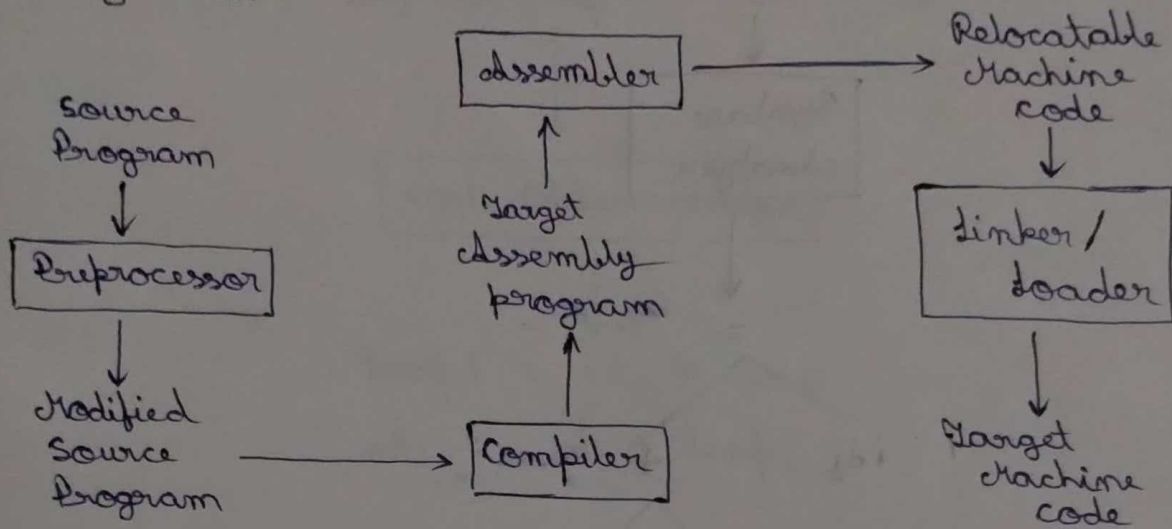


- Compiler is basically a translator. It translates a source program to a target program of a certain processor.
- Interpreter is a language processor which executes line by line.



- Models of Compilation

Analysis - Breaks up source into constituent pieces checks for syntactical and semantic errors

Synthesis - Constructs the target program and the symbol table

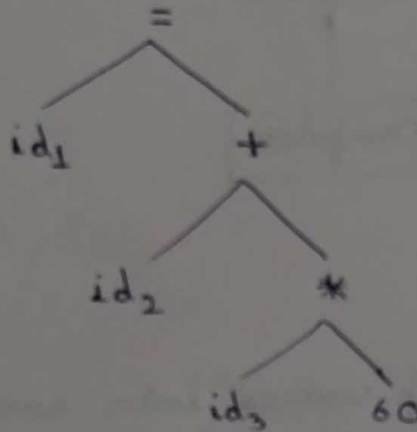
- The source program moves through a number of phases during compilation.
- Lexical Analyser takes the source program as input and produce a ~~large tree~~ long string of tokens.
- Syntax analyser takes the output from lexical analyser and produce a large tree.
- Semantic analyser takes the previous output and produces another tree.
- Similarly, Intermediate code generator takes a tree as an input and produces intermediate code.

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

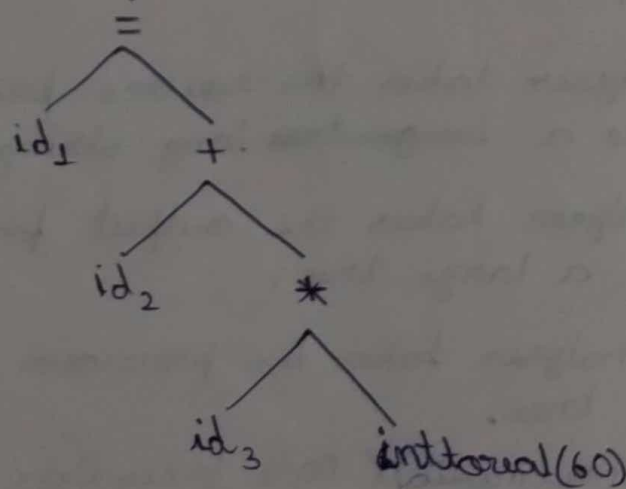
Lexical
Analyser

$$\text{id}_1 = \text{id}_2 + \text{id}_3 * 60$$

Syntax
Analyser



Semantic
Analyser



Intermediate code
generator

↓
 $temp1 = \text{intto real}(60)$
 $temp2 = id_3 * temp1$
 $temp3 = id_2 + temp2$
 $id_1 = temp3$

↓

code optimisation

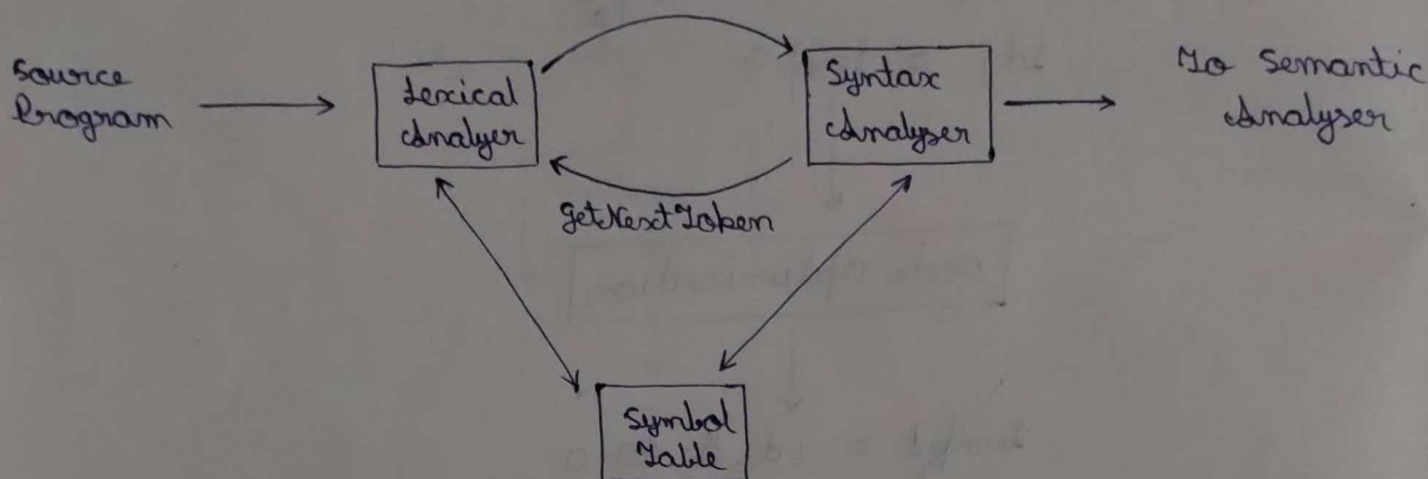
↓
 $temp1 = id_3 * 60.0$
 $id_1 = id_2 + temp1$

Textbook:

Alfred V. Aho, Ravi Sethi, Jeffrey D. Ullman

"compilers Principles Techniques and Tools"

Lexical analyser



→ Reads input characters and produces a series of tokens as output

Token : It is a pair consisting of a token name and an optional attribute value.

Pattern : The set of strings is described by a rule called pattern associated with that token.

Lexeme : It is a sequence of characters in the source program that matches the pattern for a token and is identified by a lexical analyser as an instance of a pattern.

- Tokens can be specified using :

(i) Regular expression

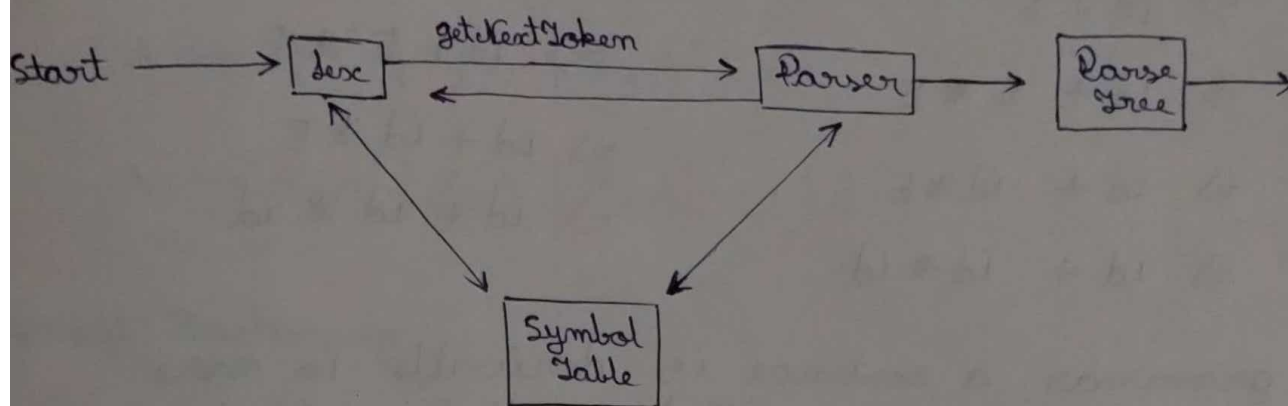
(ii) DFA

(iii) NDFA

(iv) NDFA with empty transitions

- We have to draw transition diagrams and then write programs for them.

Syntax analyzer



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

$$E \Rightarrow -E$$

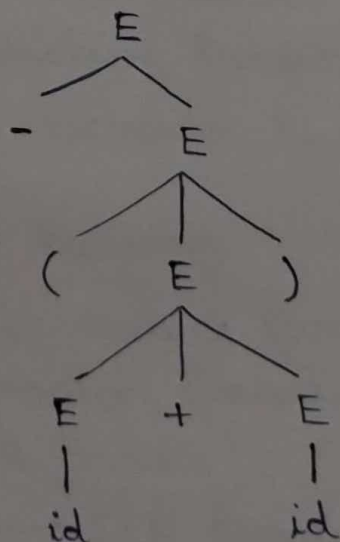
$$\Rightarrow -(E)$$

$$\Rightarrow -(E + E)$$

$$\Rightarrow -(id + E)$$

$$\Rightarrow -(id + id)$$

So, $-(id + id)$ can be derived from E .



Q. Derive $id + id * id$

Ans. : $E \Rightarrow E + E$

$$\Rightarrow id + E$$

$$\Rightarrow id + E * E$$

$$\Rightarrow id + id * E$$

$$\Rightarrow id + id * id$$

$$E \Rightarrow E * E$$

$$\Rightarrow E + E * E$$

$$\Rightarrow id + E * E$$

$$\Rightarrow id + id * E$$

$$\Rightarrow id + id * id$$

If by a grammar, a sentence is derivable in more than one way, then it will be called an 'ambiguous grammar'.

- Elimination of Left Recursion

$$A \rightarrow AX \mid \beta$$

This will generate left recursion.

This rule can be rewritten as:

$$A \rightarrow \beta A'$$

$$A' \rightarrow X A' \mid \epsilon$$

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

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

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

Remove left recursions.

$$\Rightarrow E \rightarrow TE'$$

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

$$T \rightarrow \bullet \bullet FT'$$

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

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

If a production rule is like,

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

then, after removing left recursion, it will look like,

$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \beta_3 A' \dots$$

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

- Left Factoring

If $A \rightarrow \alpha \beta_1 \mid \alpha \beta_2$ is the production rule, left factoring says that,

$$A \rightarrow \alpha A'$$

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

This is done when we cannot choose what production rule to use, for derivation.

① Approaches in parsing -

(a) Top-down Parsing

- Recursive Descent Parsing
- Predictive Parsing
- Non-recursive Predictive Parsing

(b) Bottom-up Parsing

- Shift-reduce Parsing
- Operator Precedence Parsing
- LR Parser
 - SLR Parsing Table
 - Look-ahead LR Parsing Table
 - Canonical LR Parsing Table

⊙ Top - Down Parsing

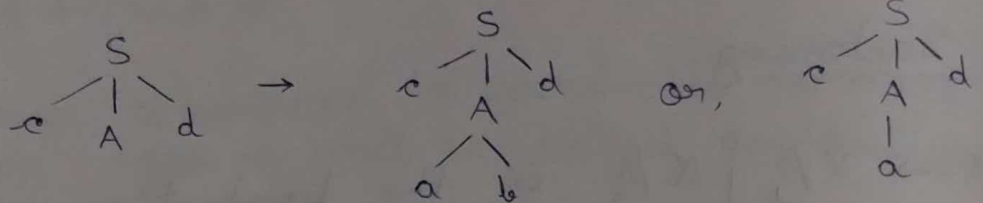
(i) Recursive Descent Parsing

Given $S \rightarrow cAd$

$A \rightarrow ab|a$

Input string, $w = cad$

Derivation:



(ii) Predictive Parsing

Given $S \rightarrow cAd$

$A \rightarrow ab|a$

Input string, $w = cad$

- In predictive parsing, there is no back-tracking.
- Modifying the production rules,

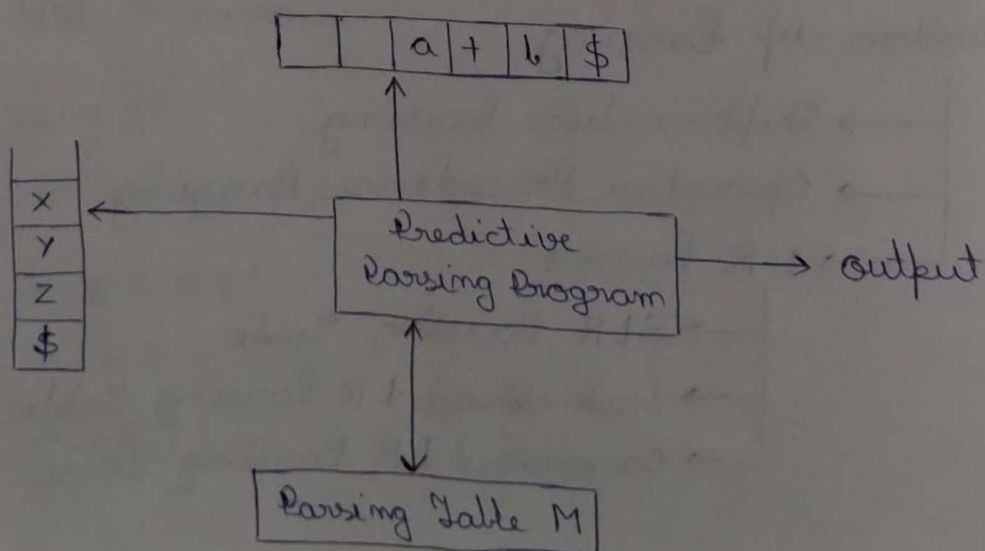
$S \rightarrow cAd$

$A \rightarrow aA'$

$A' \rightarrow b|\epsilon$

This way there is only one good way to generate string.

(iii) Non-recursive Predictive Parsing



The Parsing Table will help in choosing the production rule whereas the output will say if the string is accepted or not.

The program should implement,

- (a) if $X = a = \$$, the parser halts and announces successful completion of parsing
- (b) if $X = a \neq \$$, the parser pops-off the stack and advances the input pointer to the next input symbol
- (c) If X is a non-terminal symbol, then program consults the Parsing Table

Example :

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

Forming the parsing table :

The row-headers are non-terminal symbols whereas the column-headers are terminal symbols and right-end marker.

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

Stack	Input Buffer	Outputs
\$E	id + id * id \$	
\$E'T	id + id * id \$	$E \rightarrow TE'$
\$E'T'F	id + id * id \$	$T \rightarrow FT'$
\$E'T'id	id + id * id \$	$F \rightarrow id$
\$E'T'	+ id * id \$	
\$E'	+ id * id \$	$T' \rightarrow \epsilon$
\$E'T+	+ id * id \$	$E' \rightarrow +TE'$
\$E'T	id * id \$	
\$E'T'F	id * id \$	$T \rightarrow FT'$
\$E'T'id	id * id \$	$F \rightarrow id$
\$E'T'	* id \$	
\$E'T'F*	* id \$	$T' \rightarrow *FT'$
\$E'T'F	id \$	
\$E'T'id	id \$	$F \rightarrow id$
\$E'T'	\$	
\$E'	\$	$T' \rightarrow \epsilon$
\$	\$	$E' \rightarrow \epsilon$

- construction of 'Parsing Table' for non-recursive predictive parser :

Two functions are needed,

$FIRST(X)$
↓
Set of Terminals

$FOLLOW(A)$
↓
Set of Terminals

Rules for $FIRST()$:

1. If X is a terminal symbol, then $FIRST(X)$ is $\{X\}$
2. If $X \rightarrow \epsilon$ is a production rule, then add ϵ to $FIRST(X)$
3. If X is a non-terminal symbol and $X \rightarrow Y_1 Y_2 Y_3 \dots$ is a production rule, then place $FIRST(Y_1)$ to $FIRST(X)$.

Example :

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

$$\begin{aligned} FIRST(E) &= FIRST(TE') \\ &= FIRST(T) \\ &= FIRST(FT') \\ &= FIRST(F) \\ &= \{ (, id \} \end{aligned}$$

$$\therefore FIRST(E) \equiv FIRST(T) \equiv FIRST(F) = \{ (, id \}$$

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

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

Rules for FOLLOW():

1. Place \$ in FOLLOW(S) where S is the start symbol and \$ is the input right-end marker
2. If there is a production rule $A \rightarrow \alpha B \beta$, then everything in FIRST(β) except for ϵ , is placed in FOLLOW(B)
3. If there is a production $A \rightarrow \alpha B \dots$ or a production rule $A \rightarrow \alpha B \beta$ where FIRST(β) contains ϵ (i.e. $\beta \xrightarrow{*} \epsilon$), then everything in FOLLOW(A) is in FOLLOW(B)

In the previous example,

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

considering $F \rightarrow (E)$

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

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

FOLLOW(E') has to be found out.

$$E \rightarrow \overline{\alpha} \overline{B} \overline{\beta} \begin{matrix} T E' \\ \overline{\alpha} \overline{B} \overline{\beta} \end{matrix}$$

$$\text{By rule 3, } \text{FOLLOW}(E') = \{ \$, () \}$$

FOLLOW(T) has to be found out.

$$E \rightarrow \overline{\alpha} \overline{B} \overline{\beta} \begin{matrix} T E' \\ \overline{\alpha} \overline{B} \overline{\beta} \end{matrix}$$

FIRST(β) has empty symbol.

So, FOLLOW(E) will be in FOLLOW(T)

$$\text{Now, } E' \rightarrow + \overline{\alpha} \overline{B} \overline{\beta} \begin{matrix} T E' \\ \overline{\alpha} \overline{B} \overline{\beta} \end{matrix}$$

By rule 2, '+' will be in FOLLOW(T)

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

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

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

- LL(1) grammar

A grammar is said to be LL(1) if the parsing table does not contain multiple entries in one cell.

Example 1: $S \rightarrow i E + S S' \mid a$

$$S' \rightarrow \epsilon S \mid \epsilon$$

$$E \rightarrow b$$

Example 2: $E \rightarrow E + T \mid T$

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

$$F \rightarrow F * \mid a \mid b$$

- extension for input is '.l'.

definitions

% %

rules

% %

user code

% { <something> % } for verbatim copying

o Bottom-up Parsing

We start from the last symbol, and if we can reach the start one, the string is accepted.

(i) Shift Reduce Parsing (It's a model)

Q. Given, $S \rightarrow aABe$

$A \rightarrow Abc/b$

$B \rightarrow d$

check if 'abbcde' is acceptable.

Method: (i) Choose a sub-string which matches the right-hand side of a production rule. This will be called 'Handle'.

(ii) Keep doing it

a**b**bcde

a**Abc**de $A \rightarrow bc$ is the handle

aA**d**e $A \rightarrow Abc$ is the handle

aABe $B \rightarrow d$ is the handle

S $S \rightarrow aABe$ is the handle

Handle: Sub-string that matches the right side of a production rule and reduction is done through the non-terminal.

Handle Pruning: A right most derivation in reverse is known as handle pruning.

Q. Given, $E \rightarrow E + E$

$E \rightarrow E * E$

$E \rightarrow (E)$

$E \rightarrow id$

The input sentence is
 $id_1 + id_2 * id_3$.

Right Sentential 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	$E \rightarrow E + E$
E		

→ Stack Implementation of Shift-reduce Parser

Stack	Input Buffer	Action
\$	$id_1 + id_2 * id_3 \$$	Shift
$\$ id_1$	$+ id_2 * id_3 \$$	Reduce by $E \rightarrow id$
$\$ E$	$+ id_2 * id_3 \$$	Shift
$\$ E +$	$id_2 * id_3 \$$	Shift
$\$ E + id_2$	$* id_3 \$$	Reduce by $E \rightarrow id$
$\$ E + E$	$* id_3 \$$	Reduce by $E \rightarrow E + E$
$\$ E$	$* id_3 \$$	Shift
$\$ E *$	$id_3 \$$	Shift
$\$ E * id_3$	$\$$	Reduce by $E \rightarrow id$
$\$ E * E$	$\$$	Reduce by $E \rightarrow E * E$
$\$ E$	$\$$	Accepted

This is done manually, so reduction can be done like this or also after getting $\$ E + E * E$ in the stack. Although the later is preferable.

There are 4 actions here -

- (a) Shift
- (b) Reduce
- (c) Accept
- (d) Error

(ii) Operator Precedence Parsing

→ operator grammar

- (a) There will be no empty production at right
- (b) Has no production right side with two adjacent non-terminals

eg: $E \rightarrow EAE \mid (E) \mid -E \mid id$
 $A \rightarrow + \mid - \mid *$

This is not operator grammar.

But,

$E \rightarrow E + E \mid E - E \mid E * E \mid (E) \mid -E \mid id$
is an operator grammar.

We will use some Precedence rules here.

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

These three disjoint relations are required to form the parsing table.

on the table, row and column headers are there, which will be same, only ~~non~~ terminals.

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

formation is problem specific, that is, it's heuristically chosen.

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

Take two pointers and mark the precedence.

Then remove all non-terminals.

$\$ E + E * E \$$

\Downarrow

$\$ < + < * > \$$

So, we choose $E * E$ as handle and reduce it.

$\$ E + E \$$

Remove non-terminals again.

$\$ < + > \$$

So, we choose $E + E$ as handle and reduce it.

$\$ E \$$

(Accepted)

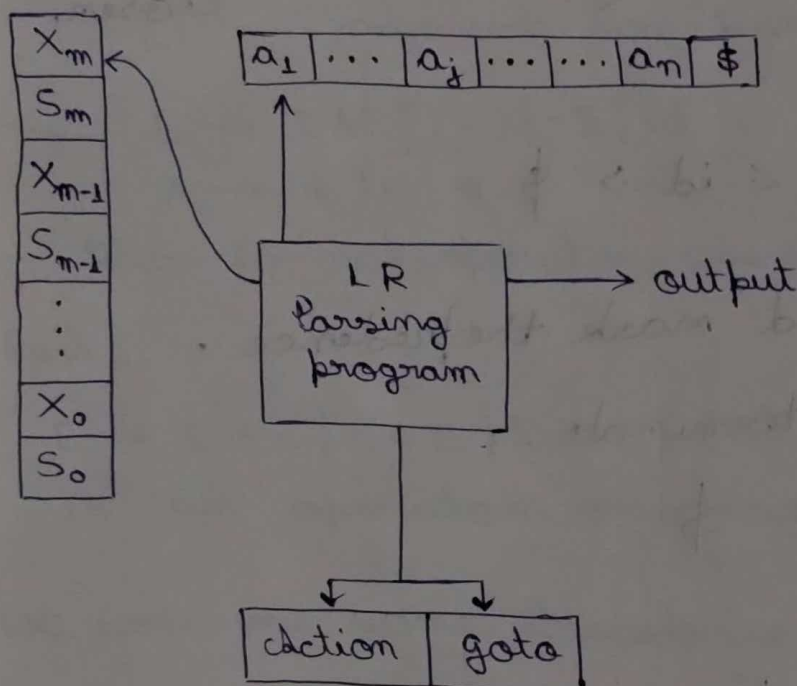
(iii) LR Parsing

LR(K) Parsing is a general model. It will be followed by all types of LR parsers.

L \rightarrow left to right scanning of the input

R \rightarrow for constructing a right most derivation in reverse

K \rightarrow the number of input symbols of lookahead that are used in making parsing decision



Whenever we have a CFG, we can apply LR(K) parsing.

There are some actions -

1. Shift [K number of symbols are pushed to stack]
2. Reduce by a grammar production rule
3. accept
4. Error

Example : $E \rightarrow E + T \mid T$
 $T \rightarrow T * F \mid F$
 $F \rightarrow (E) \mid id$

state	action						data		
	id	+	*	()	\$	E	T	F
0	S ₅			S ₄			1	2	3
1	S	S ₆				acc			
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			
6	S ₅			S ₄				9	3
7	S ₅			S ₄					10
8		S ₆			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			

Fig: Parsing Table

Stack	Input	Action
0	id + id * id \$	Shift
0 id 5	+ id * id \$	Reduce by $F \rightarrow id$
0 F 3	+ id * id \$	Reduce by $T \rightarrow F$
0 T 2	+ id * id \$	Reduce by $E \rightarrow T$
0 E 1	+ id * id \$	Shift
0 E 1 + 6	id * id \$	Shift
0 E 1 + 6 id 5	* id \$	Reduce by $F \rightarrow id$
0 E 1 + 6 F 3	* id \$	Reduce by $T \rightarrow F$
0 E 1 + 6 T 9	* id \$	Shift
0 E 1 + 6 T 9 * 7	id \$	Shift
0 E 1 + 6 T 9 * 7 id 5	\$	Reduce by $F \rightarrow id$
0 E 1 + 6 T 9 * 7 F 10	\$	Reduce by $T \rightarrow T * F$
0 E 1 + 6 T 9	\$	Reduce by $E \rightarrow E + T$
0 E 1	\$	Accept

When we reduce by a rule, twice the no. of symbols on the right side of the rule has to be popped from stack.

→ LR Grammar

Constructing the SLR Parsing Table :-

First, we have to find the LR(0) items.

LR(0) items → a dot at some position of the RHS of production rules

For $A \rightarrow XYZ$, LR(0) items are

$$A \rightarrow \cdot XYZ$$

$$A \rightarrow X \cdot YZ$$

$$A \rightarrow XY \cdot Z$$

$$A \rightarrow XYZ \cdot$$

we have to perform 'closure' and 'goto' actions on the Canonical LR(0) items.

we have to introduce a new start symbol S' .

this is called 'Augmenting' the grammar.

If there are many S productions like,

$$S \rightarrow X_1$$

$$S \rightarrow X_2 \text{ and so on,}$$

it will be difficult to know when to stop. So, we augment the grammar and parser will stop when $S' \rightarrow S$ is encountered.

eg. : $E \rightarrow E + T \mid T$

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

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

Now, we augment the grammar,

$$S' \rightarrow E$$

$$I = \{[S' \rightarrow \cdot E]\}$$

$$\text{Closure}(I) = \{[S' \rightarrow \cdot E], [E \rightarrow \cdot E + T], [E \rightarrow \cdot T], [T \rightarrow \cdot T * F], [T \rightarrow \cdot F], [F \rightarrow \cdot (E)], [F \rightarrow \cdot id]\}$$

Take production rules whenever new non-terminals appear in right-hand side.

$$I_0 = \text{closure}(I)$$

Now, we have to find $\text{goto}(I_0, E)$,
consider all productions where 'E' is on RHS.

$$\begin{aligned} I_1 &= \text{goto}(I_0, E) \\ &= \{ [S' \rightarrow E \cdot], [E \rightarrow E \cdot + T] \} \end{aligned}$$

$$\begin{aligned} \text{then } I_2 &= \text{goto}(I_0, T) \\ &= \{ [E \rightarrow T \cdot], [T \rightarrow T \cdot * F] \} \end{aligned}$$

$$\begin{aligned} \text{then } I_3 &= \text{goto}(I_0, F) \\ &= \{ [T \rightarrow F \cdot] \} \end{aligned}$$

$$\begin{aligned} \text{then } I_4 &= \text{goto}(I_0, () \\ &= \{ [F \rightarrow (\cdot E)] \} \cup \text{closure}(E) \end{aligned}$$

since E is non terminal.

$$\therefore I_4 = \{ [F \rightarrow (\cdot E)], [E \rightarrow \cdot E + T], [E \rightarrow \cdot T], [T \rightarrow \cdot T * F], [T \rightarrow \cdot F], [F \rightarrow \cdot (E)], [F \rightarrow \cdot \text{id}] \}$$

$$\begin{aligned} \text{then } I_5 &= \text{goto}(I_0, \text{id}) \\ &= \{ [F \rightarrow \text{id} \cdot] \} \end{aligned}$$

$$\begin{aligned} \text{then } I_6 &= \text{goto}(I_1, +) = \{ [E \rightarrow E + \cdot T] \} \cup \text{closure}(T) \\ &= \{ [E \rightarrow E + \cdot T], [T \rightarrow \cdot T * F], [T \rightarrow \cdot F], [F \rightarrow \cdot (E)], [F \rightarrow \cdot \text{id}] \} \end{aligned}$$

$$\begin{aligned} \text{then } I_7 &= \text{goto}(I_2, *) \\ &= \{ [T \rightarrow T * \cdot F] \} \cup \text{closure}(F) \\ &= \{ [T \rightarrow T * \cdot F], [F \rightarrow \cdot (E)], [F \rightarrow \cdot \text{id}] \} \end{aligned}$$

$$\text{then } \text{goto}(I_4, E) = I_8$$

$$= \{[F \rightarrow (E \cdot)], [E \rightarrow E \cdot + T]\}$$

$$\text{then } I_9 = \text{goto}(I_4, T)$$

$$= \{[E \rightarrow T \cdot], [T \rightarrow T \cdot * F]\} = I_2$$

$$\text{then } I_{10} = \text{goto}(I_4, F)$$

$$= \{[T \rightarrow F \cdot]\} = I_3$$

$$\text{then } I_{11} = \text{goto}(I_4, ($$

$$= \{[F \rightarrow (\cdot E)]\} \cup \text{closure}(E)$$

$$= I_4$$

$$\text{then } I_{12} = \text{goto}(I_4, \text{id})$$

$$= \{[F \rightarrow \text{id} \cdot]\} = I_5$$

$$\text{then } I_{13} = \text{goto}(I_6, T)$$

$$= \{[E \rightarrow E + T \cdot], [T \rightarrow T \cdot * F]\}$$

$$\text{then } I_{14} = \text{goto}(I_6, F)$$

$$= \{[T \rightarrow F \cdot]\} = I_3 = I_{10}$$

$$\text{then } I_{15} = \text{goto}(I_6, ($$

$$= \{[F \rightarrow (\cdot E)]\} \cup \text{closure}(E)$$

$$= I_4 = I_{11}$$

$$\text{then } I_{16} = \text{goto}(I_6, \text{id}) = \{[F \rightarrow \text{id} \cdot]\} = I_5 = I_{12}$$

$$\text{then } I_{17} = \text{goto}(I_7, F)$$

$$= \{[T \rightarrow T * F \cdot]\}$$

$$\begin{aligned}
 \text{Then, } I_{18} &= \text{goto}(I_7, () \\
 &= \{[F \rightarrow (\cdot E)]\} \cup \text{closure}(E) \\
 &= I_4 = I_{11} = I_{15}
 \end{aligned}$$

$$\begin{aligned}
 \text{Then } I_{19} &= \text{goto}(I_7, \text{id}) \\
 &= \{[F \rightarrow \text{id} \cdot]\} = I_5 = I_{12}
 \end{aligned}$$

$$\begin{aligned}
 \text{Then, } I_{20} &= \text{goto}(I_8,) \\
 &= \{[F \rightarrow (E) \cdot]\}
 \end{aligned}$$

$$\begin{aligned}
 \text{Then, } I_{21} &= \text{goto}(I_8, +) \\
 &= \{[E \rightarrow E + \cdot T]\} \cup \text{closure}(T)
 \end{aligned}$$