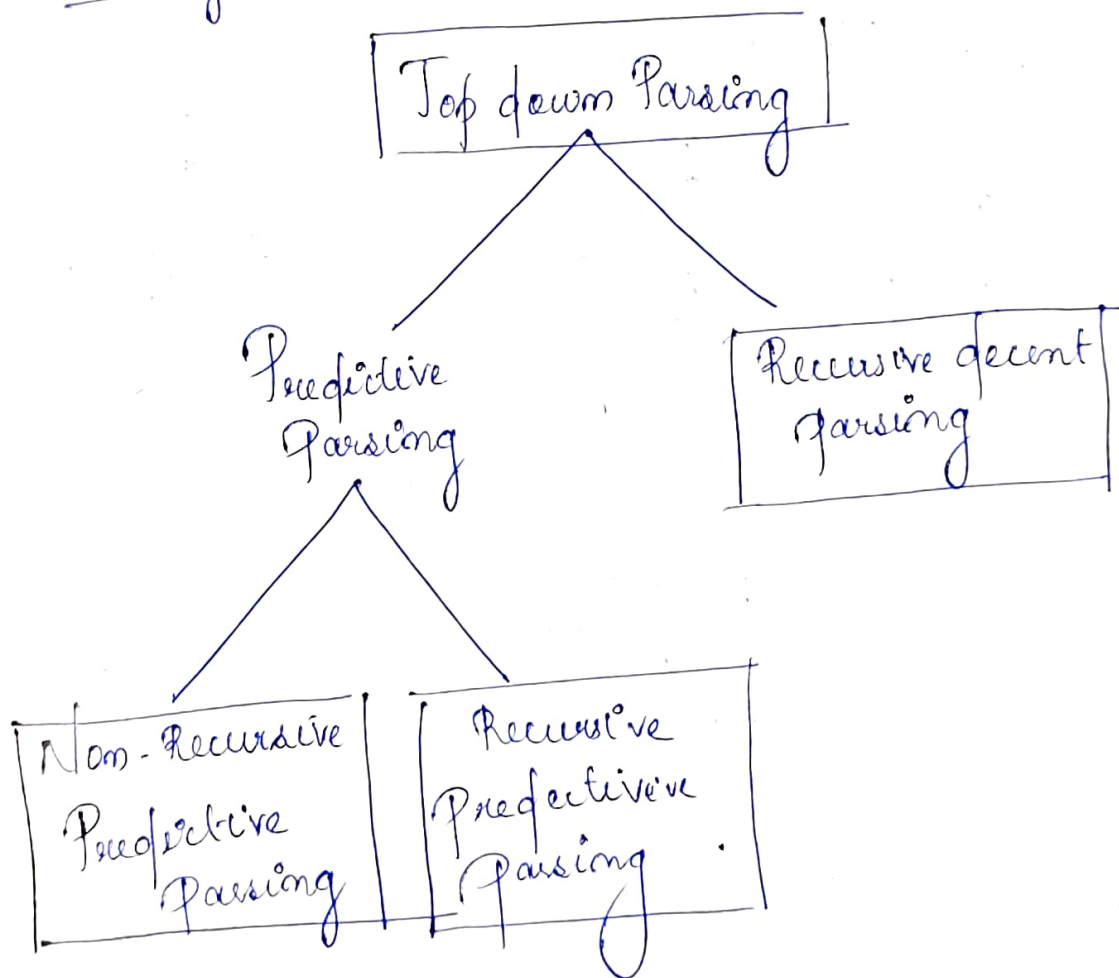
$E \rightarrow b$

→ left factoring useful for producing a grammar a su for predictive or top-down, parsing.

→ when choice between two alternative productions is not clear, we may be able to rewrite prodce to defer the decision until enough of input has been seen to make right choice.

T

## Parsing



## Recursive Descent Parser (RDP)

→ Top down parsing technique

→ Input scanning from left to right

→ Non-terminal: Implement Recursive Procedure

→ Terminal: Compare lookahead to check with input string

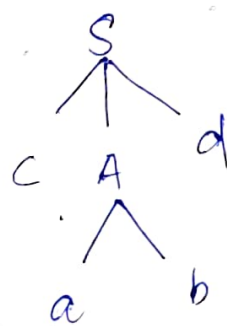
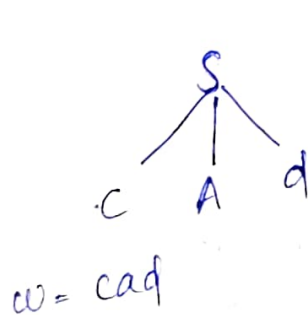
→ It need backtracking to identify the correct A-production

### Productions

$S \rightarrow cad$

$A \rightarrow ab|a$

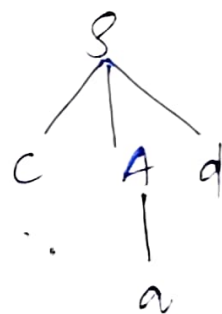
Input string:  $w = cad$



$w = cad$

b and d  
mismatch

so, backtrack



$w = cad$

Q. Production:  $E \rightarrow iE'$   
 $E' \rightarrow +iE' | e$

lookahead,  $l = \text{getchar}()$ ;

match (char l)

```
{
  if (l == l)
    l = getchar();
  else
    printf("error");
}
```

E()

```
{
  if (l == '0')
  {
    match(0);
    E'();
  }
}
```

E'()

```
{
  if (l == '+')
  {
    match(+);
    match(i);
    E'();
  }
}
```

else

```
{
  return;
}
```

main()

```
{
  E();
  if (l == '$')
    printf("Passed successfully");
}
```

## ④ Non-Recursive Predictive Parsing (Predictive Parsing)

→ Predictive parser is an efficient way of implementing recursive-descent parsing by handling the stack of activation records explicitly.

→ 2 types

① Recursive predictive Parsing.

② Non-Recursive predictive Parsing.

### Predictive Parsing Table

→ uses 2 functions.

① FIRST()

② FOLLOW()

→ ① FIRST()

① If  $x$  is a terminal then  $\text{FIRST}(x)$  is  $\{x\}$

② <sup>For</sup> eg:  $x \rightarrow a$   
 $\text{FIRST}(x) = \text{FIRST}(a) = \underline{\underline{\{a\}}}$

② If  $x \rightarrow \epsilon$  then add  $\epsilon$  to  $\text{FIRST}(x)$ .

③ If  $x$  is a non-terminal and  $x \rightarrow y_1 y_2 y_3 \dots y_n$ , then put  $a$  in  $\text{FIRST}(x)$  if for some  $i$ ,  $a$  is  $\text{FIRST}(y_i)$  and  $\epsilon$  is in all of  $\text{FIRST}(y_1) \dots \text{FIRST}(y_n)$ .

eg: 1)  $E \rightarrow TE'$

2)  $E' \rightarrow +TE' \mid \epsilon$

3)  $T \rightarrow FT'$

4)  $T' \rightarrow *FT' \mid \epsilon$

5)  $F \rightarrow (E) \mid id$

}

$F \rightarrow (E) \mid id \mid \epsilon$

$FIRST(F) = \{ (, id, \epsilon \}$

②  $E' \rightarrow +TE' \mid \epsilon$

$FIRST(E') = \{ +, \epsilon \}$

Term

Rule No: 1 of 2

③  $T' \rightarrow *FT' \mid \epsilon$

$FIRST(T') \rightarrow \{ *, \epsilon \}$

⑤  $F \rightarrow (E) \mid id$

$FIRST(F) \rightarrow \{ (, id \}$

$E \rightarrow$

①  $E \rightarrow TE'$

$FIRST(E) = FIRST(TE')$

$= FIRST(T)$

$FIRST(T) = FIRST(FT')$

$= FIRST(F)$

$= FIRST((E) \mid id)$

$= \{ (, id \}$

$FIRST(E) = \{ (, id \}$



$$\textcircled{8} \quad T \rightarrow TT'$$

$$\begin{aligned} \text{FIRST}(T) &= \text{FIRST}(T') \\ &= \underline{\underline{\{c, id\}}} \end{aligned}$$

## 2. FOLLOW()

### Rules

- ① If  $S$  is a start symbol then add  $\$$  to the FOLLOW( $S$ )
- ② If there is a production rule of form  $A \rightarrow \alpha B \beta$  then everything in  $\text{FIRST}(\beta)$  except for  $\epsilon$  is placed in  $\text{FOLLOW}(B)$ .
- ③ If there is a production  $A \rightarrow \alpha B$  or a production  $A \rightarrow \alpha B \beta$  where  $\text{FIRST}(\beta)$  contains  $\epsilon$  then everything in  $\text{FOLLOW}(A)$  in  $\text{FOLLOW}(B)$ .

Eg: 1)  $E \rightarrow TE'$

2)  $E' \rightarrow +TE' \mid \epsilon$

3)  $T \rightarrow FT'$

4)  $T' \rightarrow *FT' \mid \epsilon$

5)  $F \rightarrow (E) \mid id$

i)  $E \rightarrow TE'$

$F \rightarrow (E) \mid id$

$$\text{FOLLOW}(E) = \{ \$, ) \}$$

$$F \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \epsilon$$

$$\text{Follow}(E') = \{ \$, ) \}$$

$$\textcircled{3} T \rightarrow FT'$$

$$E \rightarrow TE'$$

Follow(

R.H.S  $\rightarrow$  variable ke first

$$\textcircled{3} T \rightarrow FT'$$

$$E \rightarrow TE'$$

$$\text{Follow}(T) = \{ +, \$, ) \}$$

$$E' \rightarrow +TE' \mid \epsilon$$

$$\text{Follow}(T) = \{ +, \$, ), \}$$

$$\textcircled{4} F' \rightarrow \text{Follow}(T') = \{ +, \$, ), \}$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \epsilon$$

$$\textcircled{5} \text{Follow}(F) = \{ *, +, \$, ) \}$$

$$T \rightarrow FT'$$

$$\textcircled{6} T \rightarrow F$$

$$T \rightarrow *FT' \mid \epsilon$$

$$\begin{aligned}
 Q. E &\rightarrow TE' \\
 E' &\rightarrow +TE' \mid \epsilon \\
 T &\rightarrow FT' \\
 T &\rightarrow *FT' \mid \epsilon \\
 F &\rightarrow (E) \mid id \mid \epsilon
 \end{aligned}$$

$$\text{FIRST}(T) \Rightarrow$$

$$\begin{aligned}
 \text{FIRST}(E) &\rightarrow \text{FIRST}(T) \\
 &\rightarrow \text{FIRST}(F) \\
 &\rightarrow \{ (, id, *, \epsilon \}
 \end{aligned}$$

$$\text{FIRST}(E') = \{ +, \epsilon \}$$

$$\begin{aligned}
 \text{FIRST}(T) &= \text{FIRST}(F) \\
 &= \{ (, id, *, \epsilon \}
 \end{aligned}$$

$$\text{FIRST}(T') = \{ *, \epsilon \}$$

$$\text{FIRST}(F) = \{ (, id, \epsilon \}$$

$$* \text{ FOLLOW}(E) = \{ \$, ) \}$$

$$* \text{ FOLLOW}(E') = \{ \$, ), \epsilon \}$$

$$* \text{ FOLLOW}(T) = \{ +, \$, ), \epsilon \}$$

$$* \text{ FOLLOW}(T') = \{ +, \$, ) \}$$

$$* \text{ FOLLOW}(F) = \{ +, \$, ), \epsilon \} \cup \{ *, +, \$, ) \}$$

⑤

# Pausing Table

E	id	+	*	(	)	\$
E'	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \xrightarrow{+} TE'$			$E' \rightarrow e$	$E' \rightarrow e$
T'	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'	$T \rightarrow e$	$T' \xrightarrow{e} FT'$	$T' \xrightarrow{*} FT'$		$T' \rightarrow e$	$T' \rightarrow e$
F	$F \rightarrow id$			$F \rightarrow (e)$		

Q.  $S \rightarrow iEtSs' | a$  ✓  
 $S' \rightarrow es | e$  ✓  
 $E \rightarrow b$

$FIRST(S) \rightarrow \{i, a\}$

$FIRST(S') \rightarrow \{e, e\}$

$FIRST(E) \rightarrow \{b\}$

$FOLLOW(S) \rightarrow \{e, \$, \}$

FOU

$FOLLOW(S') \rightarrow \{e, \$\}$

$FOLLOW(E) \rightarrow \{t\}$

## Parsing table

	i	a	e	b	t	$\epsilon$
S	$S \rightarrow iEtS'$	$S \rightarrow a$				
S'			$S' \rightarrow eS$ <del><math>S' \rightarrow \epsilon</math></del>			$S' \rightarrow \epsilon$
E				$E \rightarrow b$		

Since two productions come in one row so this grammar is not LL(1).

### LL(1)

- the first 1 stands for scanning the input from left to right
- the second 1 stands for producing leftmost derivation
- and 1 stands for using one input symbol of lookahead at each step to make parsing decision.

### Algorithm (for creating parsing table)

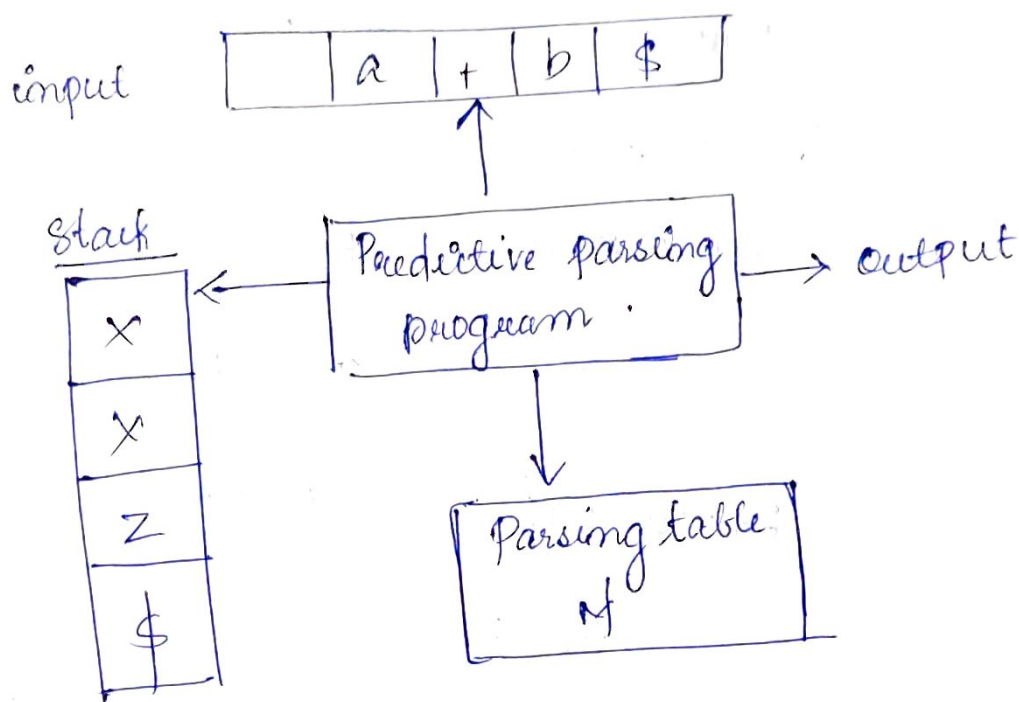
- ① For each production  $A \rightarrow \alpha$  do steps 2 & 3
- ② For each terminal  $a$  in  $\text{FIRST}(\alpha)$ , add  $A \rightarrow \alpha$  to  $M[A, a]$



③ if  $\epsilon$  is in  $FIRST(A)$ , add  $A \rightarrow \alpha$  to  $M[A, b]$  for each terminal  $b$  in  $FOLLOW(A)$ . If  $\epsilon$  is in  $FIRST(A)$ ; and  $\$$  is in  $FOLLOW(A)$ , add  $A \rightarrow \alpha$  to  $M[A, \$]$

④ make each of undefined entry of  $M$  be error.

Non-Recursive predictive parsing algorithm.



Algorithm

let  $ip$  point to the first symbol of  $w\$$ ;

repeat

let  $x$  be the top stack symbol and  $a$  the symbol pointed to

if  $x$  is a terminal or  $\$$  then

if  $x = a$  then

pop  $x$  from the stack and advance  $ip$

else error()

else /\*  $x$  is a nonterminal \*/

if  $M(x, a) = X \rightarrow y_1, y_2, \dots$  then begin

pop  $x$  from the stack;

push  $y_k, y_{k-1}, \dots, y_1$  onto the stack, with  $y_1$   
output productions  $X \rightarrow y_1, y_2, \dots, y_k$

end.

else error()

until  $X = \$/*$  stack is empty  $*/$

Q. Input:  $cd+cd*cd$

Q. Stack

Input

Output

$\$ E$

$cd+cd*cd\$$

$\$ E' T$

$cd+cd*cd\$$

$E \rightarrow TE'$

$\$ E' T' F$

$cd+cd*cd\$$

$T \rightarrow FT'$

$\$ E' T' cd$

$cd+cd*cd\$$

$F \rightarrow cd$

$\$ E' T'$

$cd+cd*cd\$$

$\$ E'$

$+cd*cd\$$

$T' \rightarrow E$

$\$ E' T' +$

$+cd*cd\$$

$E' \rightarrow +TE'$

$\$ E' T$

$cd*cd\$$

$\$ E' T' F$

$cd*cd\$$

$T \rightarrow FT'$

$\$ E' T' cd$

$cd*cd\$$

$F \rightarrow cd$

\$ E T |      \* id \$

\$ E T | F \*      \* id \$

\$ E T | F      id \$

\$ E T | id      id \$

\$ E T |      \$

\$ E ' |      \$

\$      \$

$T' \rightarrow * E T |$

$F \rightarrow id$

$T' \rightarrow \epsilon$

$E' \rightarrow \epsilon$