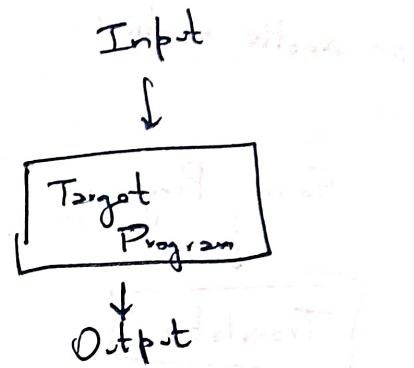
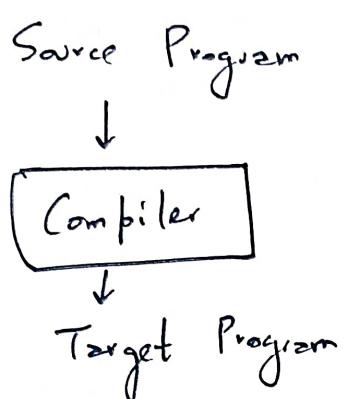
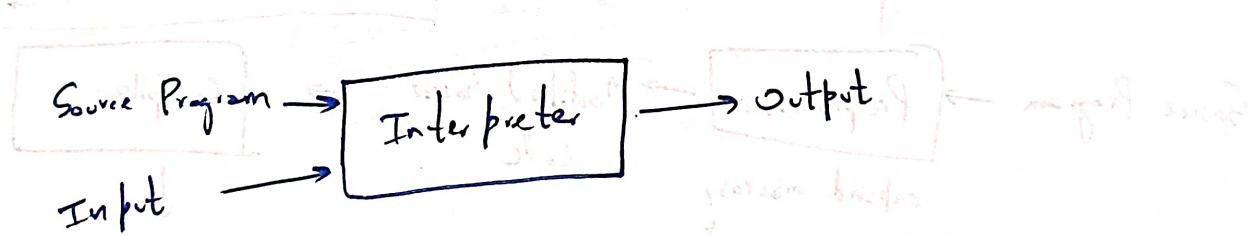


## Compiler Design

→ Before a program can be run, it must first be translated into certain form (target program) in which it can be executed by a computer.



Interpreter: Instead of producing Target Program as a translation, an interpreter appears as a translation to directly executes the operations specified in the source program on inputs supplied by user.



Assembler converts Assembly Level language to binary.

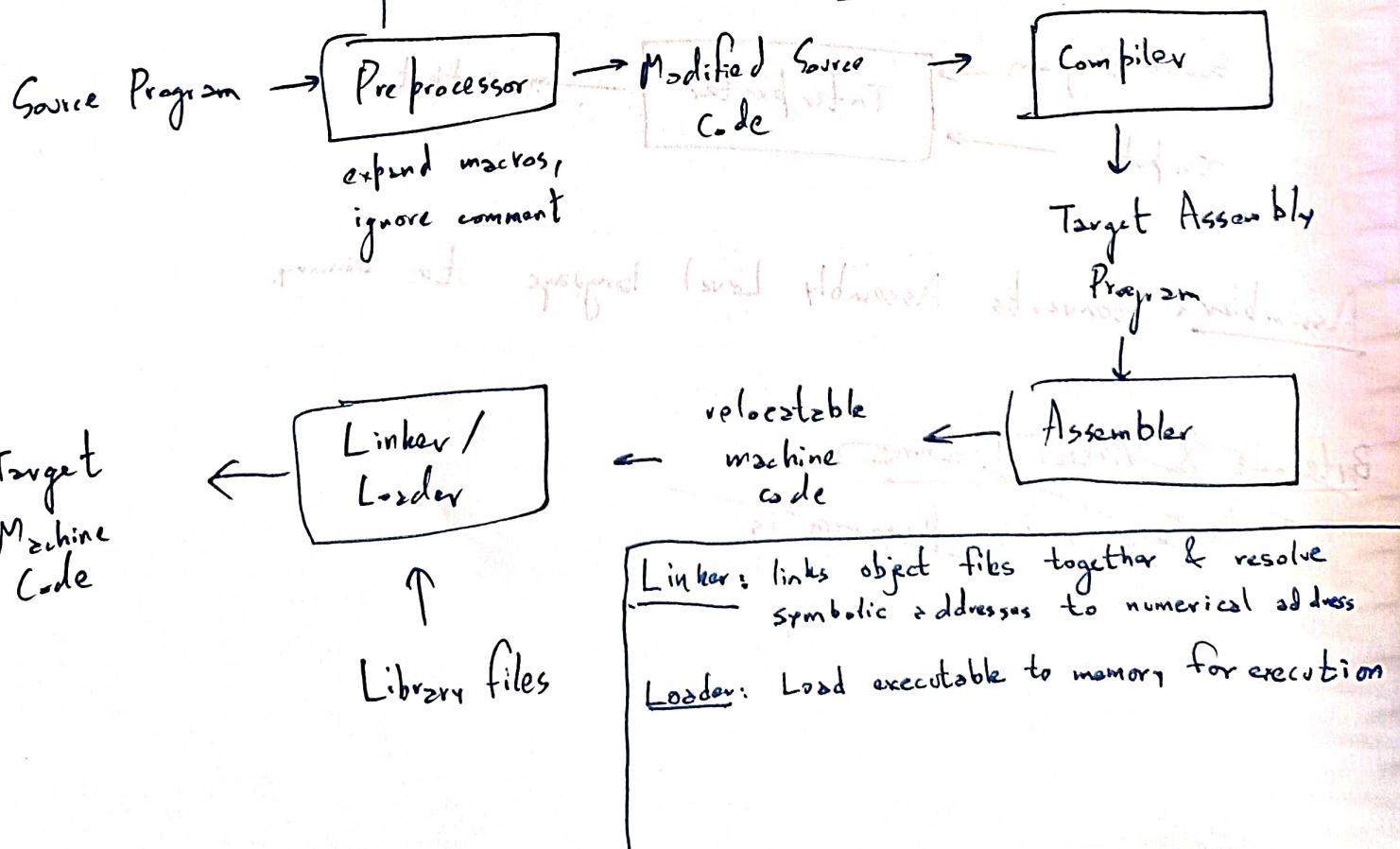
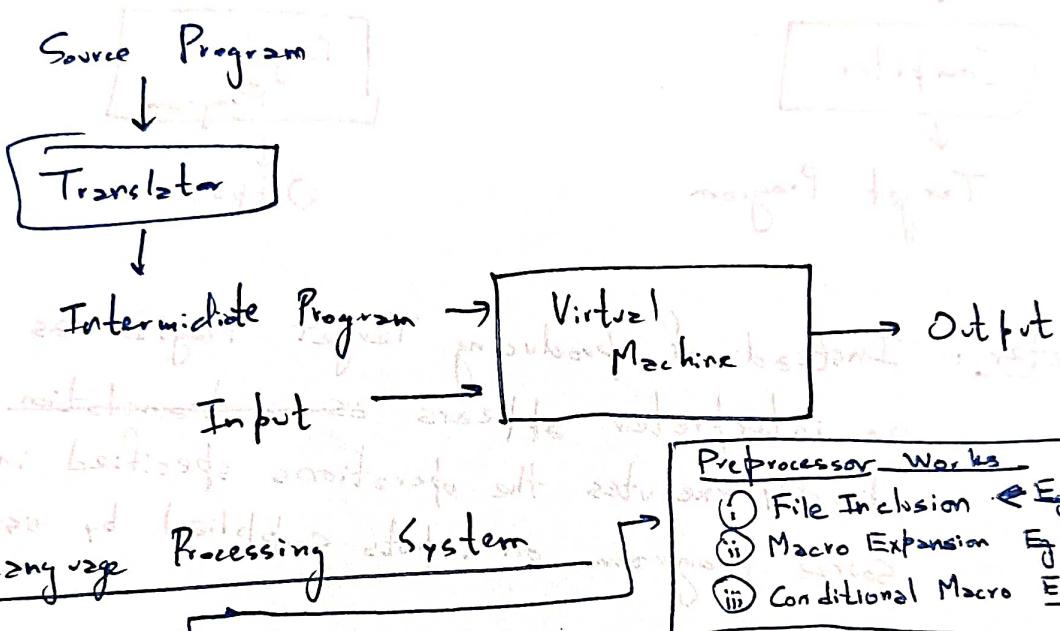
Bytecode & Virtual Machine

→ A Java source program is



## Bytecode & Virtual Machine

- A Java Source code may first be compiled into an intermediate form called bytecode, which are then interpreted by a virtual machine.
- A bytecode compiled on one machine can be interpreted on another machine.



# Analysis-Synthesis Model of Compilation

→ Two parts of compilation

## (I) Analysis

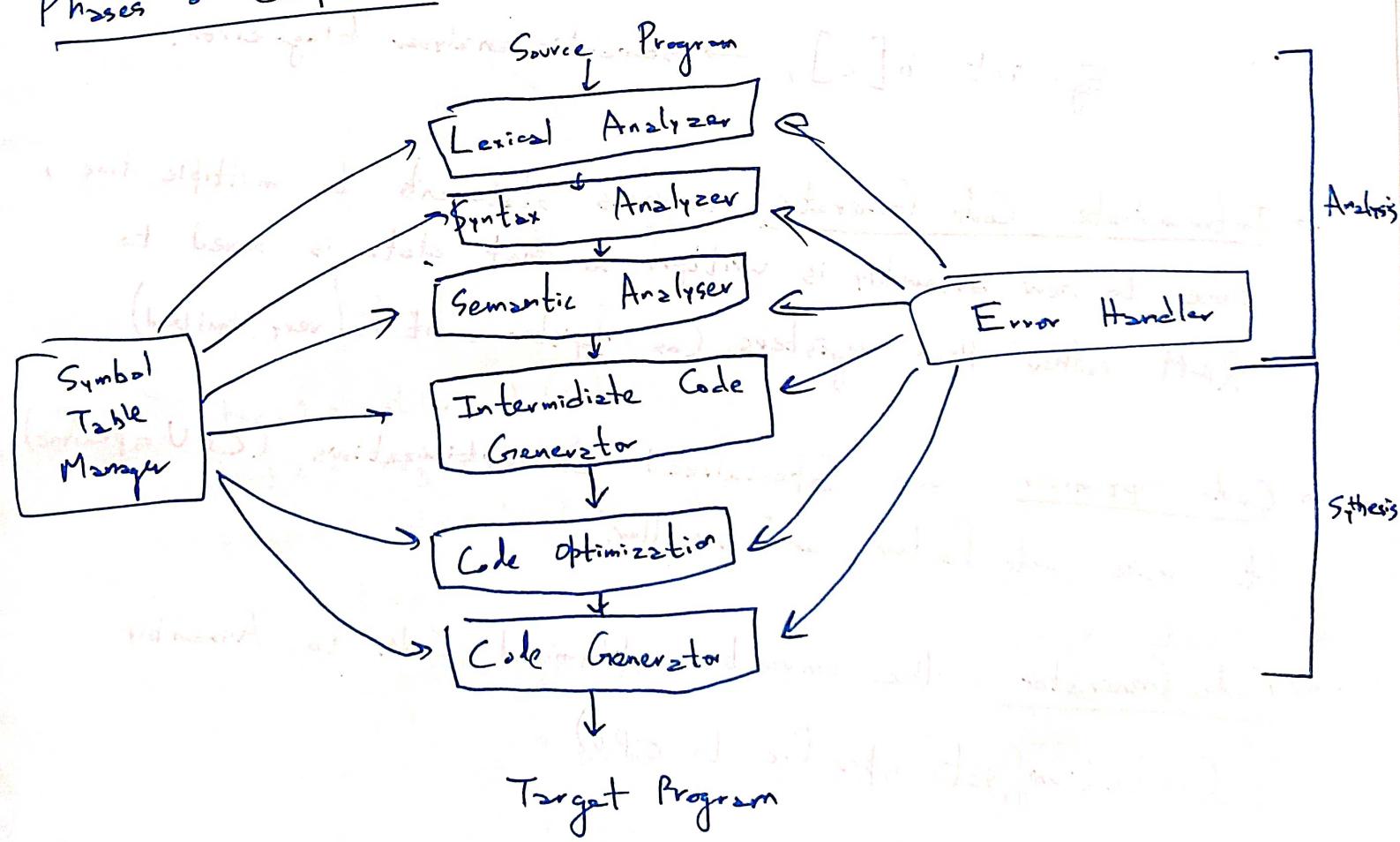
→ breaks up source program into constituent pieces & creates an intermediate representation of program.

→ If the analysis part detects that the source program is either syntactically ill formed or semantically unsound, then it must provide informative messages so that user can take corrective actions.

## (II) Synthesis

→ It constructs the desired target program from the intermediate representation & the info in symbol table.

## Phases of Compilation



→ Lexical Analyzer takes the source program as input & produces  
≥ long string of tokens.

→ Syntax Analyser checks the syntax of statements and  
generates parse tree.

→ It's also called parser.

→ Eg: yacc

[Symbol Table generated in Lexical Analysis or Syntax Analysis]

→ Semantic Analyser checks the semantics and produces  
semantic tree.

Eg: int z = b + c; → semantic analyser check if  
b & c ~~is~~<sup>possible</sup> or not  
(type checking, scope, etc)

Eg: int b[c]; → semantic analyser flag error.

→ Intermediate Code Generator converts statements to multiple line,  
close to how assembly is written so that data is saved to  
RAM rather than registers. (as register count is very limited)

→ Code Optimizer uses Specialized CPU optimizations (CPU specific)  
to make code faster and smaller.

→ Code Generator converts optimized code to Assembly  
(instruction set specific to CPU)

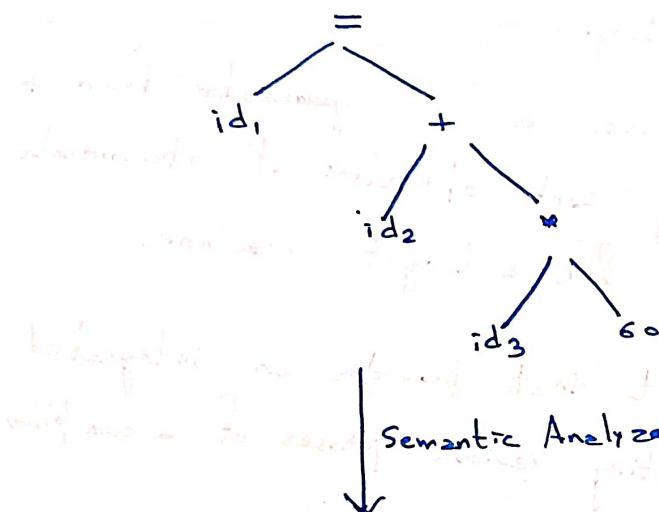
Eg

$pos = init + rate * 60$

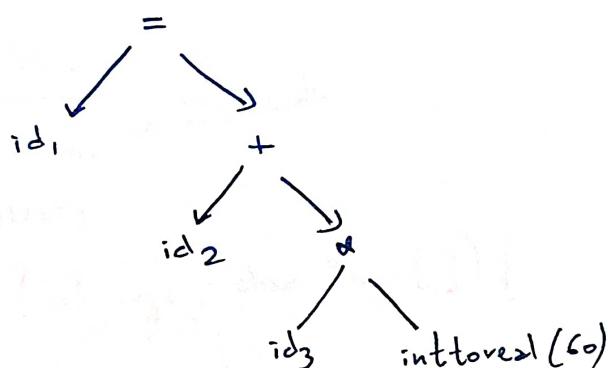
↓ Lexical Analyzer

$id_1 = id_2 + id_3 * 60$

↓ Syntax Analyzer



↓ Semantic Analyzer



↓ Intermediate code generator

temp1 = inttoreal(60)

temp2 = id3 \* temp1

temp3 = id2 + temp2

$id_1 = temp3$

↓ Code Optimization

$temp1 = id3 * 60.0$

$id_1 = id_2 + temp1$

MOVF id3, R2

MULF #60.0, R2

MOVF id2, R1

ADDF R2, R1

MOVF R1, id1

↓ Code Generator

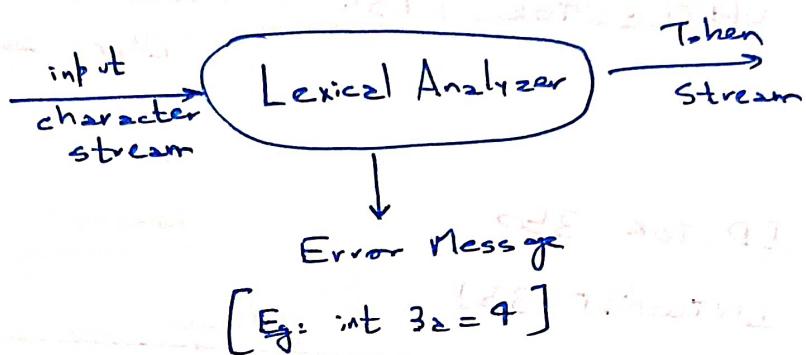
## Compiler Construction Tools

- i) Scanner Generators that does lexical analysis from regular expression.  
Eg: lex
- ii) Parser Generators that automatically produces syntax analyzer.  
Eg: yacc
- iii) Syntax Directed translation Engine
- iv) Code Generators that produces a code generator from a set of rules for translating each operation of intermediate language into machine language for a target machine.
- v) Compiler Construction toolkit that provides an integrated set of routines for constructing various phases of a compiler.

# LEXICAL ANALYZER

- converts the input program into a sequence of tokens.
- implemented using Finite Automata

while (*i*>2)  
*i* = *i*-2;



WHILE-TOK
L PAREN -TOK
ID -TOK
G T -TOK
INTCONST -TOK
R PAREN -TOK
ID -TOK
EQ -TOK
ID -TOK
MINUS -TOK
INTCONST -TOK
SEMICOLON -TOK

## Programmer's View

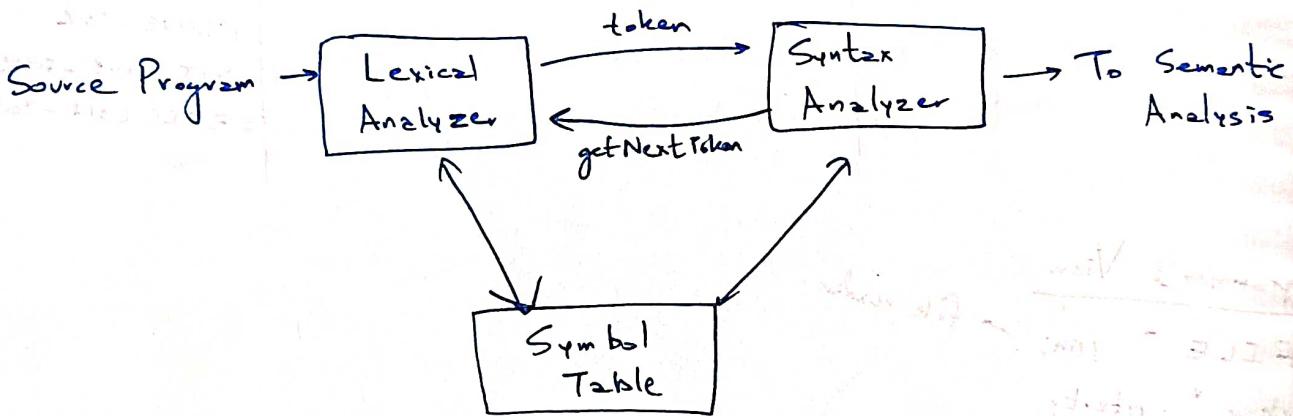
```

FILE * yyin;           → file reader
char * yytext;
void main (int argc, char * argv[])
{
    int token;
    if (argc != 2) return;

    yyin = fopen(argv[1], "r");
    while (!feof(yyin)) {
        token = yylex(); → returns the next token in yyin
        printf("%d", token);
    }
    fclose(yyin);
}
  
```

```
#define LPAREN_TOK " ("  
#define GT_TOK   '>'  
#define RPAREN_TOK ")"  
;
```

#define WHILE\_TOK 256 → from 256 cause single  
char ascii fill 255



Eg: int a,b,c; → will give ID-TOK with symbol table  
; identifier.

Specification of Token.

- Regular Expression
- Deterministic Finite Automata
- Non-Deterministic Finite Automata.
- Non-Deterministic Finite Automata with empty Transitions.

## Maximal Match Rule (or Longest Prefix Sum)

Eg: int n = 4 + ++z; → so (increment, plus)  
 ↗ these two taken (longest prefix)

## Lexical Analyzer

- Main task is to read the input characters & produce as output a sequence of tokens.
- Avoiding comments and white spaces in the form of blank, tab and newline character.
- Also correlates error message from the compiler with same source program

## Token, Patterns & Lexemes

- A token is a pair consisting of a token name and an optional attribute value.
- Token name is an abstract symbol representing a kind of lexical unit, ~~is~~
- Eg: a particular keyword, or a sequence of input characters denoting an identifier.
- This set of strings is described by a rule called pattern.
- A lexeme is a sequence of characters in the source program that matches the pattern for a token & is identified by the lexical analyzer as an instance of that token.
- These are smallest logical unit (words) of a program such as A, B, 1.0, true, +, <=, ...

Eg: `printf ("Total = %d", score);`

↓                    ↓                    ↓

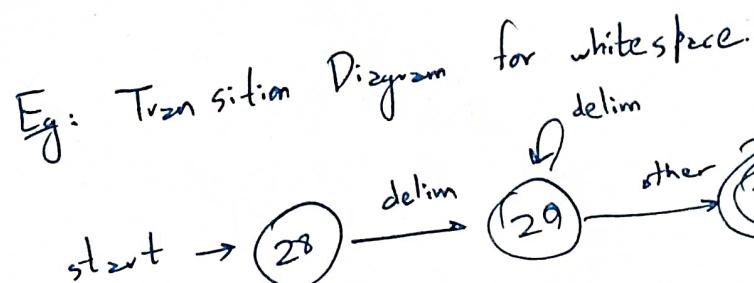
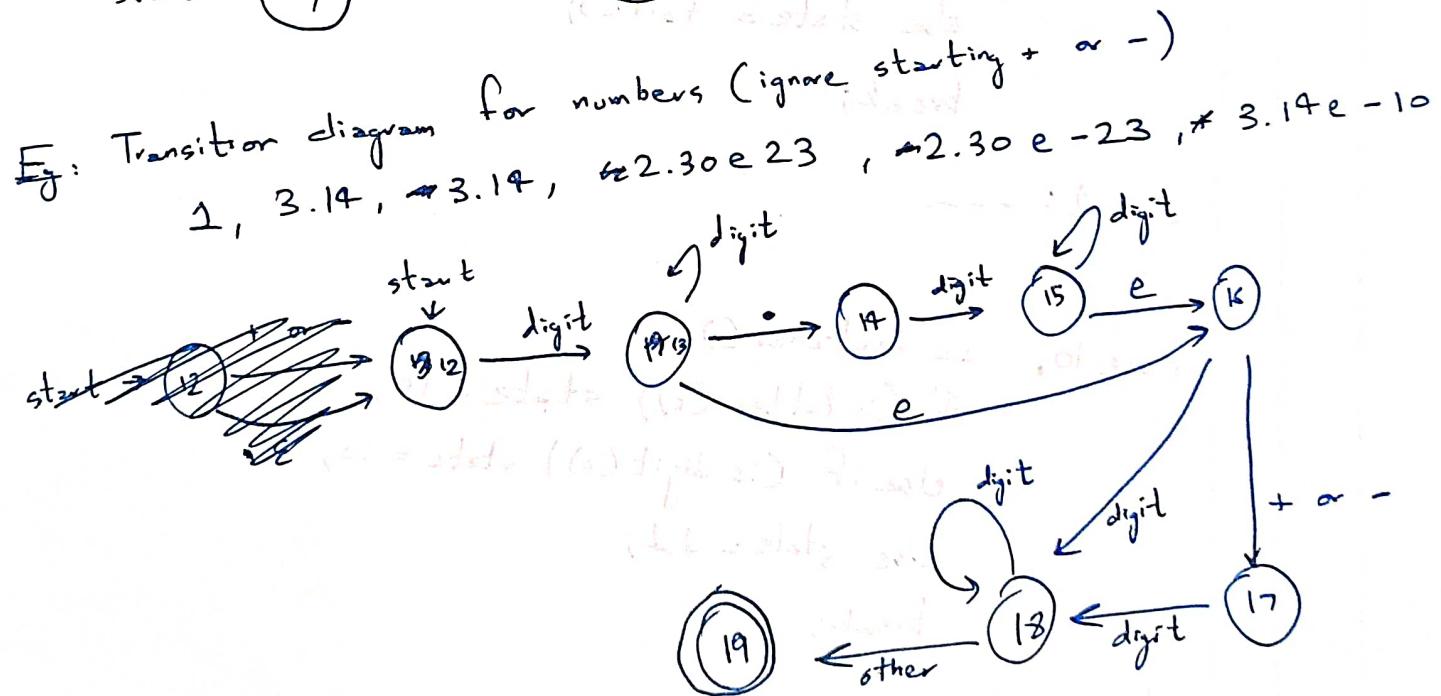
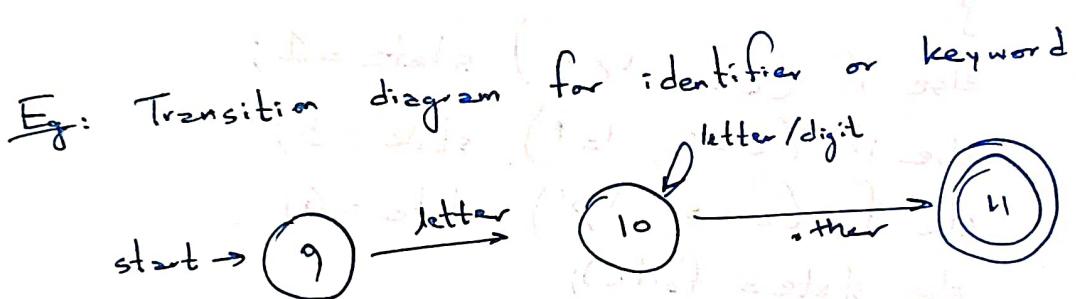
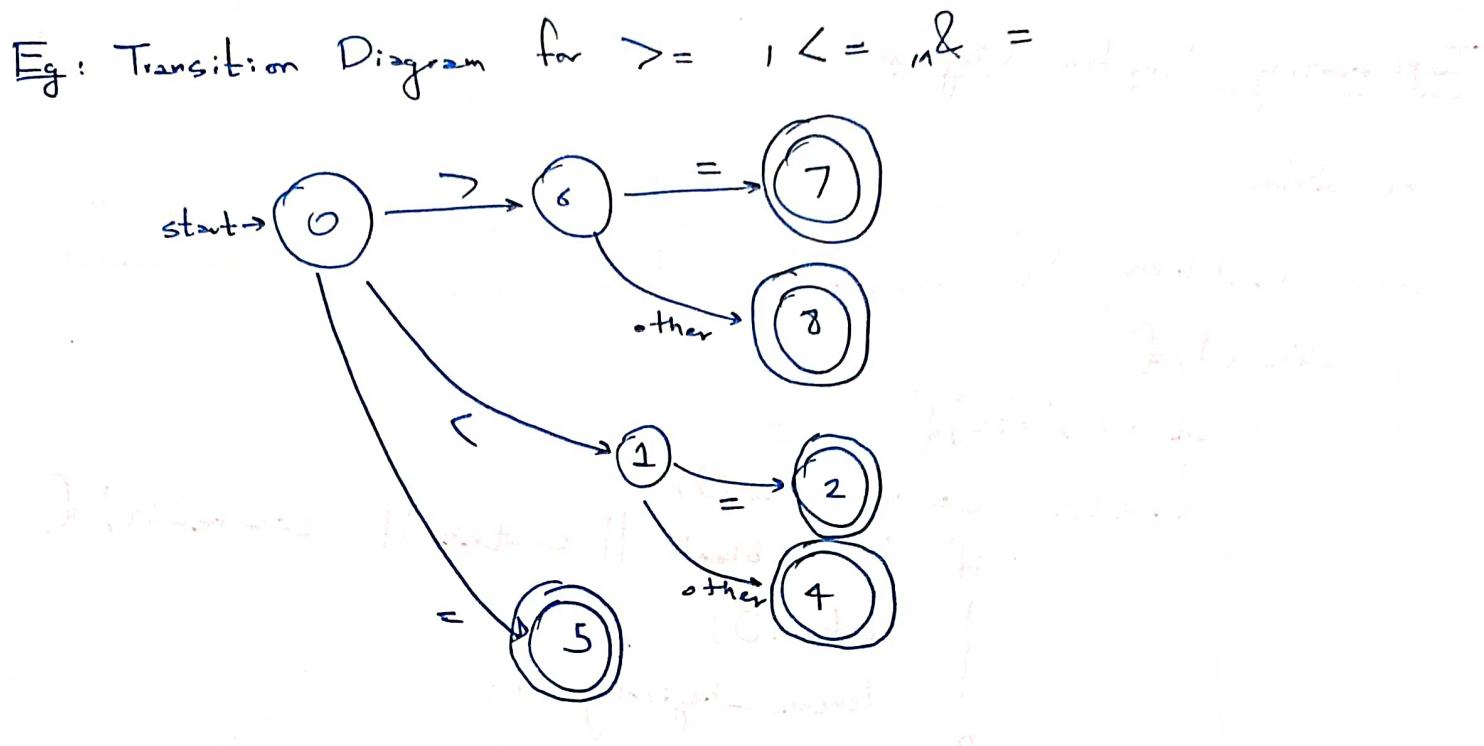
lexeme matching pattern of id      lexeme matching literal.      lexeme matching pattern of id

Token	Sample Lexemes	Informal description of pattern
I if	if	characters 'i', 'f'
II else	else	characters 'e', 'l', 's', 'e'
III comparison	<, <=, ==, !=, >, >=	< or > or <= or >= ...
IV id	pi, score, D2	letter followed by letters and digit
V number	3.14159, 0, 6.02e23	any numeric constant
VI literal	"Total = 7. d"	Total = $y.d$

Eg:  $E = M * C + 2$

$\downarrow$  tokens

- <id, ptr to symbol table entry of E>
- <mult op>
- <#id, ptr to symbol table entry of M>
- <plus op>
- <mult op>
- <#id, ptr to symbol table entry of C>
- <plus op>
- <number, (integer value 2)>



## Implementing Transition Diagram

```
int state = 0;
```

```
token nextToken () {
```

```
    while (1) {
```

```
        switch (state) {
```

```
            case 0: c = nextChar();
```

```
                if (c == blank || c == tab ||
```

```
                    c == newline) {
```

```
                state = 0;
```

```
                lexeme += beginning++;
```

```
            } else if (c == '<') state = 1;
```

```
            else if (c == '=') state = 5;
```

```
            else if (c == '>') state = 6;
```

```
            else state = fail();
```

```
            break;
```

```
        case 1: ---
```

```
        case 10: c = nextChar();
```

```
                if (isLetter(c)) state = 10;
```

```
                else if (isDigit(c)) state = 10;
```

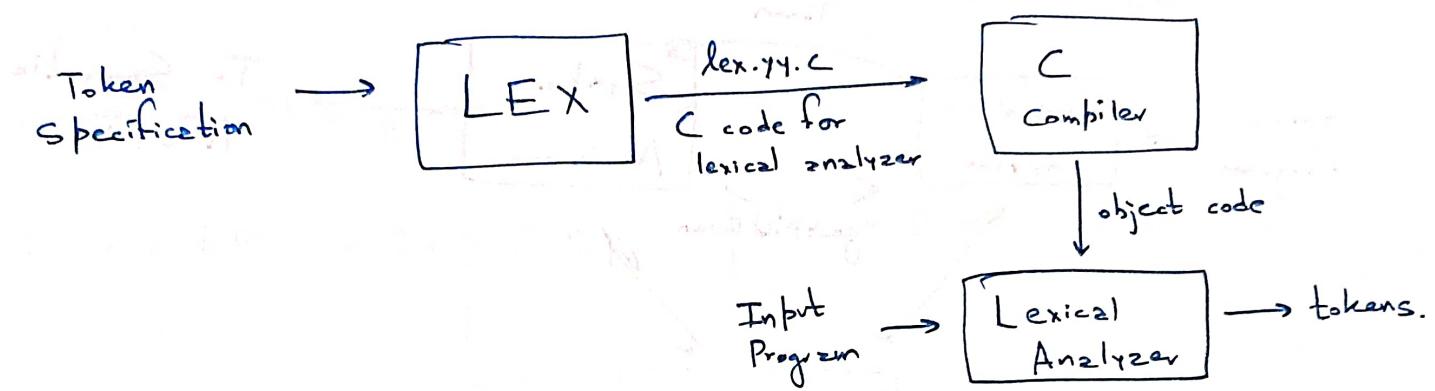
```
                else state = 11;
```

```
                break;
```

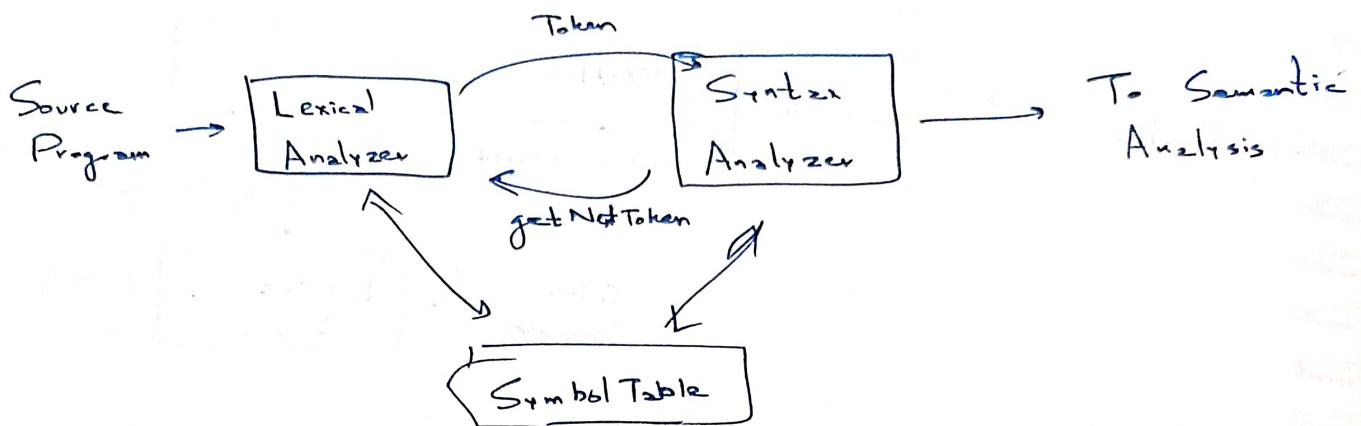
```
        }
```

```
}
```

## Lex



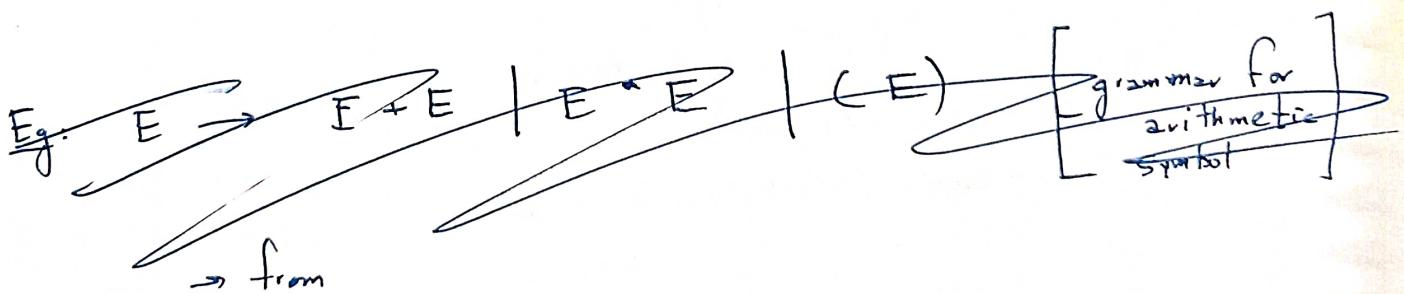
## SYNTAX ANALYZER



- checks syntax of input to determine if it is correct or not
- also keeps track of symbol via symbol table.
- Syntax can be defined using Context Free Grammar

$$G_1 = (V, \Sigma, P, S)$$

↑              ↑              ↑              ↗  
 finite variables   terminals   starting Variable.  
 finite Production Rules

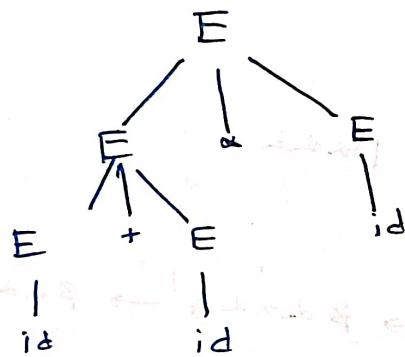


Eg:  $E \rightarrow E+E \mid E^*E \mid (E) \mid id$  [grammar for arithmetic symbol]

→ from this grammar "id + id \* id" is derivable

$$\begin{array}{c} E \rightarrow E+E \rightarrow E+E^*E \\ \downarrow \\ id + E^*E \\ \downarrow \\ id + id^*id \end{array}$$

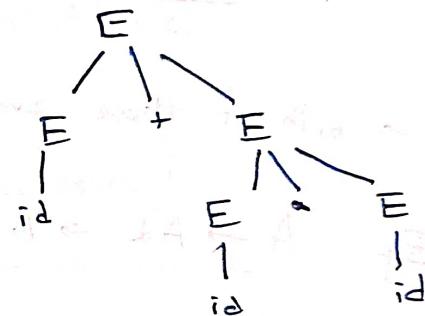
& corresponding parse tree



[leftmost derivation]

$$\begin{array}{c} E \rightarrow E+E \rightarrow E+E^*E \\ \downarrow \\ E+E^*id \\ \downarrow \\ id + id^*id \end{array}$$

& corresponding parse tree



[rightmost derivation]

→ Grammar can be ambiguous (should we use leftmost or rightmost derivation?)

→ ~~We~~ In best case, grammar should not be ambiguous.

## Elimination of Left Recursion

(by variable being in left)

→ Rules which can recursively grow longer in left can cause problem in some situations

Eg:  $A \rightarrow A\alpha | \beta$

we convert it

$$A \rightarrow \cancel{\alpha} A'$$

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

here variable on right, easy  
for top down parser.

→ both of these are same

Eg  $A \rightarrow A\alpha \rightarrow A\alpha\alpha \rightarrow A\alpha\alpha\alpha \rightarrow \beta\alpha\alpha\alpha$

$A \rightarrow \beta A' \rightarrow \beta\alpha A' \rightarrow \beta\alpha\alpha A' \rightarrow \beta\alpha\alpha\alpha A' \rightarrow \beta\alpha\alpha\alpha\alpha A'$

Q) Remove Left Recursion if present

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

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

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

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

(I)  $E \rightarrow E + T \mid T$   $\Rightarrow$   $E \rightarrow TA'$   
 $A' \rightarrow +TA' \mid \epsilon$

(II)  $T \rightarrow T * F \mid F$   $\Rightarrow$   $T \rightarrow FA'$   
 $A' \rightarrow *FA' \mid \epsilon$

(III)  $F \rightarrow (E) \mid id$  [not left recursive]

→ In general for

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid A\alpha_3 \mid \dots \mid \beta_1 \mid \beta_2 \mid \beta_3 \mid \dots$$



$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots$$

$$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \alpha_3 A' \mid \dots \mid \epsilon$$

## Left Factoring of Grammar

consider

$$\begin{array}{l} A \rightarrow \alpha B_1 \mid \alpha B_2 \\ \downarrow \\ A \rightarrow \alpha A' \quad \cancel{\mid \alpha} \\ A' \rightarrow B_1 \mid B_2 \end{array}$$

(ambiguous)

→ here  $B_1$  &  $B_2$  will be taken so choice will be based on next input

Eg:

$$\begin{array}{l} \text{statement if expression} \quad \text{statement} \\ \text{then} \quad \text{else} \quad \text{statement} \\ S \rightarrow i E t S \quad | \quad i E t S e S \quad | \quad a \quad \cancel{| \quad a} \\ E \rightarrow b \\ \downarrow \\ S \rightarrow i E t S A' \quad | \quad a \\ A' \rightarrow \epsilon \quad | \quad e S \\ E \rightarrow b \end{array}$$

## Different Types of Parsers

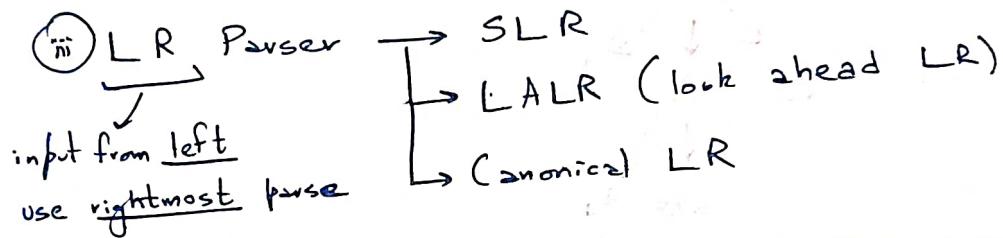
### (I) Top Down Parser

- i) Recursive Descent Parsing.
- ii) Predictive Parser
- iii) Non-Recursive Predictive Parser.

### (II) Bottom-Up Parsing.

- i) Shift Reduce Parser

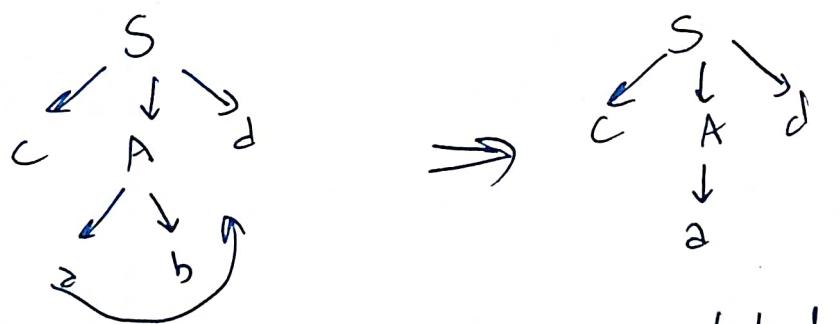
- ii) Operator Precedence Parser



Eg: consider  $S \rightarrow c A d$ ,  $w = "c z d"$

$$A \rightarrow z b | a$$

- i) Recursive Descent  $\rightarrow$  Recursively check every case



did not match  
back track  
& try other option

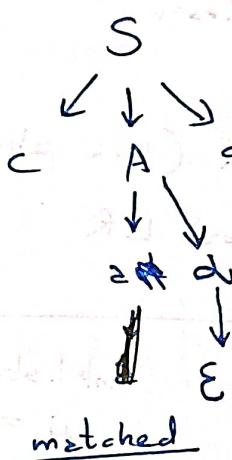
matched.

ii) Predictive Parser: use left factorization & then just keep matching.

$$S \rightarrow cAd \quad S \rightarrow cAd$$

$$A \rightarrow ab|a \quad \Rightarrow \quad A \rightarrow ad \quad \left. \begin{array}{l} \\ d \rightarrow b|\epsilon \end{array} \right\} \text{left factor}$$

$$w = "c a d"$$



## Model of a Table Driven Nonrecursive Predictive Parser

i) Input

ii) Stack

iii) Predictive Parsing Program

iv) Parsing Table ( $M$ )

v) Output

Eg: consider following Grammar

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

→ it's corresponding Parsing Table (how will come later) be

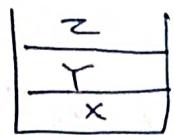
Non terminal ↓	Terminal id	+	*	(	)	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow id$			$F \rightarrow ( E )$		

consider parsing

" id + id \* id "

Stack	Input	Output	Note
\$E	id + id * id \$		Initial Step always like this
E\$ E'T	id + id * id \$	$E \rightarrow TE'$	# E/id in M
\$ E'T'F	id + id * id \$	$T \rightarrow FT'$	T/id in M
\$ E'T'id	id + id * id \$	$F \rightarrow id$	F/id in M
\$ E'T'	+ id * id \$		Terminal in stack, remove input tok
\$ E'	+ id * id \$	$T' \rightarrow \epsilon$	T'/+ in M
\$ E'T+	+ id * id \$	$E' \rightarrow +TE'$	E'/+ in M
\$ E'T	id * id \$		Terminal in stack remove +
\$ E'T'F	id * id \$	$T \rightarrow FT'$	# T/id in M
\$ E'T'id	id * id \$	$F \rightarrow id$	F/id in M
\$ E'T'	* id \$		To ken/input rem
\$ E'T'F*	* id \$	$T' \rightarrow \alpha FT'$	T'/* in M
\$ E'T'F	id \$		To ken/input rem
\$ E'T'id	id \$	$F \rightarrow id$	F/id in M
\$ E'T'	\$		To ken/Input Rem
\$ E'	\$	$T' \rightarrow \epsilon$	\$ / T' in M.
\$	\$	$E' \rightarrow \epsilon$	\$ / E' in M
Parsing Done	Parsing Done		To ken / In Rem

$\rightarrow \#$  stack



i) if  $x = a = \$$  terminal, the parser halts & announces successful completion of parsing

ii) if  $x = a \neq \$$  terminal, the parser pops off the stack & advances the pointer to the next input token

iii) If  $x$  is non terminal, the program consults  $M[x, a]$  of parsing table  $M$ . The entry will either be an  $X$ -production of grammar or nothing (in this case, we throw error)

Eg: if  $M[x, a] \Rightarrow \{X \rightarrow UVW\}$ , the parser replaces  $X$  on top of stack by  $WVU$  ( $V$  on top)

### LL(1) Grammar

→ A grammar whose parsing table has no multiply defined entries is said to be LL(1).

$\rightarrow L \ L \ (1)$

- use one input symbol of lookahead at each step for decision.
- use leftmost derivation.
- scan input from L to R

## Construction of Parsing Table

→ we use two functions FIRST() & FOLLOW()

$$G = (V, \Sigma, P, S)$$

(I) FIRST(X) → compute for all grammar symbol  $x \in V$

- i) if  $X$  is a terminal, then  $\text{FIRST}(X) = \{X\}$
- ii) if  $X \rightarrow \Sigma$  is a production, then add  $\Sigma$  to  $\text{FIRST}(X)$
- iii) if  $X$  is a non-terminal &  $X \rightarrow Y_1, Y_2, \dots, Y_n$  is a production, add  $\text{FIRST}(Y_i)$  to  $\text{FIRST}(X)$

Eg:  $E \rightarrow TE'$        $V = \{E, T, E', F, T'\}$   
 $E' \rightarrow +TE' | \epsilon$        $\Sigma = \{+, *, id, (, )\}$   
 $T \rightarrow FT'$   
 $T' \rightarrow *FT' | \epsilon$   
 $F \rightarrow (E) | id$

②  $\text{FIRST}(F) = \text{FIRST}((E)) \cup \text{FIRST}(id) [F \rightarrow (E) | id]$   
 $= \text{FIRST}(C) \cup \{id\} [iii \& i]$

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

③  $\text{FIRST}(E) = \text{FIRST}(TE') [E \rightarrow TE']$   
 $= \text{FIRST}(T) [iii]$   
 $= \text{FIRST}(FT') [T \rightarrow FT']$   
 $= \text{FIRST}(F) [iii]$

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

$$\boxed{\text{FIRST}(T) = \{(, id\}}$$

NOTE: CFG will have one non-terminal in  
1 LLC

$$\begin{aligned}
 \textcircled{c} \quad \text{FIRST}(E') &= \text{FIRST}(+TE') \cup \text{FIRST}(\epsilon) \quad [E' \rightarrow +TE' \mid \epsilon] \\
 &= \text{FIRST}(+) \cup \{\epsilon\} \quad [\text{iii} \& \text{ii}] \\
 &= \{+\} \cup \{\epsilon\} \quad [\text{i}] \\
 \boxed{\text{FIRST}(E') = \{+, \epsilon\}}
 \end{aligned}$$

$$\begin{aligned}
 \textcircled{d} \quad \text{FIRST}(T') &= \text{FIRST}(xFT') \cup \text{FIRST}(\epsilon) \quad [T' \rightarrow xFT' \mid \epsilon] \\
 &= \text{FIRST}(x) \cup \{\epsilon\} \quad [\text{iii} \& \text{ii}] \\
 &= \{x\} \cup \{\epsilon\} \quad [\text{i}] \\
 \boxed{\text{FIRST}(T') = \{x, \epsilon\}}
 \end{aligned}$$

- II compute FOLLOW( $X$ ) for all non-terminals  $A, X \in V$
- Place  $\$$  in FOLLOW( $S$ ) [ $S$  is the start symbol of grammar]  
[\$\\$ is the input right end marker]
  - If there is a production  $A \rightarrow \alpha B \beta$ , then everything in  $\text{FIRST}(\beta)$  except  $\epsilon$  will be placed in FOLLOW( $B$ )
  - If there is a production  $A \rightarrow \alpha B$  or  $A \rightarrow \alpha B \beta$  where  $\text{FIRST}(\beta)$  contains  $\epsilon$ , then everything in FOLLOW( $A$ ) will be placed in FOLLOW( $B$ )

Eg: ② FOLLOW(E)

$\rightarrow E$  is starting symbol  $\rightarrow \{\$\}$  [\text{i}]

$\rightarrow$  rules where  $E$  in RHS

①  $\frac{F}{A} \rightarrow \frac{E}{\alpha B} \frac{}{\beta}$   $[\text{FIRST}(E) = \{\alpha\}]$

~~→ rule ②  $\rightarrow \{\}\}$~~

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

### b) FOLLOW(E')

Rules  $\rightarrow$  ①  $E \rightarrow \frac{A}{\alpha} \frac{E'}{\beta} \frac{B}{\beta}$   $[FIRST(B) = \{\epsilon\}]$  [iii]

$$\boxed{\begin{aligned} FOLLOW(E') &= FOLLOW(E) \\ &= \{\$, )\} \end{aligned}}$$

~~②~~  $\rightarrow$  ②  $E' \rightarrow \frac{+}{\alpha} \frac{TE'}{\beta} \frac{B}{\beta}$   $[FIRST(B) = \{\epsilon\}]$  [iii]

$$FOLLOW(E') = FOLLOW(E')$$

### c) FOLLOW(T)

rules with T in RHS  $\rightarrow$  ①  $E \rightarrow \frac{A}{\alpha} \frac{T}{\beta} \frac{E'}{\beta}$   $[FIRST(E') = \{+, \epsilon\}]$

$$\text{rule } \text{iii} \rightarrow FOLLOW(ET) = FOLLOW(E) \\ = \{\$, )\}$$

~~$\text{rule ii} \rightarrow FOLLOW(T) =$~~

②  $E' \rightarrow \frac{+}{\alpha} \frac{TE'}{\beta} \frac{B}{\beta}$   $[FIRST(B) = \{\epsilon\}]$

$$\text{rule } \text{iii} \rightarrow FOLLOW(E') = FOLLOW(E) = \{\$, )\}$$

$$\text{rule ii} \rightarrow FOLLOW(T) = FIRST(E') = \cancel{\{\epsilon\}} - \cancel{\{E\}} \\ = \{+\}$$

$$\therefore FOLLOW(T) = \{\$, ), +\}$$

### d) FOLLOW(T')

Rule ①  $T \rightarrow \frac{FT'}{\alpha} \frac{T}{\beta} \frac{B}{\beta}$   $[FIRST(B) = \{\epsilon\}]$

$$\text{rule iii} \rightarrow FOLLOW(T') = FOLLOW(T) = \{\$, ), +\}$$

②  $T' \rightarrow \frac{FT'}{\alpha} \frac{B}{\beta} \frac{B}{\beta}$   $[FIRST(B) = \{\epsilon\}]$

$$\text{rule iii} \rightarrow FOLLOW(T') = FOLLOW(T')$$

### e) FOLLOW(F)

→ Rules with F in RHS

$$\textcircled{I} \quad T \xrightarrow[A]{B} FT' \quad [\text{FIRST}(B) = \text{FIRST}(T') = \{\star, \epsilon\}]$$

$$\text{rule } \textcircled{iii} \rightarrow \text{FOLLOW}(F) = \text{FOLLOW}(T) = \{\$, ), +\}$$

$$\begin{aligned} \text{rule } \textcircled{ii} \rightarrow \text{FOLLOW}(F) &= \text{FIRST}(T') - \{\epsilon\} \\ &= \{\star\} \end{aligned}$$

$$\therefore \boxed{\text{FOLLOW}(F) = \{\$, ), +, \star\}}$$

### III Use FIRST() & FOLLOW() to make parsing Table,

$$\text{Eg: } \textcircled{a} \quad E \rightarrow TE'$$

$$\text{FIRST}(E) = \{C, \text{id}\} \rightarrow \text{add rule to } [E, C] \& [E, \text{id}]$$

$$\textcircled{b} \quad E' \xrightarrow{+} E'$$

~~consider only this rule~~

$$\text{FIRST}(E') = \text{FIRST}(+TE') = \{+\} \rightarrow \text{add rule to } [E', +]$$

$$\textcircled{c} \quad E' \rightarrow \epsilon$$

$$\text{FIRST}(E') = \{\epsilon\} \rightarrow \text{no addition possible}$$

$$\boxed{\text{FOLLOW}(E') = \{\$, )\}} \rightarrow \text{rule in } [E', \$] \& [E', )]$$

$$\textcircled{d} \quad T \rightarrow FT'$$

$$\text{FIRST}(T) = \{C, \text{id}\} \rightarrow [T, \epsilon C], [T, \text{id}]$$

$$\textcircled{e} \quad T' \rightarrow \star FT'$$

$$\text{here } \text{FIRST}(T') = \{\star\} \rightarrow [T', \star]$$

$$\textcircled{f} \quad T' \rightarrow \epsilon$$

$$\text{FIRST}(T')_{\text{here}} = \{\epsilon\} \Rightarrow \text{FOLLOW}(T') = \{\$, ), +\}$$

$$\hookrightarrow [T', \$], [T', )]$$

$$\textcircled{g} \quad F \rightarrow (E)$$

$$\text{FIRST}(F) = \{(\}$$

$$[T', +]$$

$$\textcircled{h} \quad F \rightarrow \text{id}$$

$$\text{FIRST}(\star F) = \{\text{id}\}$$

hence resulting table

Non Terminal \ Terminals	id	+	*	(	)	\$
E	$E \rightarrow TE'$		$E \rightarrow TE'$			
$E'$		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
$T'$		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow id$			$F \rightarrow ( \text{ } \overline{id} )$		

### Implementation Detail of Non Recursive Predictive Parser

- Table can be implemented in Dictionary / Array / etc.
- Stack is already present
- Input Buffer can be thought as an array

Q)  $S \rightarrow ETS S' | \epsilon$  compute FIRST(), FOLLOW() of every non terminals, make parsing Table, is this grammar LL(1)?  
 $S' \rightarrow eS | \epsilon$   
 $E \rightarrow b$

## Bottom Up PARSING

→ we go from leaf to root [i.e. we go reverse]

Eg.  $S \rightarrow aABe$

$$A \rightarrow Abc \mid b \quad w = "abbcdde"$$

$$B \rightarrow d$$

String	Handle
<u>a b b c d e</u>	$A \rightarrow b$
<u>a A b c d e</u>	$A \rightarrow Abc$
<u>a A d e</u>	$B \rightarrow d$
<u>a A B e</u>	$S \rightarrow aABe$
<u>s</u>	[accept]

→ The Substring we choosing here are called handle

→ Choosing Proper Handle is the challenge

## Handle Pruning

→ Technique in which we do Right Most Derivation in Reverse.

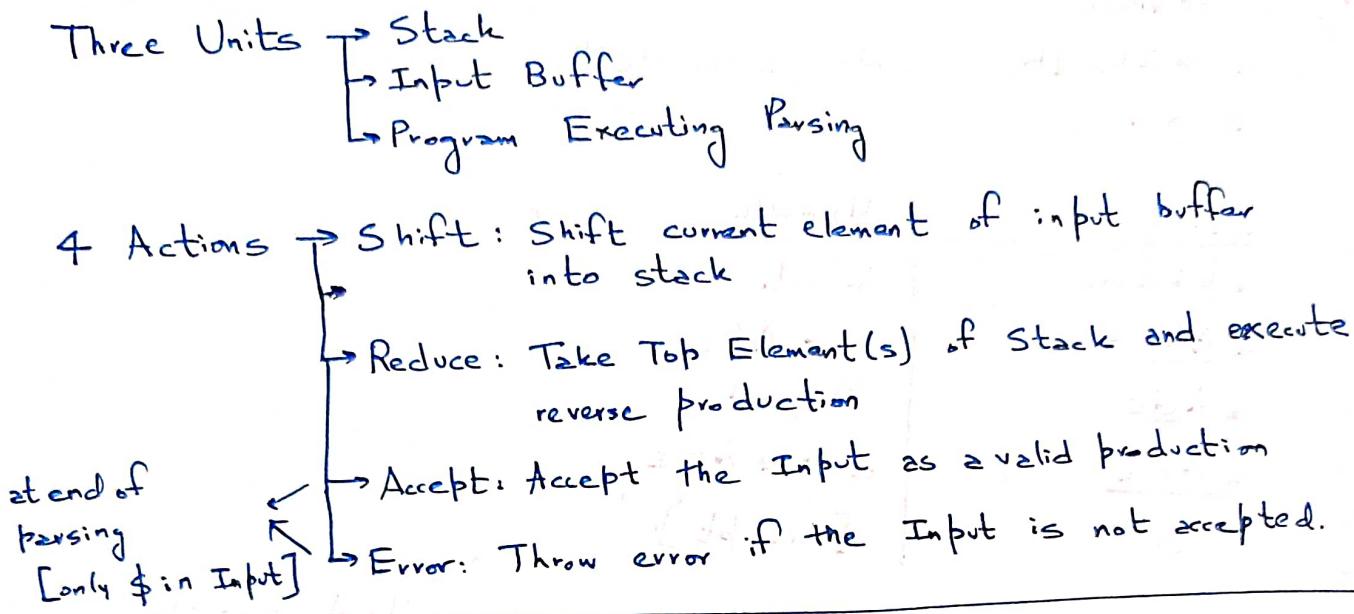
Eg:  $E \rightarrow E + E$  → E is start symbol  
 $E \rightarrow E^* E$  →  $w = "id_1 + id_2 * id_3"$   
 $E \rightarrow C(E)$   
 $E \rightarrow id$

Right Sequential Form	Handle	Reducing Production
$id_1 + id_2 * id_3$	$id_1$	$E \rightarrow id$
$E + id_2 * id_3$	$id_2$	$E \rightarrow id$
$E + E * id_3$	$id_3$	$E \rightarrow id$
$E + E * E$	$E^* E$	$E \rightarrow E^* E$
$E + E$	$E + E$	$E \rightarrow E + E$
$E$		

[ look in reverse way, you can see RMD ]

## Stack Implementation of Shift Reduce Parser

→ Programmatic Method of Implementing Handle Pruning

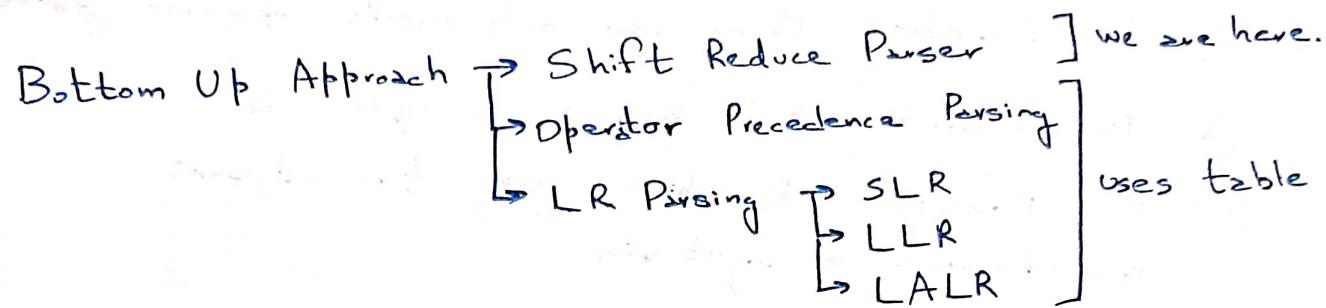


Eg:  $E \rightarrow E + E$   
 $E \rightarrow E * E$   
 $E \rightarrow (E)$   
 $E \rightarrow id$

$$w = id_1 + id_2 * id_3$$

Stack	Input	Action
\$	id <sub>1</sub> + id <sub>2</sub> * id <sub>3</sub> \$	Shift
\$ id <sub>1</sub>	+ id <sub>2</sub> * id <sub>3</sub> \$	Reduce by $E \rightarrow id$
\$ E	+ id <sub>2</sub> * id <sub>3</sub> \$	Shift
\$ E +	id <sub>2</sub> * id <sub>3</sub> \$	Shift
\$ E + id <sub>2</sub>	* id <sub>3</sub> \$	Reduce by $E \rightarrow id$
\$ E + E	* id <sub>3</sub> \$	Shift [not Reduce cause RMD]
\$ E + E *	id <sub>3</sub> \$	Shift
\$ E + E * id <sub>3</sub>	\$	Reduce by $E \rightarrow id$
\$ E + E * E	\$	Reduce by $E \rightarrow E^* E$
\$ E + E E	\$	Reduce by $E \rightarrow E + E$
\$ E	\$	Accept

NOTE → Ambiguity, about choosing Shift or Reduce is still present  
→ Solution to this is by making a Parsing Table



- Each method coming forward will have different method of generating Parsing Table.
- Parsing Procedure will be same as we did now.

## Operator Precedence Parsing

### Operator Grammar

→ Operator Precedence Parser can only handle operator grammar

- (a) No Production Rule have  $\epsilon$  in right side
- (b) No Production rule which have two adjacent non-terminals on right side

Eg:  $E \rightarrow EAE \mid (E) \mid -E \mid id$

consecutive non term

$A \rightarrow + \mid - \mid \ast \mid / \mid \uparrow$

↑  
convert

empty  
Not Operator Grammar

Eg:  $E \rightarrow E+E \mid E-E \mid E \cdot E \mid E/E \mid E \uparrow E \mid -E \mid id$

operator grammar

(E)

### Operator Precedence Parsing Table

Relation	Meaning
$a < b$	$a$ yields precedence to $b$
$a = b$	$a$ has same precedence as $b$
$a > b$	$a$ takes precedence over $b$

→ These three relations are used to generate parsing table

Eg: assume following parsing table. (how will come later)

	id	+	*	\$	
id		>	>	>	
+	<	>	<	>	
*	<	*	>	>	
\$	<	<	<	<	

and we parsing  $w = " id + id * id "$

(I) add delimiter to both side

$\$ id + id * id \$$

↓ using table, add <, > & =

$\$ <id> + <id> * <id> \$$

↓ we use handle (rule) which are between <>  
 $[E \rightarrow id]$

$\$ E + <id> * <id> \$$

↓  $[E \rightarrow id]$

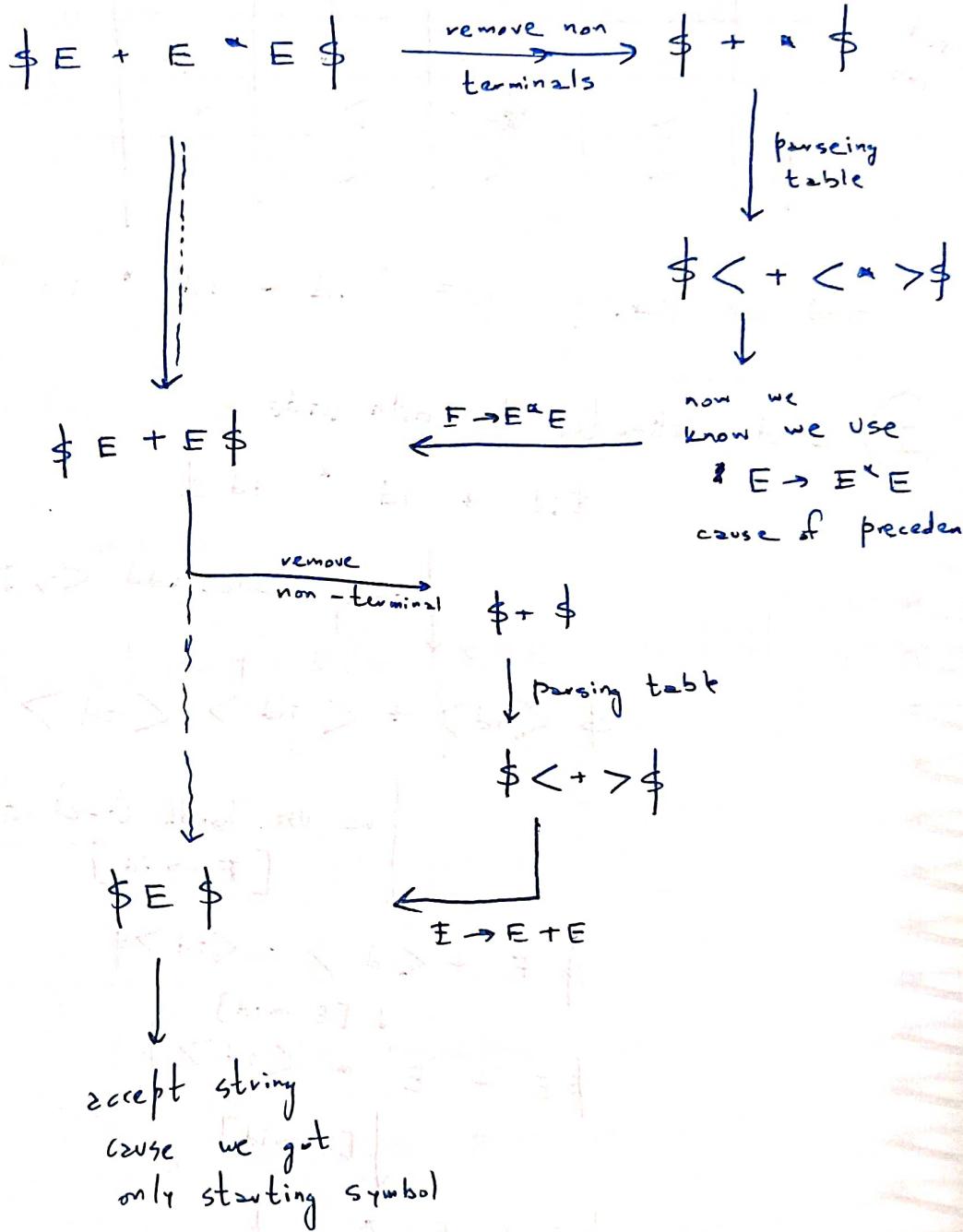
$\$ E + E * <id> \$$

↓  $[E \rightarrow id]$

$\$ E + E * E \$$

$\$ E + E^* E \$$

→ now there is a collision between what to parse, in this case, we follow



## Rules to follow while making Operator Precedence Parsing Table

- ① If operator  $\textcircled{O}_1$  has higher precedence (according to BODMAS) than operator  $\textcircled{O}_2$ , then  $\textcircled{O}_1 > \textcircled{O}_2$  &  $\textcircled{O}_2 < \textcircled{O}_1$
- ② If  $\textcircled{O}_1$  &  $\textcircled{O}_2$  have equal precedence, then  
 $\textcircled{O}_1 > \textcircled{O}_2$  &  $\textcircled{O}_2 > \textcircled{O}_1$  (follow left associativity)
- ③ Relations b/w \$, id & operators & brackets
- i)  $\textcircled{O} < \text{id}$ ;  $\text{id} > \textcircled{O}$
  - ii)  $\textcircled{O} < (\text{ } )$ ;  $(\text{ }) > \textcircled{O}$
  - iii)  $) > \textcircled{O}$ ;  $\textcircled{O} < )$
  - iv)  $\textcircled{O} > \$$ ;  ~~$\textcircled{O} < \$$~~
  - v)  $\$ < (\text{ })$ ;  ~~$\textcircled{O} < \$$~~
  - vi)  $\$ < \text{id}$
  - vii)  $( < )$
  - viii)  $\text{id} > \$$
- etc (use common sense & bodmas)

### Notes

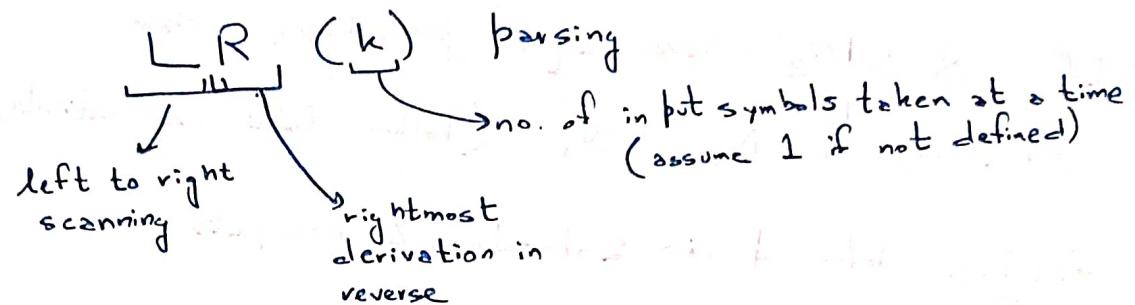
→ This Parsing can only handle a small class of grammars  
(Called operator grammar)

→ There is no way to define unary operators

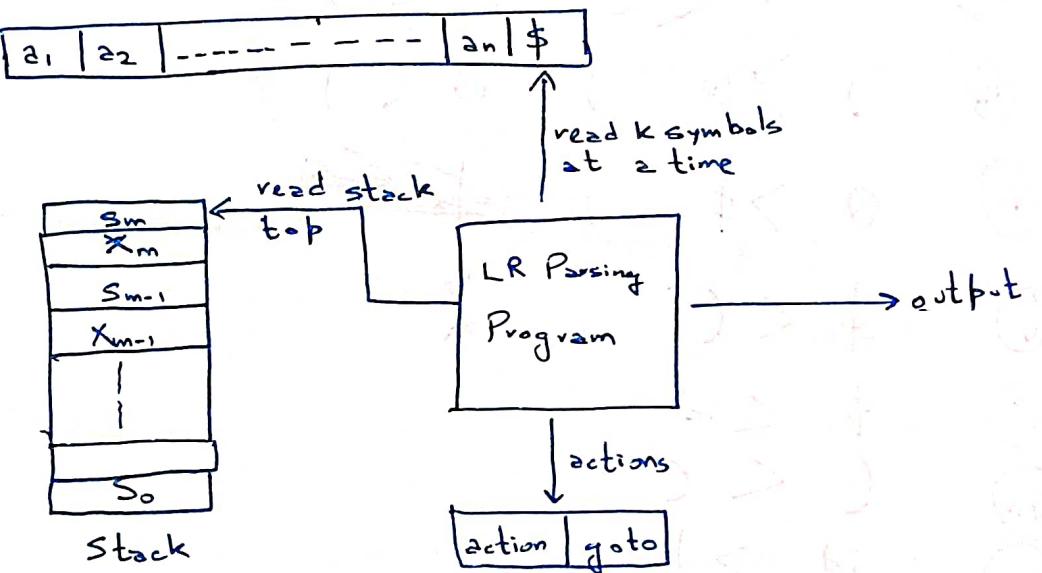
Solution → LR Parsing Table which can handle any grammar.

## LR Parsing

- used in Yacc
- any CFG (hence any programming language) can be parsed using this



## LR Parsing Model



→ Three different methods of constructing parsing Table

- Simple LR [SLR]
- Canonical LR [CLR]
- Look Ahead LR [LALR]

- SLR simple but least powerful
- CLR complex but most powerful
- LALR is intermediate

## How is Parsing Done?

Eg: consider grammar:

- ①  $E \rightarrow E + T$
- ②  $E \rightarrow T$
- ③  $T \rightarrow T * F$
- ④  $T \rightarrow F$
- ⑤  $F \rightarrow (E)$
- ⑥  $F \rightarrow id$

→ consider following SLR parsing table for the grammar (how coming later)

Stack State	action (terminals)						goto (non-term)		
	id	+	*	(	)	\$	E	T	F
0	$S_5$				$S_4$		1	2	3
1		$S_6$				accept			
2		$r_2$	$S_7$		$r_2$	$r_2$			
3		$r_4$	$r_4$		$r_4$	$r_4$			
4	$S_5$			$S_4$			8	2	3
5		$r_6$	$r_6$		$r_6$	$r_6$			
6	$S_5$			$S_4$				9	3
7	$S_5$			$S_4$					10
8		<del><math>r_6 S_6</math></del>			$S_{11}$				
9		$r_1$	$S_7$		$r_1$	$r_1$			
10		$r_3$	$r_3$		$r_3$	$r_3$			
11		$r_5$	$r_5$		$r_5$	$r_5$			

$r_j \rightarrow$  reduce stack element using rule no.  $j$

$s_j \rightarrow$  shift curr input to stack and then post stack state  $j$

acc → accept

blank → throw error

consider  $\omega = "id \Rightarrow id + id"$

Stack	input buffer	Action
0	id * id + id \$	$s_5$
0, id, 5	* id + id \$	$r_6 [F \rightarrow id]$
0, F, 3 (used goto at 0)	* id + id \$	$r_4 [T \rightarrow F]$
0, T, 2 (used goto at 0,T)	* id + id \$	$s_7$
0, T, 2, *, 7	id + id \$	$s_5$
0, T, 2, *, 7, id, 5	+ id \$	$r_6 [F \rightarrow id]$
0, T, 2, *, 7, F, 10 (goto, *, T, F)	+ id \$	$r_3 [T \rightarrow T * F]$
0, T, 2 (goto, 0, T)	+ id \$	$r_2 [E \rightarrow T]$
0, E, 1 (goto, 0, E)	+ id \$	$s_6$
0, E, 1, +, 6	id \$	$s_5$
0, E, 1, +, 6, id, 5	\$	$r_6 [F \rightarrow id]$
0, E, 1, +, 6, F, 3 (goto, *, F)	\$	$r_4 [T \rightarrow F]$
0, E, 1, +, 6, T, 9 (goto, *, T, T)	\$	$r_1 [E \rightarrow E + T]$
0, E, 1 (goto, 0, E)	\$	accept

- after reduce, we also add corresponding goto of non-term being added
- if rule has d variables in RHS, we pop 2d symbols from stack

### NOTE

- we didn't have to remove left recursion or ambiguity
- we add/read variables in symbol table while we parse

## Constructing SLR Parsing Table

LR(0) item  $\rightarrow$  " . " in the RHS of production rule

Eg:  $A \rightarrow XYZ$

$\Downarrow$

$$A \rightarrow \cdot X Y Z$$

$$A \rightarrow X \cdot Y Z$$

$$A \rightarrow X Y \cdot Z$$

$$A \rightarrow X Y Z \cdot$$

→ To make SLR Parsing Table, we need canonical LR(0) items by augmenting the grammar

→ To make canonical LR(0) items, we use closure & goto

(I) ~~we augment grammar by introd~~

→ assume we making parsing table for previous grammar.

- ①  $E \rightarrow E + T$
- ②  $E \rightarrow T$
- ③  $T \rightarrow T * F$
- ④  $T \rightarrow F$
- ⑤  $F \rightarrow (E)$
- ⑥  $F \rightarrow id$

(I) we augment grammar by introducing new starting variable

$$E' \rightarrow E$$

→ if new rule is used, we can safely announce that syntax of sentence is correct [remember this is bottom up parsing]

→ let I = set of one item =  $\{[E' \rightarrow \cdot E]\}$

→ we use two functions closure(I) & goto(I) to find  
canonical LR(0) items

$I_0 = \text{Closure}(I) \doteq \left\{ [E' \rightarrow .E], \dots \right. \rightarrow \text{initial item} \\
[ E \rightarrow .E + T ], \dots \left. \right\} \rightarrow \text{any other rule which have } E \text{ as start symbol} \\
[ E \rightarrow .T ], \dots \left. \right\} \text{ as closure has } .E \\
[ T \rightarrow .T * F ], \dots \left. \right\} \rightarrow \text{any other rule which have } T \text{ as start symbol as closure have } .T \\
[ T \rightarrow .F ], \dots \left. \right\} \\
[ F \rightarrow .(E) ], \dots \left. \right\} \rightarrow \text{any other rule which have } F \text{ as starting symbol as } .F \text{ in closure} \\
[ F \rightarrow .id ] \left. \right\}$

(keep • in start of rule)

→ after . , we can get terminal or non-terminal , for all those symbols ,  
 we keep finding goto( ) till we find all ~~not~~ unique goto's ,  
 (there will be "loops") , we stop that path then )

$$\text{II} \quad \text{goto } (I_0, E) = \left\{ [E' \rightarrow E_0], \quad \dots \rightarrow \text{all items which had } .E \right. \\ \left. [E \rightarrow E_0 + T] \right\} \quad \begin{matrix} \text{in RHS, } .E \text{ shifted one} \\ \text{pos right} \end{matrix} \\ = I_1 \quad (\text{unique})$$

$$\text{III) } \text{goto } (I_0, T) = \{ [E \rightarrow T_0], \dots = I_2 \text{ (Unique)}$$

$$[T \rightarrow T_0 \# F] \}$$

$$\text{IV} \quad \text{goto } (I_0, F) = \{ [T \rightarrow F_s] \} = I_3 \quad (\text{unique})$$

V  $\text{goto}(I_0, C) = \{ [F \rightarrow (\cdot E)],$  after non-terminal, add closures of E  
~~non-terminal,  
also add closure(I<sub>0</sub>)~~

$[E \rightarrow \cdot E + T],$  after non-terminal, add closures of E  
 $[E \rightarrow \cdot T],$  after non-terminal, add closures of T  
 $[T \rightarrow \cdot T * F],$  after non-terminal, add closure of F  
 $[T \rightarrow \cdot F],$  after non-terminal, add closure of F  
 $[F \rightarrow \cdot (E)],$   
 $[F \rightarrow \cdot \text{id}] \} = I_4 \text{ (unique)}$

VI  $\text{goto}(I_0, \text{id}) = \{ [F \rightarrow \text{id}\cdot] \} = I_5 \text{ (unique)}$

VII  $\text{goto}(I_1, +) = \{ [E \rightarrow E + \cdot T],$  after non-terminal, add closures of non-term (from production rules)  
 $[T \rightarrow \cdot T * F],$   
 $[T \rightarrow \cdot F],$   
 $[F \rightarrow \cdot (E)],$   
 $[F \rightarrow \cdot \text{id}] \} = I_6 \text{ (unique)}$

VIII  $\text{goto}(I_2, *) = \{ [T \rightarrow T * \cdot F],$  adding closure of F  
 $[F \rightarrow \cdot (E)],$   
 $[F \rightarrow \cdot \text{id}] \} = I_7 \text{ (unique)}$

IX  $\text{goto}(I_4, E) = \{ [F \rightarrow (E\cdot)],$   
~~E~~  
 $[E \rightarrow E\cdot + T], \} = I_8 \text{ (unique)}$

X  $\text{goto}(I_4, T) = I_2$

XI  $\text{goto}(I_4, F) = I_3$

(XII)  $\text{goto } (I_4, C) = I_4$

(XIII)  $\text{goto } (I_4, \text{id}) = I_5$

(XIV)  $\text{goto } (I_6, T) = \{ [E \rightarrow E + T], [T \rightarrow T \cdot F] \} = I_9 \text{ (unique)}$

(XV)  $\text{goto } (I_6, F) = I_3$

(XVI)  $\text{goto } (I_6, C) = I_4$

(XVII)  $\text{goto } (I_6, \text{id}) = I_5$

(XVIII)  $\text{goto } (I_7, F) = \{ [T \rightarrow T^* F] \} = I_{10} \text{ (unique)}$

(XIX)  $\text{goto } (I_7, C) = I_4$

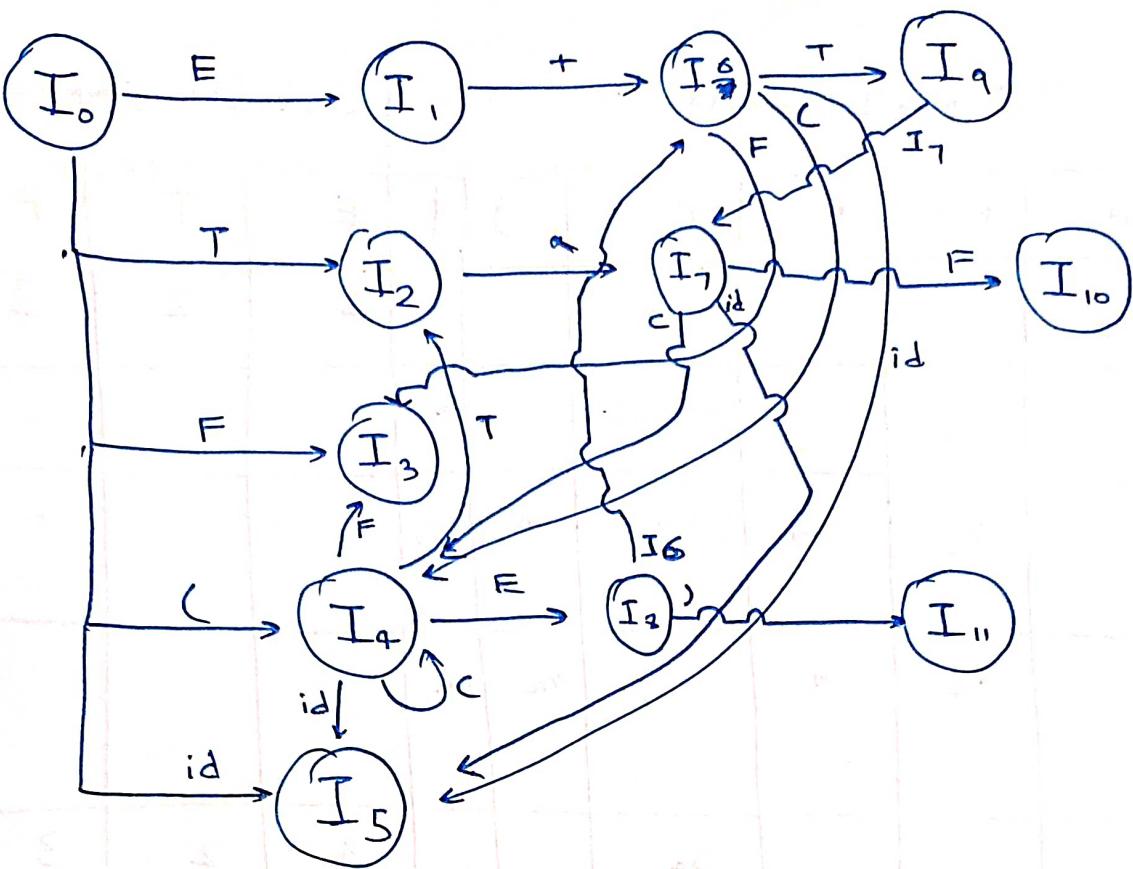
(XX)  $\text{goto } (I_7, \text{id}) = I_5$

(XXI)  $\text{goto } (I_8, )) = \{ [F \rightarrow (E)] \} = I_{11} \text{ (unique)}$

(XXII)  $\text{goto } (I_8, +) = I_6$

(XXIII)  $\text{goto } (I_9, *) = I_7$

→ we essentially tracked all possible ways the grammar can go  
& found 11 different unique states



→ now we make table using goto

$$\text{goto } (I_\alpha, \beta) = I_\gamma \Rightarrow \text{action } [\alpha, \beta] = S_\alpha \text{ if } \beta \text{ non-term}$$

$$\qquad\qquad\qquad \text{goto } [\alpha, \beta] = \gamma \text{ if } \beta \text{ term}$$

→ in cases  $A \rightarrow \alpha .$  for any state  $I_\beta$ , we use FOLLOW(A) to find where to put  $r_{\beta k}$ , where  $k$  is the production rule no.

$$\text{Eg } I_2 \Rightarrow [E \rightarrow T.] \Rightarrow \text{FOLLOW}(E) = \{ \$, +, ) \}$$

↳ rule ②

$\downarrow$  starting symbol       $\downarrow$   $E \rightarrow E + T$        $\downarrow$   $F \rightarrow (E)$   
 (D I Y)

$$\begin{aligned} \text{so action } [2, \$] &= r_2 \\ \text{action } [2, +] &= r_2 \\ \text{action } [2, )] &= r_2 \end{aligned}$$

→ If ~~at start~~ goto possible  $\rightarrow S$

→ If . at end  $\rightarrow \text{FOLLOW}(.)$

$\downarrow \$, \text{ACC}$  if we have augmented rule  $[E' \rightarrow E]$   
in goto set

State	Action (Terminals)						Gotos (Non-Terminal)		
	id	+	*	c	)	\$	E	T	F
0	$s_5$				$s_4$			1	2
1		$s_6$					ACC		
2 (FOLLOW) $\{+, \epsilon, \$, F\}$		$r_2$	$s_7$			$r_2$	$r_2$		
3 FOLLOW(T) = $\{+, \epsilon, \$, F\}$		$r_4$	$r_4$			$r_4$	$r_4$		
4	$s_5$				$s_4$		8	2	3
5 FOLLOW(F) = $\{+, \epsilon, \$, T\}$		$r_6$	$r_6$			$r_6$	$r_6$		
6	$s_5$				$s_4$			9	3
7	$s_5$				$s_4$				10
8	$s_6$					$s_{11}$			
9 FOLLOW(E) = $\{+, \epsilon, \$, F\}$		$r_1$	$s_7$			$r_1$	$r_1$		
10 FOLLOW(T) = $\{+, \epsilon, \$, F\}$		$r_3$	$r_3$			$r_3$	$r_3$		
11 FOLLOW(F) = $\{+, \epsilon, \$, F\}$		$r_5$	$r_5$			$r_5$	$r_5$		

[Same as before]

- There can be cases when one location can have multiple steps which causes ambiguity
  - yacc follows the first step it found in case of ambiguity
- DRAWBACK → we have to use FOLLOW()  $\Rightarrow \overline{F}$
- SOLUTION → ~~LALR~~ Canonical LR table

## Canonical LR

### Construction of Canonical LR parsing table

→ we include extra info with items called look ahead

→ General form of an item becomes

$[A \rightarrow \alpha \cdot \beta, z]$  → This item is called  
 ↑    LR(1) item  
 look                                       ↑  
 ahead                                      amount of  
    extra info  
    used

where  $A \rightarrow \alpha \beta$  is a production rule  
&  $z$  is terminal/right end marker.

→ because of Look Ahead, we don't need to find FOLLOW( $\beta$ ) anymore.

### How look ahead gets calculated?

→ for an item  $[A \rightarrow \alpha \cdot \beta \beta, z]$ , look ahead for its closure = FIRST( $\beta z$ )

Eg:  $[S' \xrightarrow{\text{Look ahead}} \cdot S, \$]$   
 $\rightarrow$  closure for its look ahead = FIRST( $S \$$ ) = { $\$$ }

Eg:  $[S \xrightarrow{\text{Look ahead}} \cdot C, \$]$   
 $\rightarrow$  closure for its look ahead = FIRST( $C \$$ ) = { $\$$ }

→ look ahead for its closure = FIRST( $C C z$ ) = closure( $C$ )  
 $=$  FIRST( $C z$ )  
 $=$  FIRST( $C$ )

According to grammar rule

→ This is followed ~~only~~ only for closures & not gotos

Eg:  $S' \rightarrow S$  [augment rule]

①  $S \rightarrow CC$

②  $C \rightarrow cC$

③  $C \rightarrow d$

① augment grammar

②  $I_0 = \{ [S' \rightarrow .S, \$] \}$  initial state

③  $I_0 = \text{closure}(I) = \{ [S' \rightarrow .S, \$],$  same as before

$[S \rightarrow .CC, \$],$  this

$[C \rightarrow .cC, c/d],$  look ahead from

$[C \rightarrow .d, c/d]$

$\{ \} \Rightarrow$

④  $\text{goto}(I_0, S) = \{ [S' \rightarrow S., \$] \} = I_1 \text{ (unique)}$  look ahead of before

⑤  $\text{goto}(I_0, C) = \{ [S \rightarrow C.C, \$],$

$[C \rightarrow .cC, \$],$

~~$[C \rightarrow .d, \$]$~~   $[C \rightarrow .d, \$] \} = I_2 \text{ (unique)}$

⑥  $\text{goto}(I_0, c) = \{ [C \rightarrow c.C, c/d],$

$[C \rightarrow .cC, c/d],$

$[C \rightarrow .d, c/d] \} = I_3 \text{ (unique)}$

⑦  $\text{goto}(I_0, d) = \{ [C \rightarrow d., c/d] \} = I_4 \text{ (unique)}$

$$\textcircled{VIII} \quad \text{goto}(I_2, C) = \{ [S \rightarrow CC\cdot, \$] \} = I_5$$

$$\textcircled{IX} \quad \text{goto}(I_2, C) = \{ [C \rightarrow C\cdot C, \$], \\ [C \rightarrow \cdot C C, \$], \\ [C \rightarrow \cdot d, \$] \} = I_6 \text{ (unique)}$$

$$\textcircled{X} \quad \text{goto}(I_2, d) = \{ [C \rightarrow d\cdot, \$] \} = I_7 \text{ (unique)}$$

$$\textcircled{XI} \quad \text{goto}(I_3, C) = \{ [C \rightarrow CC\cdot, c/d] \} = I_8 \text{ (unique)}$$

$$\textcircled{XII} \quad \text{goto}(I_3, C) = I_3$$

$$\textcircled{XIII} \quad \text{goto}(I_{\cancel{3}}, d) = I_4$$

~~$$\textcircled{XIV} \quad \text{goto}(I_6, C) = \{ \cancel{C \rightarrow \cdot C} \} = I_6$$~~

$$= \{ [C \rightarrow CC\cdot, \$] \} = I_9$$

$$\textcircled{XV} \quad \text{goto}(I_6, C) = I_6$$

$$\textcircled{XVI} \quad \text{goto}(I_6, d) = I_7$$

→ now we found canonical LR(1) items, now we make Table

→ How to fill reduce rules?

$$\text{Eg: } I_4 = \{ [C \rightarrow d\cdot, c/d] \} \Rightarrow \cancel{r_2} \rightarrow r_3 @ c/d$$

rule 3

→ as before if goto has augment rule, we add ACC

State	action (term)					goto (non-term)
	c	d	\$	s	c	
0	$s_3$	$s_4$		1		2
1			Acc			
2	$s_6$	$s_7$				5
3	$s_3$	$s_4$				8
4	$r_3$	$r_3$				
5			$r_1$			
6	$s_6$	$s_7$				9
7			$v_3$			
8	$r_2$	$r_2$				
9			$r_2$			

→ here we only need to calculate FIRST() in closure  
 hence simpler than SLR

## LALR

→ Steps

- ① construct LR(1) items (same as CLR)
- ② union the states that have same type of production rules but different look ahead

Eg in previous grammar

$$I_3 = \{ [C \rightarrow c.C, c/d], [C \rightarrow .cC, c/d], [C \rightarrow .d, c/d] \}, \quad I_6 = \{ [C \rightarrow c.C, \$], [C \rightarrow .cC, \cancel{\$}], [C \rightarrow .d, \$] \}$$

we union them as  $I_{36} = I_3 \cup I_6$

$$= \{ [C \rightarrow c.C, c/d/\$], [C \rightarrow .cC, c/d/\$], [C \rightarrow .d, c/d/\$] \}$$

similarly,

$$I_4 = \{ [C \rightarrow d., c/d] \}, \quad I_7 = \{ [C \rightarrow d., \$] \}$$

$$\text{Union} \rightarrow I_{47} = I_4 \cup I_7 = \{ [C \rightarrow d., c/d, \$] \}$$

similarly

$$I_{89} = \{ [C \rightarrow c.C., c/d/\$] \}$$

→ resultant table

state	action				goto
	c	d	\$	s	c
0	$S_{36}$	$S_{47}$		1	2
1			ACC		
2	$S_{36}$	$S_{47}$			5
36	$S_{36}$	$S_{47}$			89
47	$r_3$	$r_3$	$r_3$		
5			$r_1$		
89	$r_2$	$r_2$	$r_2$	$v_2$	

→ This reduces the states of CLR

## Error - Recovery Strategies

### Types of Parser

#### I Universal

→ Can parse any grammar

i) Cocke- Younger - Kasami Parsing Method

ii) Earley's Algorithm

→ Inefficient to be used in production of grammar combiners

#### II Top- Down

→ builds parse tree

→ Starts parsing from top (root) to the bottom (leaves)

#### III Bottom- Up

→ starts from leaves and work their way upto the root.

→ Both in Top-Down & Bottom-Up , we

→ scan the input from left to right

→ scan one symbol at a time.

### What should happen when parser finds error in user input?

I Stop immediately & signal error

or II Record error but try to continue.

→ In I), user must recompile from scratch after trivial fix

→ In II), user might be overwhelmed by a whole series of error messages, all caused by essentially the same problem.

## Types of Errors

### i) Lexical Error

→ misspelling an identifier, keyword or operator

### ii) Syntactic Error

→ Eg: arithmetic expression with unbalanced parenthesis

### iii) Semantic Error

Eg: operator applied to an incompatible operand

### iv) Logical Error

Eg: infinitely recursive call.

## Error - Recovery Strategies

### I) Panic Mode

### II) Phrase Level Recovery

### III) Error Production

### IV) Global Correction.

## Panic Mode Error Recovery

→ Idea: skip input symbol until  $\geq$  symbol comes which is in a set of synchronizing token

Eg:  $z = b + c$  //no semicolon  
 $d = e + f ;$

→ The compiler will discard ("d = e + f") all subsequent tokens till  $\geq$  semicolon is encountered.

## Phrase Level Recovery

- Performs local correction on the input to repair the error.
- Changes input streams by inserting missing token.

Eg int id 5;  $\Rightarrow$  int id = 5;  
                                ^ added =

## Error Production

- Augment grammar with productions for erroneous constructs

## Global Corrections

- Choose a minimal sequence of changes to obtain a global least-cost correction.

## Error Recovery in Predictive Parsing

### Error

- An Error detected during predictive parsing when the terminal on the top of stack does not match with the next input symbol.
- ii) when non-terminal  $A$  is on top of the stack,  $a$  is next input symbol, & the parsing table entry  $M[A, a]$  is empty.
- Some Heuristics for the choice of Synchronizing Set are as follows :
  - i) As a starting point, place all symbols in  $\text{FOLLOW}(A)$  into the synchronizing set for the non-terminal  $A$ . Skip the tokens until an element of  $\text{FOLLOW}(A)$  is seen & pop  $A$  from the stack, it is likely that parsing can continue [Panic Mode]
  - ii) If we add symbols in  $\text{FIRST}(A)$  to the synchronizing set for non-terminal  $A$ , then it may be possible to resume parsing according to  $A$  if a symbol in  $\text{FIRST}(A)$  appears in input.
  - iii) If a terminal on top of a stack cannot be matched, a simple idea is to pop the terminal, issue a message saying that the terminal was inserted and continue the parsing.  
[Phrase Level]

$$\text{Eg: } E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

$$\text{FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(F) = \{C, id\}$$

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

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

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

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

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

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

(DIY using  
rules mention  
before)

Non Terminal	id	+	*	(	)	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$	synch	synch
E'		$E' \rightarrow + TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$	synch		$T \rightarrow FT'$	synch	synch
T'		$T' \rightarrow \epsilon$	$T' \rightarrow * FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow id$	synch	synch	$F \rightarrow (E)$	synch	synch

~~synch added using FOLLOW() of resp. non-terminal~~  
~~→ synchronizing symbols~~

- Synch: the driver pops current non-terminal A & skips input till synch token OR skips ~~it~~ until one of FIRST(A) is found (i.e production rule present)
  - If the parser looks up entry  $M[A_{1,2}]$  and finds it blank, the input symbol a is skipped.
  - If the entry is synch, the non terminal on top of stack is popped.
  - If z token on top of stack does not match the input symbol, then we pop the token from stack. (i.e add it to input)
- Eg: consider parsing  $w = "id + id \&"$  using prev parsing table

	Term	id	+	*	(	)	\$
Non-Term							
E	$E \rightarrow TE'$				$E \rightarrow TE'$	synch	synch
$E'$		$E' \rightarrow TE'$				$E' \rightarrow E$	$E' \rightarrow E$
T	$T \rightarrow FT'$	synch			$T' \rightarrow FT'$	synch	synch
$T'$		$T' \rightarrow E$	$T' \rightarrow FT'$			$T' \rightarrow E$	$T' \rightarrow E$
F	$F \rightarrow id$	synch	synch	$F \rightarrow (E)$	synch	synch	

Stack	input	output
\$ E	) id * + id \$	synch error, skip ) (not pop cause nothing in stack) then
\$ E	id * + id \$	$E \rightarrow TE'$
\$ E' T	id * + id \$	$T \rightarrow FT'$
\$ E' T' F	id * + id \$	$F \rightarrow id$
\$ E' T' id	id * + id \$	terminal match
\$ E' T'	* + id \$	$T' \rightarrow FT'$
\$ E' T' F *	* + id \$	terminal match
\$ E' T' F	+ id \$	synch error, pop F
\$ E' T	+ id \$	synch error, pop T
\$ E'	+ id \$	$E' \rightarrow + TE'$
\$ E' T +	+ id \$	terminal match
\$ E' T	id \$	$T \rightarrow FT'$
\$ E' T' F	id \$	$F \rightarrow id$
\$ E' T'	\$	$T' \rightarrow \epsilon$
\$ E'	\$	$E' \rightarrow \epsilon$
\$	\$	Accept

## Phrase Level Recovery in Predictive Parsing

~~Fix~~ Change input stream by inserting missing tokens

Eg: "id id"  $\Rightarrow$  "id a id"

Non Term	id	+	a	(	)	\$
E	$E \rightarrow TE'$			$E' \rightarrow \pi E'$	synch	synch
E'		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$	synch		$T \rightarrow FT'$	synch	synch
T'	Insert "	$T' \rightarrow \epsilon$	$T' \rightarrow \alpha FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow id$	synch	synch	$F \rightarrow (E)$	synch	synch

## Error Production Recovery in Predictive Parsing

- include productions for common error
- these production will make parser pop ~~non~~ terminal, essentially "adding token" in input stream"

Eg: add  $T' \rightarrow FT'$  in  $M[T', id]$

<del>Non Term</del>	id	+	*	(	)	$\neq$
E	$E \rightarrow TE'$			$E' \rightarrow TE'$	synch	synch
$E'$		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow E$
T	$T \rightarrow FT'$	synch		$T \rightarrow FT'$	synch	synch
$T'$	$T' \rightarrow FT'$	$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow id$	synch	synch	$F \rightarrow (E)$	synch	synch

NOTE → Phrase Level Recovery & Recovery via adding Error Production are very powerful, but needs manual additions

## Error Handling in LR Parsing [Phrase Level]

- for each error entry in table (blanks), insert a pointer to a particular error procedure, which assumes the most likely cause and takes appropriate action
- NOTE: also need manual intervention.

## Construction of Parsing Table

→ we use two functions FIRST() & FOLLOW()

$$G = (V, \Sigma, P, S)$$

(I) FIRST(X) → compute for all grammar symbol  $x \in V$

- i if  $x$  is a terminal, then  $\text{FIRST}(x) = \{x\}$
- ii if  $X \rightarrow \epsilon$  is a production, then add  $\epsilon$  to  $\text{FIRST}(X)$
- iii if  $X$  is a non-terminal &  $X \rightarrow Y_1 Y_2 \dots Y_n$  is a production, add  $\text{FIRST}(Y_i)$  to  $\text{FIRST}(X)$

$$\text{Eg: } E \rightarrow TE'$$

$$V = \{E, T, E', F, T'\}$$

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

$$\Sigma = \{+, *, \text{id}, (, )\}$$

$$T \rightarrow FT'$$

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

$$F \rightarrow (E) \mid \text{id}$$

$$\begin{aligned} \text{(a) } \text{FIRST}(F) &= \text{FIRST}((E)) \cup \text{FIRST}(\text{id}) \quad [F \rightarrow (E) \mid \text{id}] \\ &= \text{FIRST}(C) \cup \{\text{id}\} \quad [\text{iii} \& \text{i}] \end{aligned}$$

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

$$\begin{aligned} \text{(b) } \text{FIRST}(E) &= \text{FIRST}(TE') \quad [E \rightarrow TE'] \\ &= \text{FIRST}(T) \quad [\text{iii}] \\ &= \text{FIRST}(FT') \quad [T \rightarrow FT'] \\ &= \text{FIRST}(F) \quad [\text{ii}] \end{aligned}$$

$$\boxed{\text{FIRST}(E) = \{C, \text{id}\}}$$

$$\boxed{\text{FIRST}(T) = \{C, \text{id}\}}$$

NOTE: CFG will have one non-terminal in LHS