

---

# **Chapter 5**

# **Compilers: Analysis**

# **Phase**

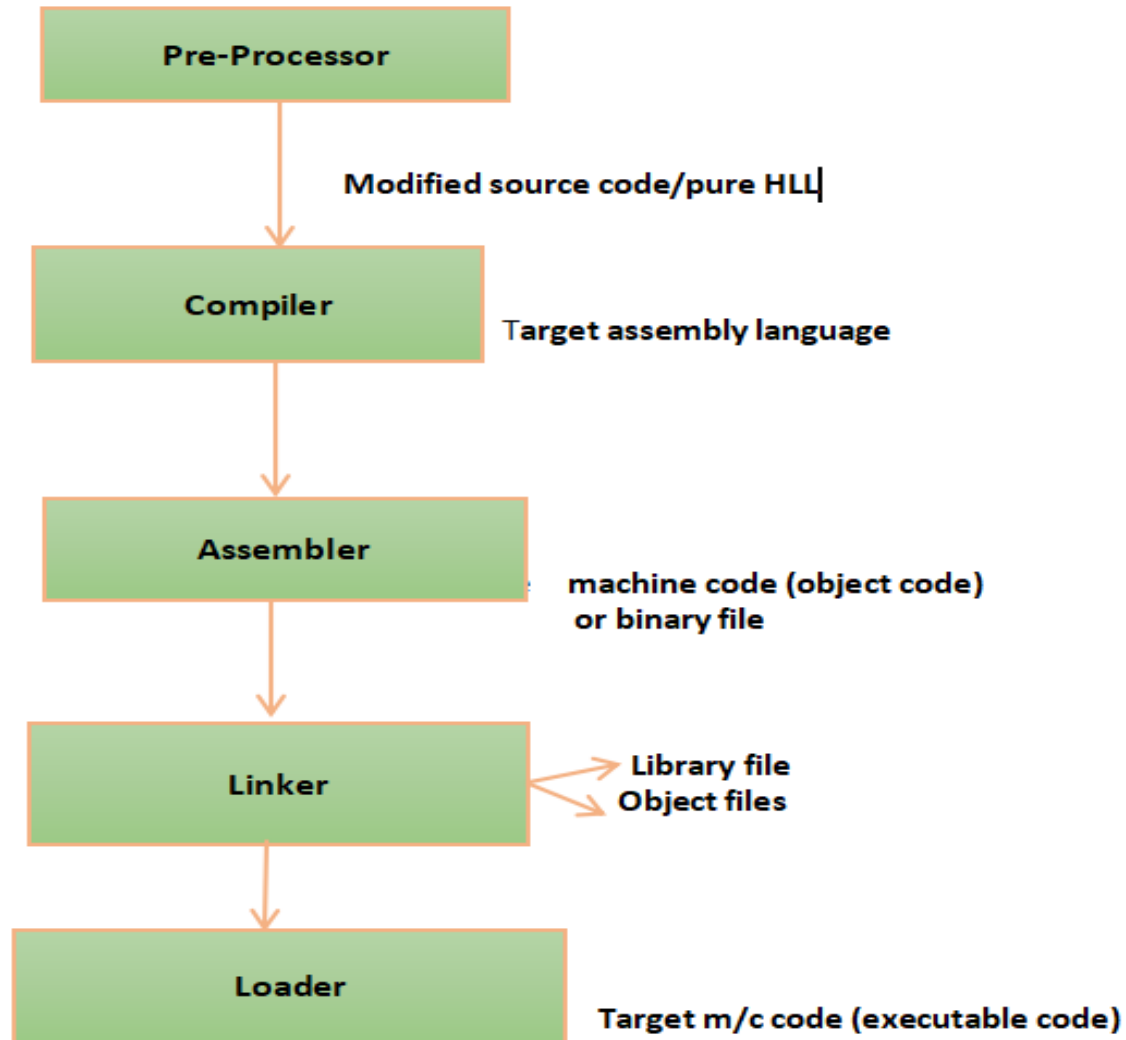
# Syllabus

---

	<b>Compilers: Analysis Phase</b>
1	Introduction to compilers
2	Phases of compilers: Lexical Analysis- Role of Finite State Automata in Lexical Analysis
3	Design of Lexical analyser, data structures
4	Syntax Analysis- Role of Context Free Grammar in Syntax analysis,
5	Types of Parsers:
6	Top down parser- LL(1)
7	Bottom up parser- Operator precedence parser,
8	SLR
9	Semantic Analysis
9	Syntax directed definitions

# Background

---

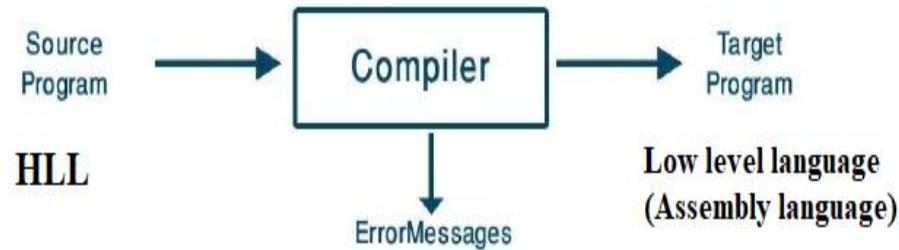


# Introduction to Compilers

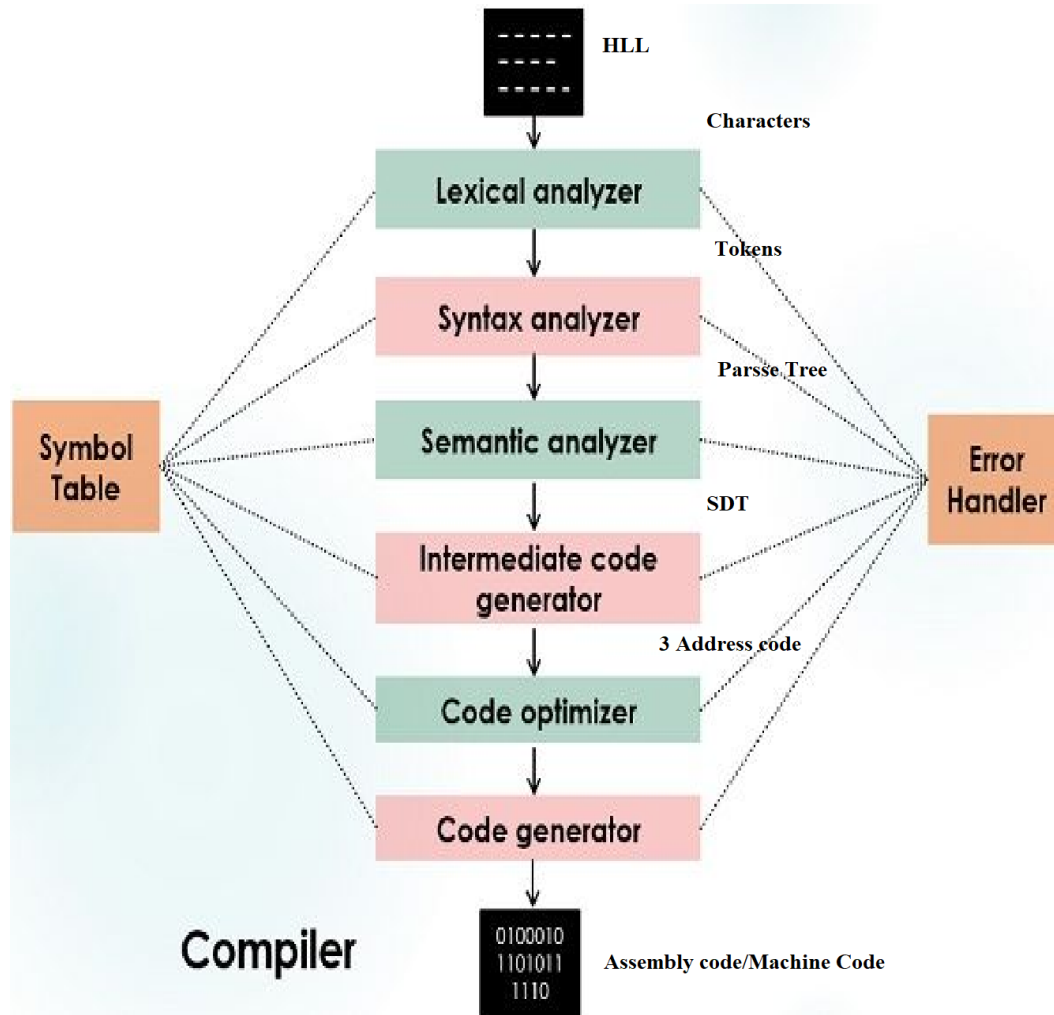
---

A compiler translates source code written in HLL to target code written in assembly language or low level language.

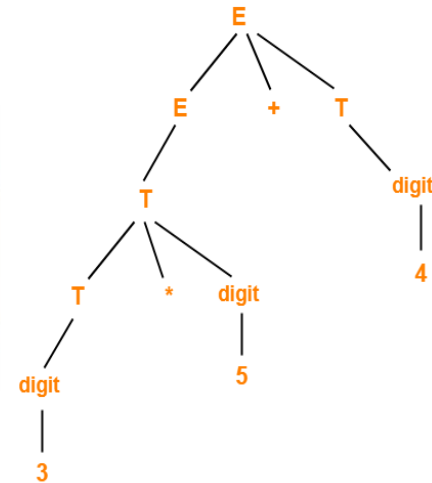
It reports the user about error(s) in source program if any.



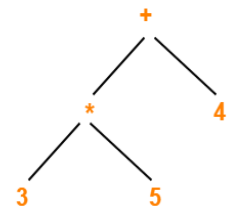
# Phases of Compiler



$$y = m * x + c$$



Parse Tree



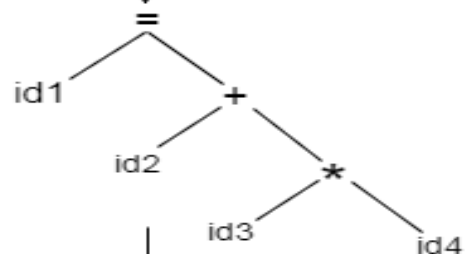
Syntax Tree

Sum:= Old sum + Rate \* 50

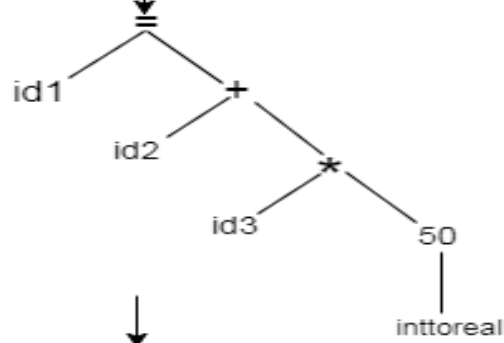
Lexical Analyzer

id1 = id2 + id3 \* id4

Syntax analyzer



Semantic analyzer



Intermediate code generator

```
temp1: = inttoreal(50)
temp2: = id3*temp1
temp3: = id2*temp2
id1: = temp3
```

Code optimization

```
temp1: = id3* 50.0
id1: = id2 + temp1
```

Code generation

```
MOVF id3,R2
MULF #50.0,R2
MOVF id2,R1
ADDF R2,R1
MOVF R1,id1
```

# Lexical Analysis

---

LEXICAL ANALYSIS is the very first phase in the compiler designing.

It is also known as scanner

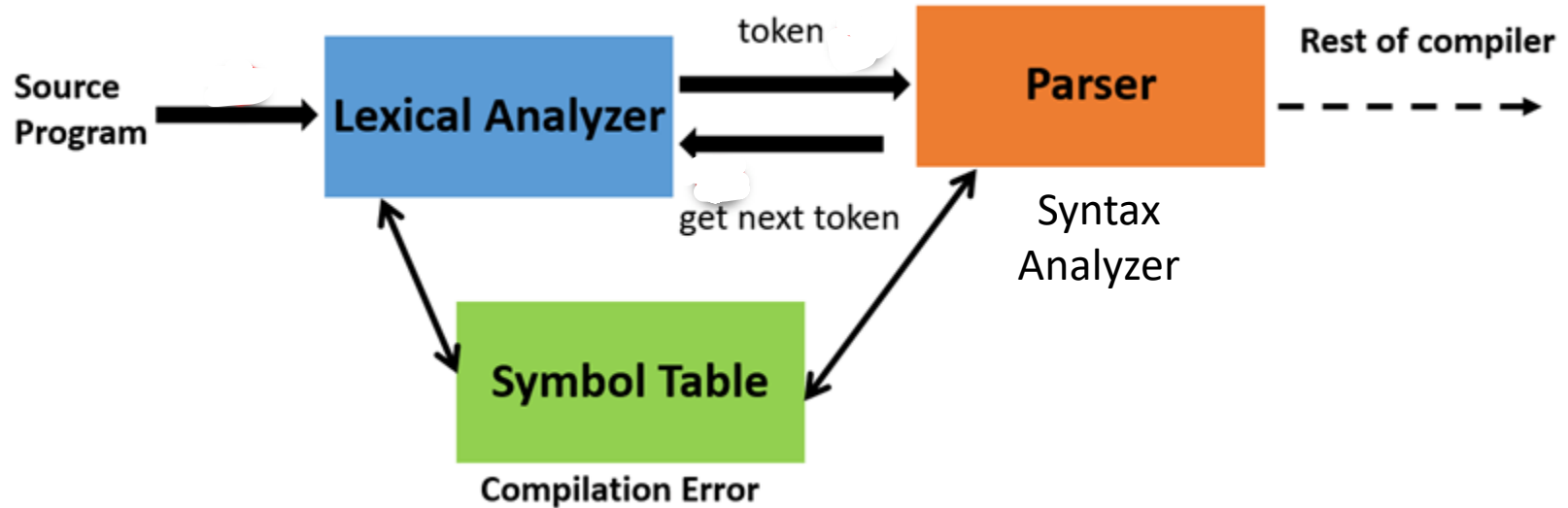
- keywords- lexer, tokenizer, scanner

## ➤ Functions:

- A Lexer reads the source code which is written in the form of sentences
- It converts a sequence of characters into a sequence of tokens (**tokenization**)
- It removes any extra space or comment written in the source code (//, #)
- If the lexical analyzer detects that the token is invalid, it generates an error- *lexical errors* – *spelling error, exceeding length of identifier, appearance of illegal characters. It is found during execution of program*
- It reads character streams from the source code, checks for legal tokens, and pass the data to the syntax analyzer when it demands.
- Helps to identify tokens into the symbol table

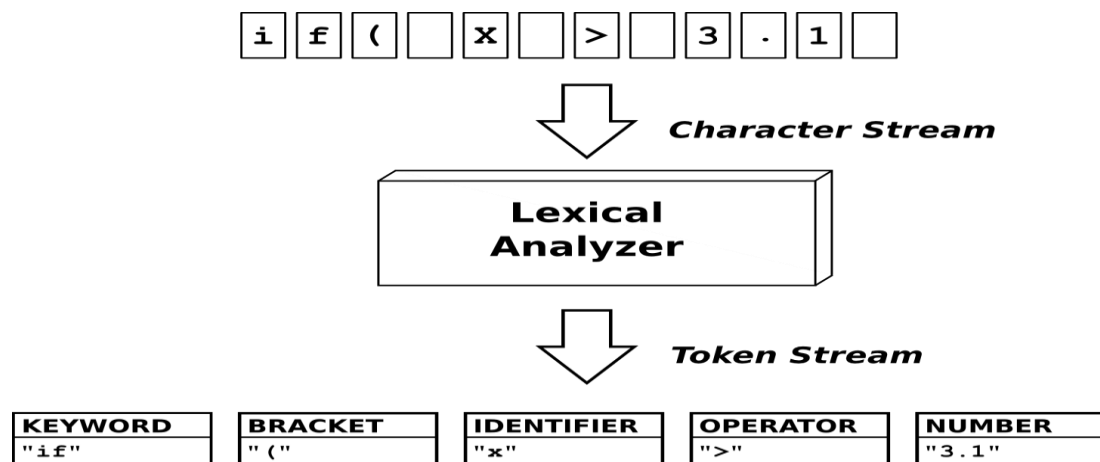
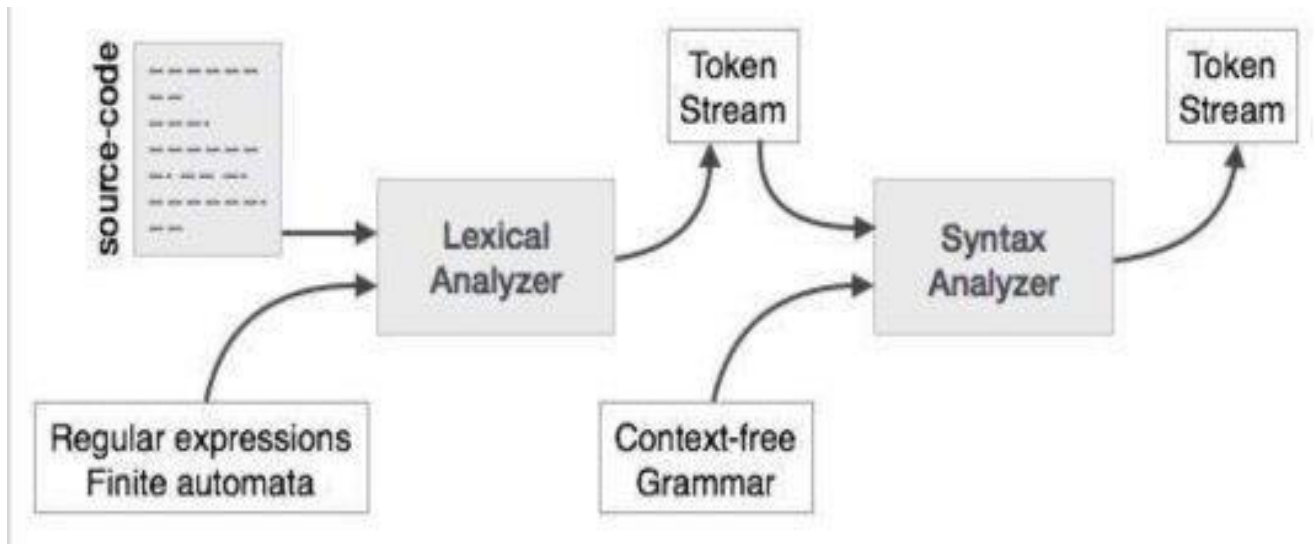
E.g.: *Can be used by web browser to format a webpage with parse data from javascript, html and css.*

# Lexical Analysis



- Lexical analyzer reads the input source program, scans the characters and produce a sequence of tokens that the parser can use for syntactic analysis.
- “get next token” is a command sent from the parser to the lexical analyzer





---

## **Language**

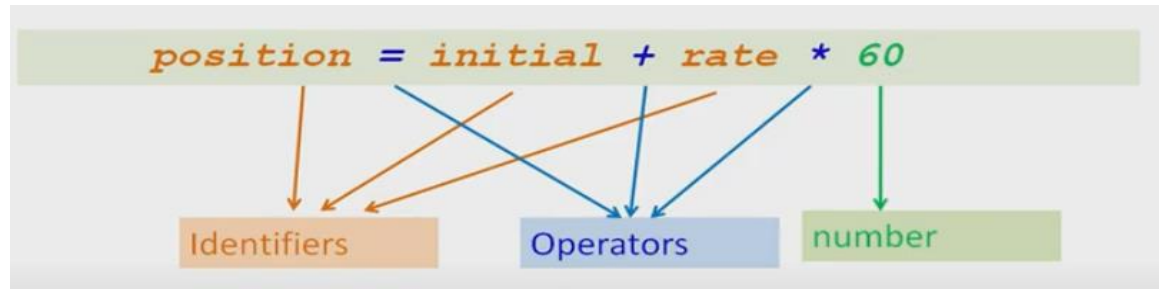
A language is considered as a finite set of strings over some finite set of alphabets. Computer languages are considered as finite sets, and mathematically set operations can be performed on them. Finite languages can be described by means of regular expressions.

### **Longest Match Rule**

When the lexical analyzer reads the source-code, it scans the code letter by letter; and when it encounters a whitespace, operator symbol, or special symbols, it decides that a word is completed.

`int intvalue;`

While scanning both lexemes till 'int', the lexical analyzer cannot determine whether it is a keyword `int` or the initials of identifier `int value`.



`position = initial + rate * 60`

gets translated as:

- The IDENTIFIER `position`
- The ASSIGNMENT SYMBOL `=`
- The IDENTIFIER `initial`
- The PLUS OPERATOR `+`
- The IDENTIFIER `rate`
- The MULTIPLICATION OPERATOR `*`
- The NUMERIC LITERAL `60`

# Recognition of Tokens

---

Complete set of tokens form the set of terminal symbols used in the grammar for the parser.

In most of the languages, the tokens fall into these categories:

- Keywords,
- operators,
- identifiers,
- constants,
- literal strings,
- punctuation

The regular definitions for tokens are as follows:

**if**  $\rightarrow$  if

**then**  $\rightarrow$  then

**else**  $\rightarrow$  else

**relop**  $\rightarrow$  < | <= | = | > | >=

**id**  $\rightarrow$  letter (letter|digit)\*

**num**  $\rightarrow$  digit<sup>+</sup> (.digit<sup>+</sup>)? (E(+|-)?digit<sup>+</sup>)?

**delim**  $\rightarrow$  blank | tab | newline

**ws**  $\rightarrow$  delim<sup>+</sup>

# FA and RE in lexical analyzer

---

Lexical analysis is process of recognizing tokens from the input.

Following are the steps:

1. Store the input in the input buffer
2. The token is read and regular expressions are built for corresponding token.
3. Regular expression is converted into a finite automata.
4. For each state of FA, a function is designed and each input along with transitional edges correspond to the input parameters of these functions.
5. The set of such functions ultimately creates lexical analyser program.

# Steps of lexical analysis

---

- Tokenization
- Give only lexeme related Error messages: Exceeding length, unmatched string, illegal characters. No msgs related to syntax, semantics
- Eliminate comments, white spaces (Tab, blank space, newline)

E.g. Find no. of tokens in the foll:

```
int a=20, b=30;
```

No. of tokens= 9

```
Printf("i=%d",i);
```

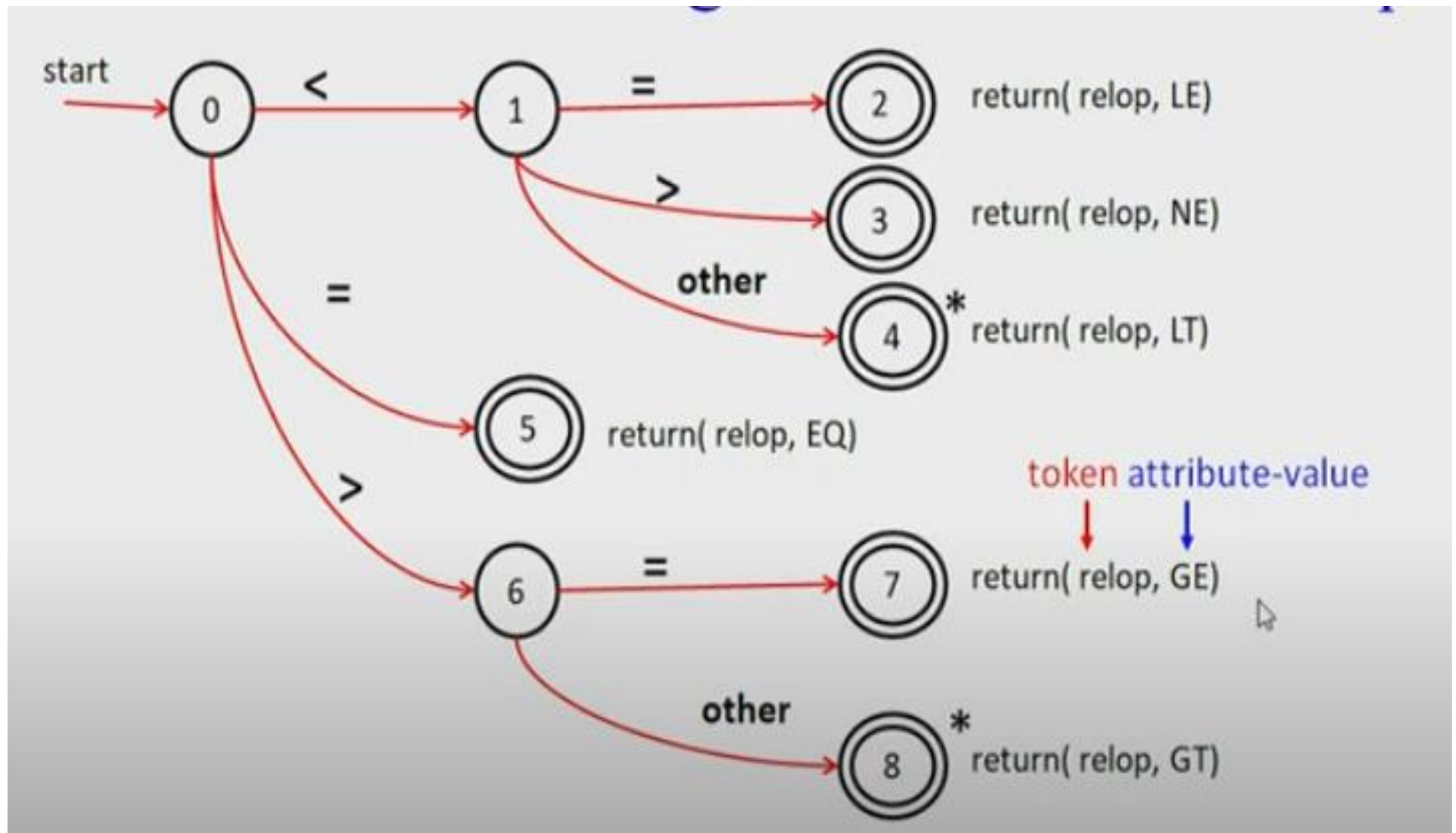
No. of tokens = 7

---

---

LEXEMES	TOKEN NAME	ATTRIBUTE VALUE
Any <i>ws</i>	–	–
<b>if</b>	<b>if</b>	–
<b>then</b>	<b>then</b>	–
<b>else</b>	<b>else</b>	–
Any <i>id</i>	<b>id</b>	Pointer to table entry
Any <i>number</i>	<b>number</b>	Pointer to table entry
<b>&lt;</b>	<b>relop</b>	<b>LT</b>
<b>&lt;=</b>	<b>relop</b>	<b>LE</b>
<b>=</b>	<b>relop</b>	<b>EQ</b>
<b>&lt;&gt;</b>	<b>relop</b>	<b>NE</b>
<b>&gt;</b>	<b>relop</b>	<b>GT</b>
<b>&gt;=</b>	<b>relop</b>	<b>GE</b>

# Transition diagram for relational operators

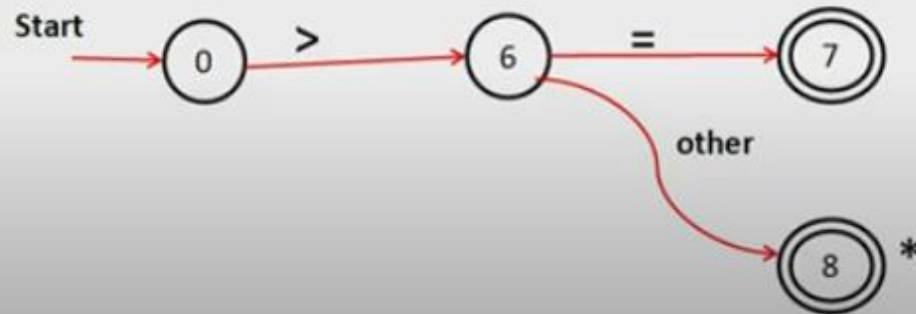




# Transition Diagram

---

- Lexical analysis uses transition diagram to keep track of information about characters that are seen as the forward pointer scans the input
- Transition diagrams are also called finite automata



# Transition diagram for identifiers and keywords



`gettoken( )`: returns token (**id**, **if**, **then**,...) if it looks the symbol table

`install_id( )`: return 0 if keyword or a pointer to the symbol table entry if **id**

Input String	<code>gettoken()</code>	<code>install_id()</code>
for	for	0
if	if	0
else	else	0
count	id	ptr to sym table entry for count
num4	id	ptr to sym table entry for num4

# Example

```
if x <= 10
then
    x
else
    y
```

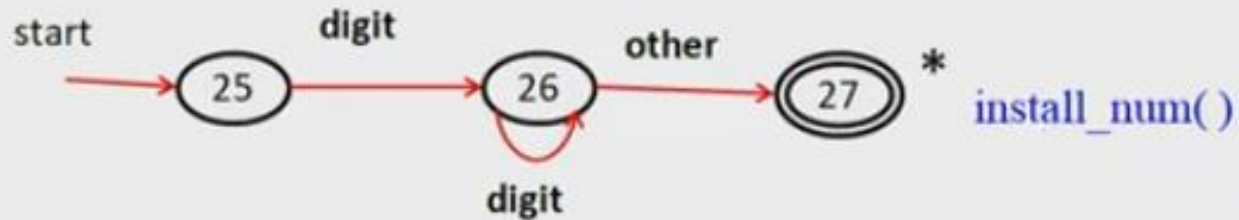
Source program

Lexical  
Analyzer

Token	Attribute
IF	NULL
ID	Ptr to x in ST
RELOP	LE
NUM	Ptr to 10 in ST
THEN	NULL
ID	Ptr to x in ST
ELSE	NULL
ID	Ptr to y in ST

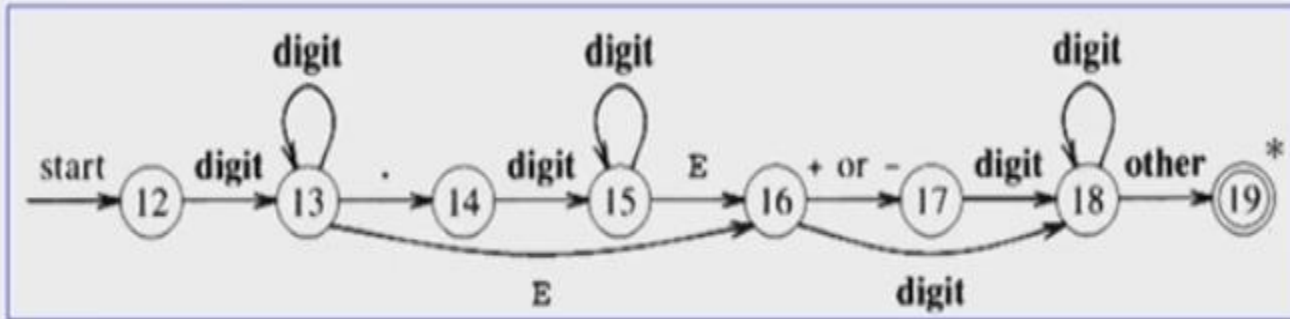
Tokens and Attributes

# Unsigned integer numeric constant



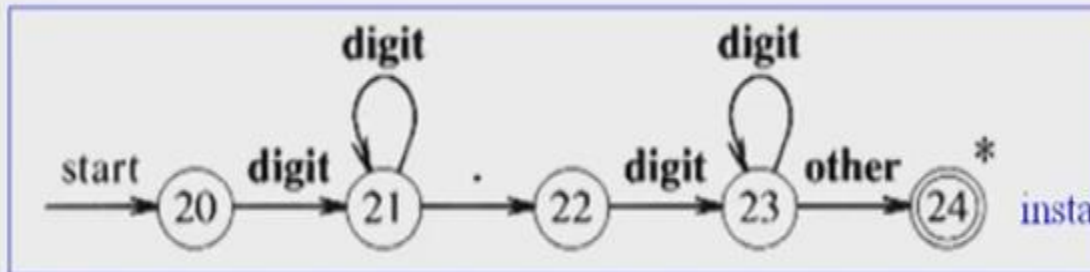
- It accepts the unsigned integer numeric constants
- Example- 314, 9, 41267 .....
- `install_num()`: Enters the number in a literal table and returns a pointer to that entry

# Unsigned Numbers



`install_num()`

e.g. 12.3E4, 10E15

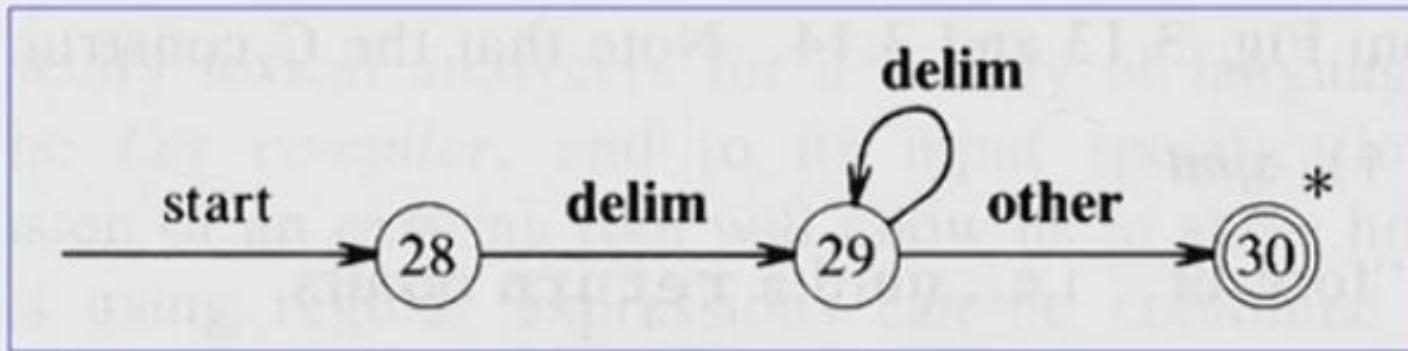


`install_num()`

e.g. 12.3, 10.15

# Whitespace

---



- It accepts the whitespaces (newline, tab space and space)
- Returns the token as **ws**

# Tokens, Lexemes and patterns

---

A **lexeme** is a sequence of characters that are included in the source program according to the matching pattern of a token. It is nothing but an instance of a token.

The **token** is a sequence of characters which represents a unit of information in the source program.

E.g. Identifier, keyword, constants (literals), operators, separators, special characters.

A **pattern** is a description which is used by the token. In the case of a keyword which uses as a token, the pattern is a sequence of characters.

Token	Sample Lexemes	Informal description of pattern
if	if	if
while	while	while
Relation	<, <=, =, >, >=	< or <= or = or > or >=
Id	count, sun, i, j, pi, D2	Letter followed by letters and digits
Num	0, 12, 3.1416, 6.02E23	Any numeric constant

# Tokens, Lexemes and patterns

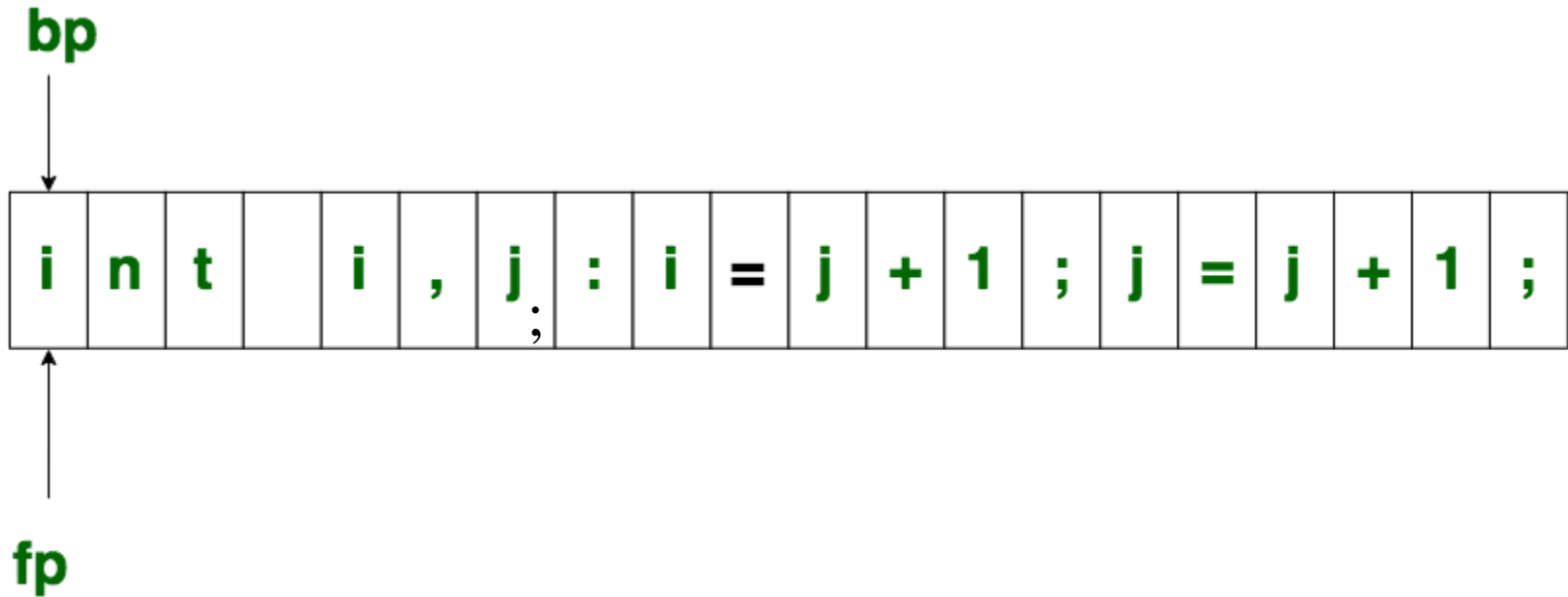
---

Token Type	Token Value	Informal Description
Integer constant	Numeric value	Numbers like 3, -5, 12 without decimal pts.
Floating constant	Numeric value	Numbers like 3.0, -5.1, 12.2456789
Reserved word	Word string	Words like if, then, class...(Keywords)
Identifiers	Symbol table index	Words not reserved starting with letter or _ and containing only letters, _, and digits
Relations	Operator string	<, <=, ==, ...
Operators	Operator string	=, +, -, ++, ...
Char constant	Char value	'A', ...
String	String	"this is a string", ...



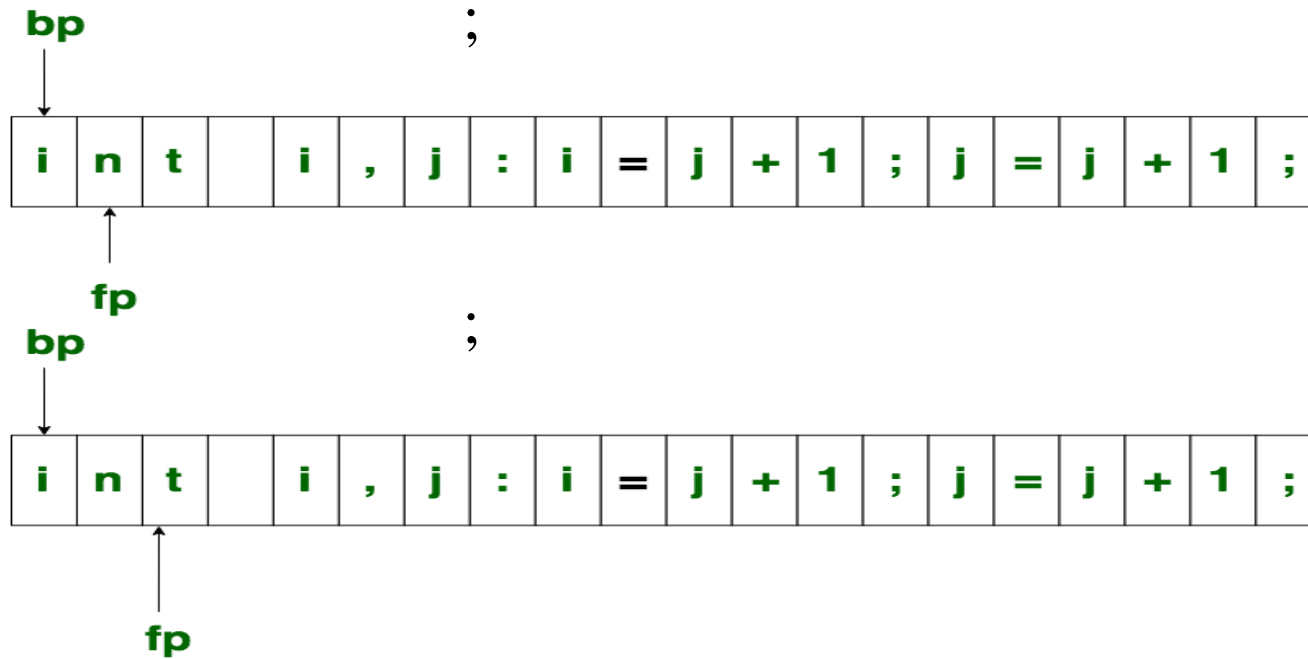
# Example

Ex. `int i,j; i=j+1; j=j+1;`



**Initial Configuration**

Initially both the pointers point to the first character of the input string as shown below



### Input Buffering

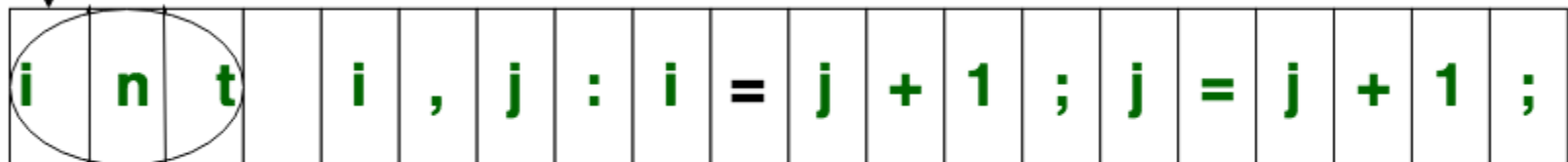
- The forward ptr moves ahead to search for end of lexeme. As soon as the blank space is encountered, it indicates end of lexeme.
- In above example as soon as ptr (fp) encounters a blank space the lexeme “int” is identified.

---

**Token**

**bp**

**;**



**fp**

**Input buffering**

# Practice problems

```
int main() {  
    x = y + z;  
    int x, y, z;  
    printf("Sum %.d %.d", x);  
}
```

No. of tokens= 26

```
main() {  
    int a = 10;  
    char b = "abc";  
    int c = 30;  
    char a[] = "xyz";  
    int /* Comment */ t = 40.5;  
}
```

No. of tokens= 33

```
main() {  
    a = b++ + -- -- ++ + ==;  
    printf("%.d %.d", a, b);  
}
```

No. of tokens= 25

# Data Structure-Symbol Table

---

The symbol table will contain the following types of information for the input strings in a source program:

- The lexeme (*input string*) itself
- Corresponding token
- Its semantic component (e.g., variable, operator, constant, functions, procedure, etc.)
- Data type
- Pointers to other entries (when necessary)

---

## Data stored in the Symbol Table

Let's take a look at the kinds of information stored in the symbol table for different types of elements within a program.

	Token	Symbol	Data		
<u>Lexeme</u>	<u>Class</u>	<u>Type</u>	<u>Type</u>	<u>Value</u>	<u>Scope</u>
while	TokWhile	Keyword	None	0	Global
+	TokPlus	Operator	None	0	Global
x	TokIdentifier	Variable	Integer	0	Sample
8	TokConstant	Literal	Integer	8	Global

---

## The basic operations of the Symbol Table

- There are several operations that are performed on the Symbol Table, the most important being:
- Adding symbols - The reserved words, standard identifiers and operators are placed in the Symbol Table during its initialization.
  - New lexemes are added when the scanner encounters them, and they are assigned a token class.
  - Similarly, the semantic analyzer adds the appropriate properties and attributes that belong to the lexeme.

---

There are several different ways to organize the symbol table:

The most obvious is to organize it as an *array of records*, but it would either require a linear search or constant sorting.

begin	tokbegin	stkeyword	dtnone	...	global
call	tokcall	stkeyword	dtnone	...	global
called	tokidentifier	stvariable	dtreal	...	sample
procedure	tokprocedure	stkeyword	dtnone	...	global
xmin	tokidentifier	tokconstant	dtinteger	...	sample



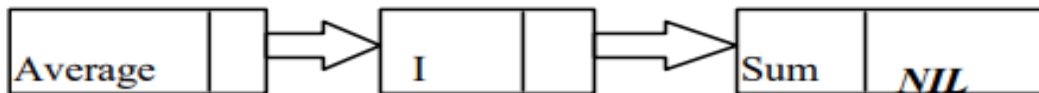
---

## Possible organization methods (continued)

We could use an unordered list. This would make efficient use of memory and it would be quick and easy to install names, but retrieval would be unacceptably slow.

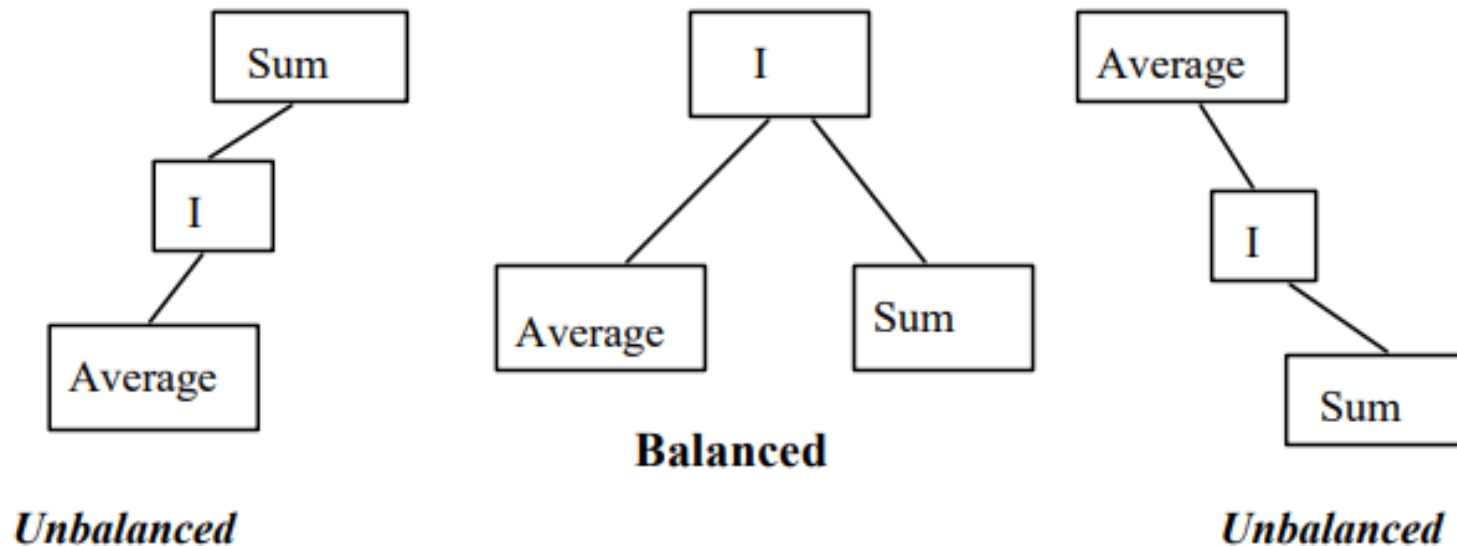


An ordered list would make storing and retrieving names faster, but not by that much.



---

**Binary search trees** would provide some compromises on these structures above, but only when the tree is balanced. This rarely happens on its own and takes time for balancing to be done by the insertion routine.



---

## Organization of the Symbol Table

We need a method for storing the lexemes in an efficient manner. Using an array of strings is not efficient and would limit the length of key words and identifiers.

1	c	a	l	l				
2	d	e	c	l	a	r	e	
3	a	v	e	r	a	g	e	s
4	x							

FORTRAN and other earlier programming languages restricted identifiers to a specific length for exactly this reason. (FORTRAN limited identifiers to 6 letters; C originally ignored anything after the first 8 letters.)

---

# The string table and name table

String table

b	e	g	i	n	c	a	l	l	d	e	c	l	a	r	e	
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	--

Name table

0	0	5	0	-1
1	5	4	1	-1
2	9	7	2	-1

Start  
position

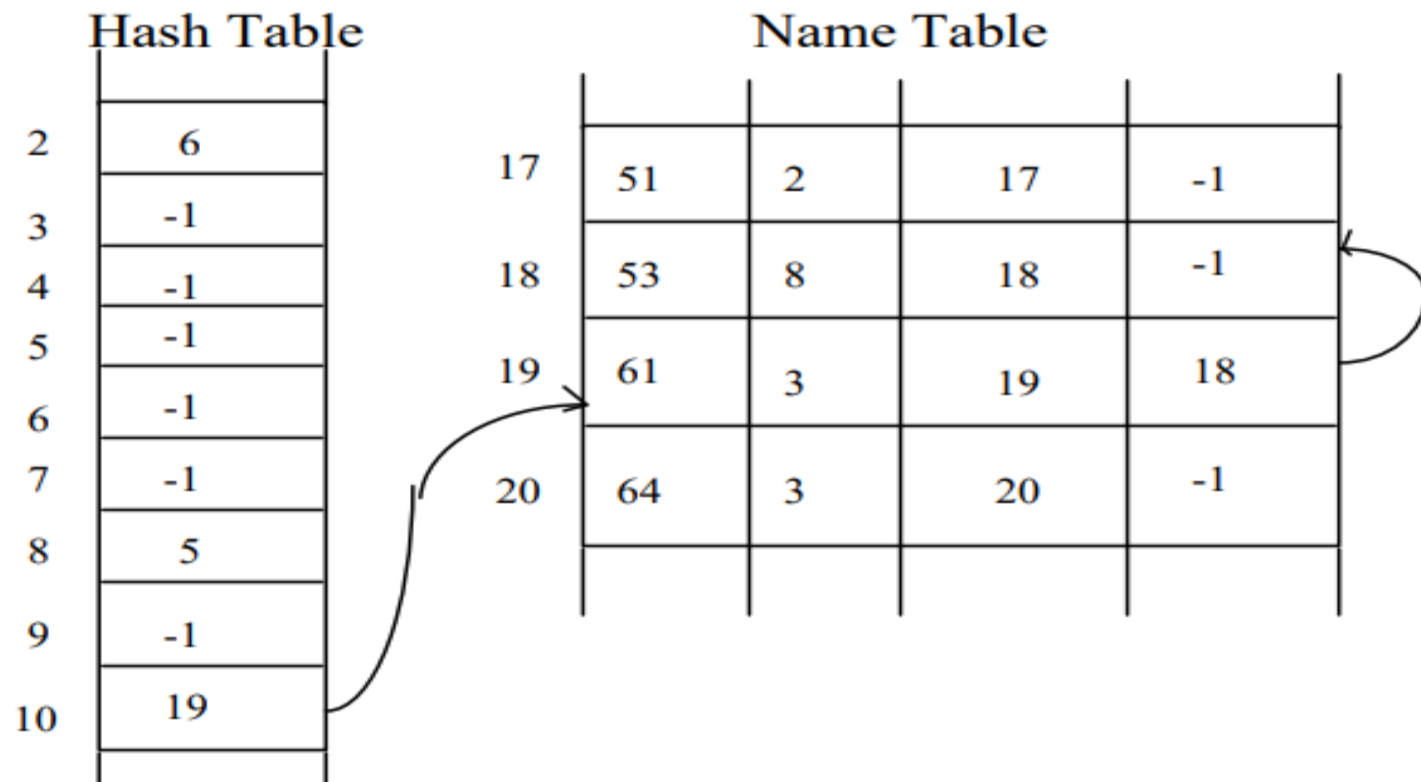
Length

Pointer to attribute  
table entry

Token

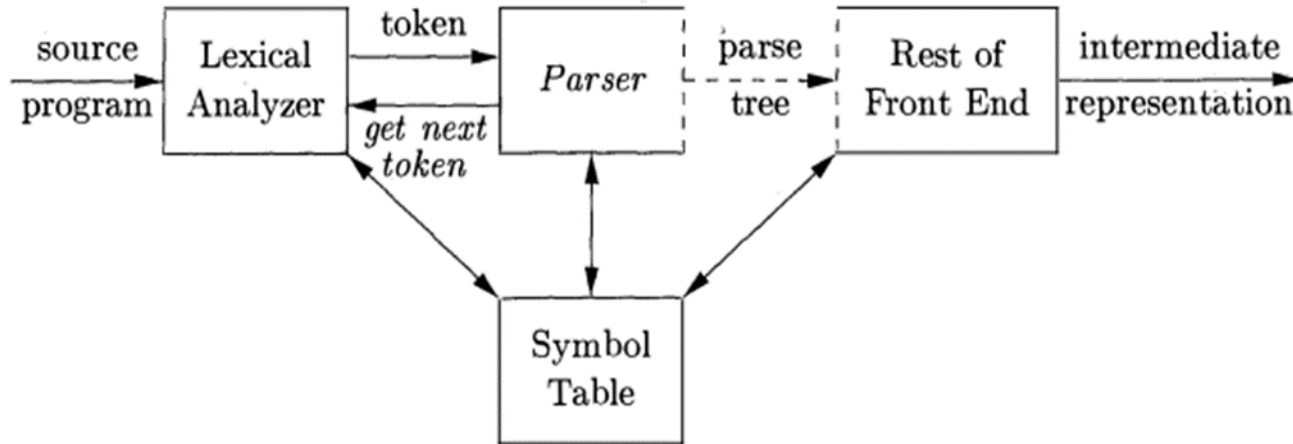
---

## Hash Table and the Name Table



# Syntax Analysis

Syntax Analyzer follows the production rules defined by the Context free grammar.



Context Free Grammar is defined as :  $G (V, T, P, S)$

$V$  is a set of non-terminal symbols.

$T$  is a set of terminals where  $V \cap T = \text{NULL}$ .

$P$  is a set of rules,  $P: V \rightarrow (V \cup T)^*$ , i.e., the left-hand side of the production rule  $P$  does not have any right context or left context.

$S$  is the start symbol.

# Syntax Analysis

---

**A derivation tree or parse tree** is an ordered rooted tree that graphically represents the semantic information a string derived from a context-free grammar.

## **Leftmost and Rightmost Derivation of a String**

**Leftmost derivation** – A leftmost derivation is obtained by applying production to the leftmost variable in each step.

**Rightmost derivation** – A rightmost derivation is obtained by applying production to the rightmost variable in each step.

Let any set of production rules in a CFG be

$X \rightarrow X+X \mid X*X \mid X \mid a$

over an alphabet  $\{a\}$ .

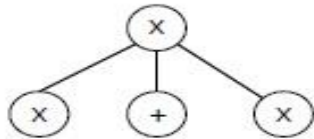
---

$$X \rightarrow X+X \mid X*X \mid X \mid a$$

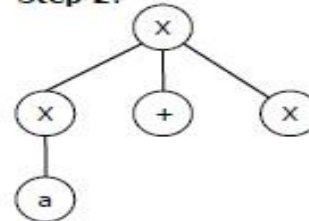
The leftmost derivation for the string "a+a\*a" may be –

$$X \rightarrow X+X \rightarrow a+X \rightarrow a + X*X \rightarrow a+a*X \rightarrow a+a*a$$

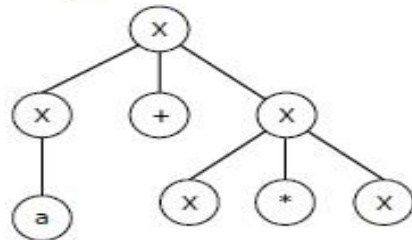
**Step 1:**



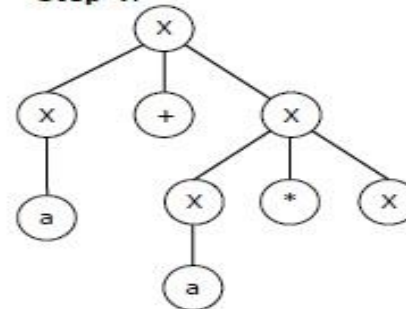
**Step 2:**



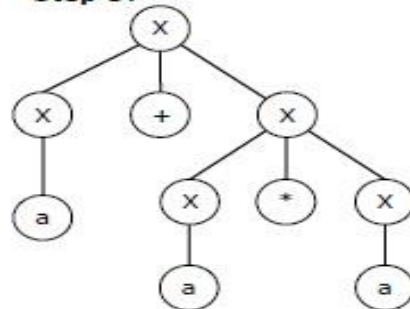
**Step 3:**



**Step 4:**



**Step 5:**





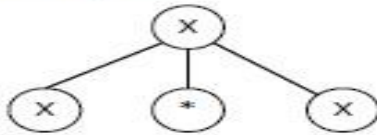
---

$$X \rightarrow X+X \mid X*X \mid X \mid a$$

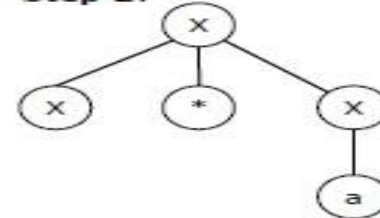
The rightmost derivation for the above string "a+a\*a" may be –

$$X \rightarrow X*X \rightarrow X*a \rightarrow X+X*a \rightarrow X+a*a \rightarrow a+a*a$$

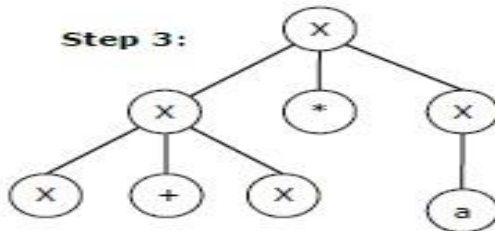
Step 1:



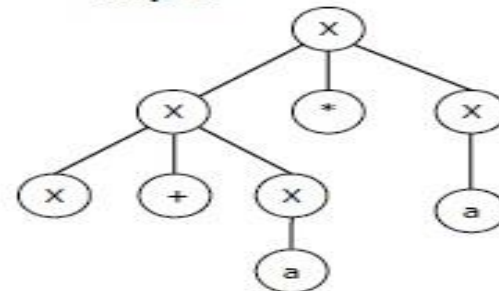
Step 2:



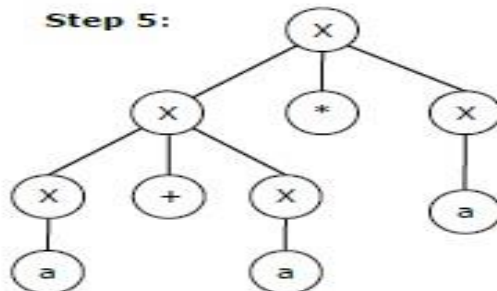
Step 3:



Step 4:



Step 5:



# Ambiguous grammar

---

A grammar is said to be ambiguous if for any string generated by it, it produces more than one-

- Parse tree
- Or derivation tree
- Or syntax tree
- Or leftmost derivation
- Or rightmost derivation

Ambiguous Grammar creates confusion for parser

Check the Grammar :  $E \rightarrow E + E / E \times E / id$

# $E \rightarrow E + E / E \times E / id$ ambiguous grammar

---

$E \rightarrow E + E$

$\rightarrow id + E$

$\rightarrow id + E \times E$

$\rightarrow id + id \times E$

$\rightarrow id + id \times id$

$E \rightarrow E \times E$

$\rightarrow E + E \times E$

$\rightarrow id + E \times E$

$\rightarrow id + id \times E$

$\rightarrow id + id \times id$

**Leftmost Derivation-01**

$E \rightarrow E + E$

$\rightarrow E + E \times E$

$\rightarrow E + E \times id$

$\rightarrow E + id \times id$

$\rightarrow id + id \times id$

**Leftmost Derivation-02**

$E \rightarrow E \times E$

$\rightarrow E \times id$

$\rightarrow E + E \times id$

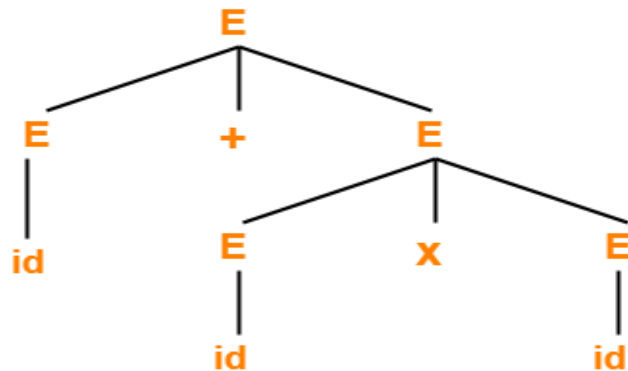
$\rightarrow E + id \times id$

$\rightarrow id + id \times id$

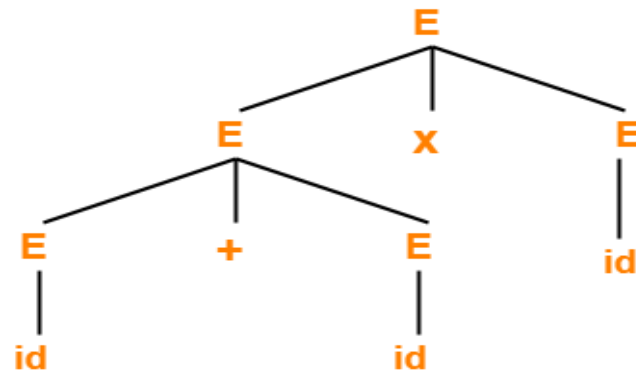
**Rightmost Derivation-01**

**Rightmost Derivation-02**

# $E \rightarrow E + E \mid E \times E \mid id$ ambiguous grammar

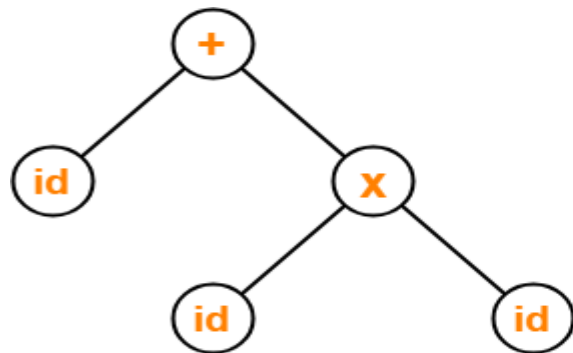


Parse Tree-01

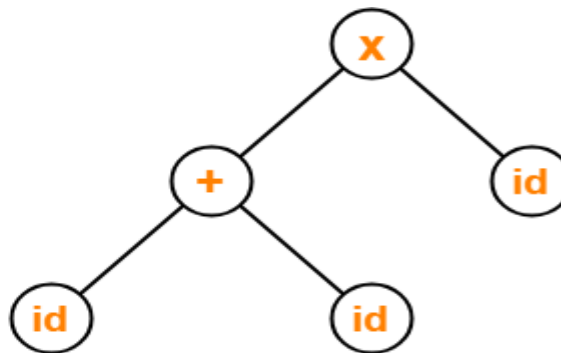


Parse Tree-02

Since two parse trees exist for string  $w$ , therefore the grammar is ambiguous.



Syntax Tree-01



Syntax Tree-02

Consider a string  $w$  generated by the grammar-

$$w = id + id \times id$$

Since two syntax trees exist for string  $w$ , therefore the grammar is ambiguous.

---

## Left and Right Recursive Grammars

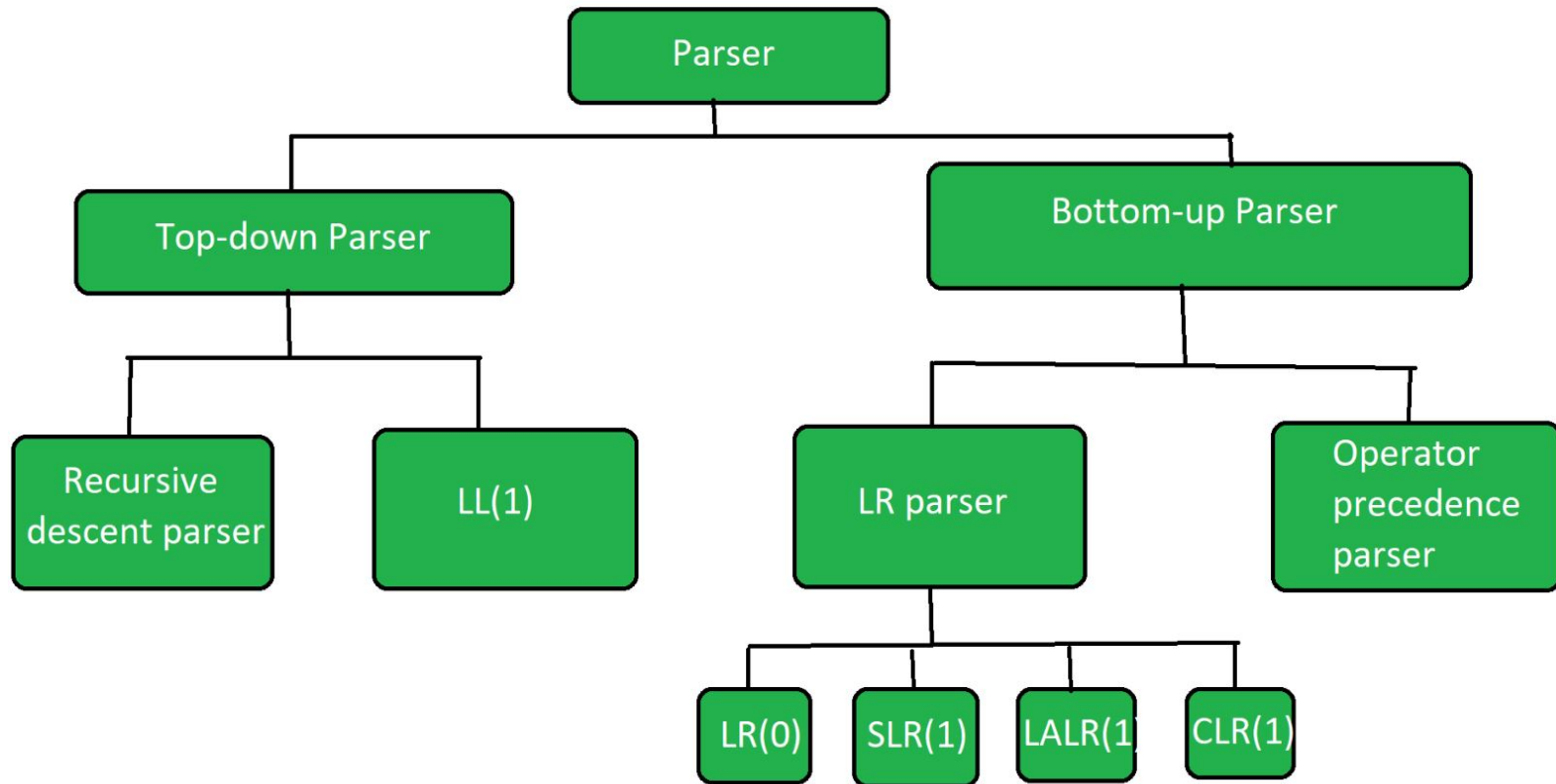
In a context-free grammar  $G$ , if there is a production in the form  $X \rightarrow Xa$  where  $X$  is a non-terminal and 'a' is a string of terminals, it is called a left recursive production. The grammar having a left recursive production is called a left recursive grammar.

And if in a context-free grammar  $G$ , if there is a production is in the form  $X \rightarrow aX$  where  $X$  is a non-terminal and 'a' is a string of terminals, it is called a right recursive production. The grammar having a right recursive production is called a right recursive grammar.

# PARSING

---

It is a process of deriving string from a given grammar.



# Rules for calculating first()

---

First and Follow sets are needed so that the parser can properly apply the needed production rule at the correct position.

$\text{First}(\alpha)$  is a set of terminal symbols that begin in strings derived from  $\alpha$ .

Consider the production rule-

$A \rightarrow abc / def / ghi$

Then, we have-

$\text{First}(A) = \{ a, d, g \}$

Rules For Calculating First Function-

## **Rule-01:**

For a production rule  $X \rightarrow \epsilon$ ,

$\text{First}(X) = \{ \epsilon \}$

## **Rule-02:**

For any terminal symbol 'a',

$\text{First}(a) = \{ a \}$

---

### Rule-03:

For a production rule  $X \rightarrow Y_1Y_2Y_3$ ,

Calculating  $\text{First}(X)$

If  $\epsilon \notin \text{First}(Y_1)$ , then  $\text{First}(X) = \text{First}(Y_1)$

If  $\epsilon \in \text{First}(Y_1)$ , then  $\text{First}(X) = \{ \text{First}(Y_1) - \epsilon \} \cup \text{First}(Y_2Y_3)$

Calculating  $\text{First}(Y_2Y_3)$

If  $\epsilon \notin \text{First}(Y_2)$ , then  $\text{First}(Y_2Y_3) = \text{First}(Y_2)$

If  $\epsilon \in \text{First}(Y_2)$ , then  $\text{First}(Y_2Y_3) = \{ \text{First}(Y_2) - \epsilon \} \cup \text{First}(Y_3)$

Similarly, we can make expansion for any production rule  $X \rightarrow Y_1Y_2Y_3\dots Y_n$ .



---

Consider the production rule-

$$A \rightarrow abc / def / ghi$$

Then, we have-

$$\text{First}(A) = \{ a , d , g \}$$

---

Calculate the first and follow functions for the given grammar-

$S \rightarrow aBDh$

$B \rightarrow cC$

$C \rightarrow bC / \epsilon$

$D \rightarrow EF$

$E \rightarrow g / \epsilon$

$F \rightarrow f / \epsilon$

### First Functions-

- $\text{First}(S) = \{ a \}$
- $\text{First}(B) = \{ c \}$
- $\text{First}(C) = \{ b, \epsilon \}$
- $\text{First}(D) = \{ \text{First}(E) - \epsilon \} \cup \text{First}(F) = \{ g, f, \epsilon \}$
- $\text{First}(E) = \{ g, \epsilon \}$
- $\text{First}(F) = \{ f, \epsilon \}$

### Follow Functions-

- $\text{Follow}(S) = \{ \$ \}$
- $\text{Follow}(B) = \{ \text{First}(D) - \epsilon \} \cup \text{First}(h) = \{ g, f, h \}$
- $\text{Follow}(C) = \text{Follow}(B) = \{ g, f, h \}$
- $\text{Follow}(D) = \text{First}(h) = \{ h \}$
- $\text{Follow}(E) = \{ \text{First}(F) - \epsilon \} \cup \text{Follow}(D) = \{ f, h \}$
- $\text{Follow}(F) = \text{Follow}(D) = \{ h \}$

# Rules for calculating follow()

---

$\text{Follow}(\alpha)$  is a set of terminal symbols that appear immediately to the right of  $\alpha$ .

## Rule-01:

For the start symbol  $S$ , place  $\$$  in  $\text{Follow}(S)$ .

## Rule-02:

For any production rule  $A \rightarrow \alpha B$ ,  
 $\text{Follow}(B) = \text{Follow}(A)$

## Rule-03:

For any production rule  $A \rightarrow \alpha B \beta$ ,  
If  $\epsilon \notin \text{First}(\beta)$ , then  $\text{Follow}(B) = \text{First}(\beta)$   
If  $\epsilon \in \text{First}(\beta)$ , then  $\text{Follow}(B) = \{ \text{First}(\beta) - \epsilon \} \cup \text{Follow}(A)$

---

$S \rightarrow AaAb \mid BbBa$

$A \rightarrow \epsilon$

$B \rightarrow \epsilon$

$\text{Follow}(A) = \{a, b\}$

$\text{Follow}(B) = \{b, a\}$

$S \rightarrow ABC$   $\text{Follow}(A) = \text{First}(B)$  if it is  $\epsilon$ , So,  $\text{First}(C)$  then  $\text{follow}(S) = \$$

$A \rightarrow DEF$

$B \rightarrow \epsilon$

$C \rightarrow \epsilon$

$D \rightarrow \epsilon$

$E \rightarrow \epsilon$

$F \rightarrow \epsilon$

---

Calculate the first and follow functions for the given grammar-

$S \rightarrow aBDh$

$B \rightarrow cC$

$C \rightarrow bC / \epsilon$

$D \rightarrow EF$

$E \rightarrow g / \epsilon$

$F \rightarrow f / \epsilon$

### **First Functions-**

- $\text{First}(S) = \{ a \}$
- $\text{First}(B) = \{ c \}$
- $\text{First}(C) = \{ b, \epsilon \}$
- $\text{First}(D) = \{ \text{First}(E) - \epsilon \} \cup \text{First}(F) = \{ g, f, \epsilon \}$
- $\text{First}(E) = \{ g, \epsilon \}$
- $\text{First}(F) = \{ f, \epsilon \}$

### **Follow Functions-**

- $\text{Follow}(S) = \{ \$ \}$
- $\text{Follow}(B) = \{ \text{First}(D) - \epsilon \} \cup \text{First}(h) = \{ g, f, h \}$
- $\text{Follow}(C) = \text{Follow}(B) = \{ g, f, h \}$
- $\text{Follow}(D) = \text{First}(h) = \{ h \}$
- $\text{Follow}(E) = \{ \text{First}(F) - \epsilon \} \cup \text{Follow}(D) = \{ f, h \}$
- $\text{Follow}(F) = \text{Follow}(D) = \{ h \}$

	First()	Follow()
S->ABCDE	{a, b, c}	{ $\$$ }
A->a/ $\epsilon$	{a, $\epsilon$ }	{b, c}
B->b/ $\epsilon$	{b, $\epsilon$ }	{c}
C->c	{c}	{d, e, $\$$ }
D->d/ $\epsilon$	{d, $\epsilon$ }	{e, $\$$ }
E->e/ $\epsilon$	{e, $\epsilon$ }	{ $\$$ }

	First()	Follow()
S->ACB/CbB/Ba	{d, g, h, $\epsilon$ , b, a}	{ $\$$ }
A->da/BC	{d, g, h, $\epsilon$ }	{h, g, $\$$ }
B->g/ $\epsilon$	{g, $\epsilon$ }	{ $\$$ , a, h, g}
C->h/ $\epsilon$	{h, $\epsilon$ }	{g, $\$$ , b, h}

	First()	Follow()
S->Bb/Cd	{a,b,c,d}	{ $\$$ }
B->aB/ $\epsilon$	{a, $\epsilon$ }	{b}
C->cC/ $\epsilon$	{c, $\epsilon$ }	{d}

# LL1 parser

---

It is a non recursive predictive parsing

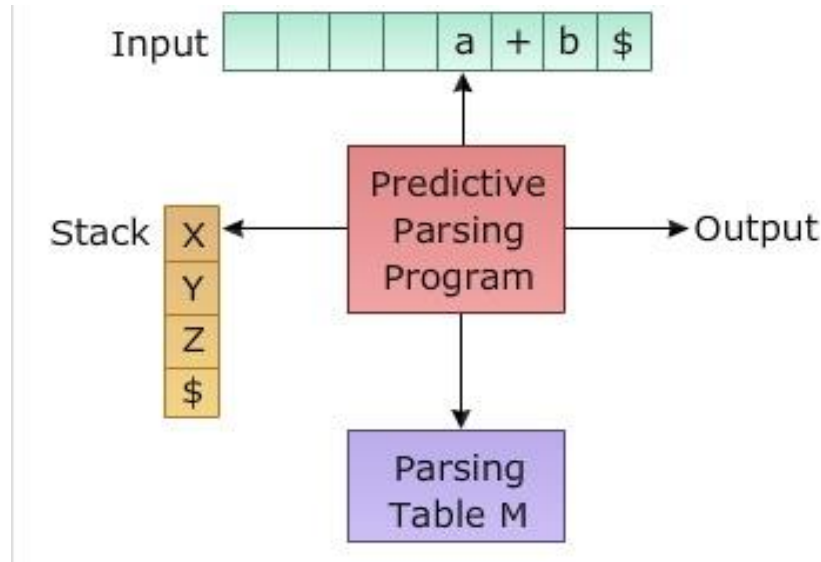
It is a top-down parser.

In this, L - Scanning input from left to right,

L - Producing a leftmost derivation and

1- One input symbol of lookahead at each step

- A grammar  $G$  is LL(1) iff whenever  $A \rightarrow \alpha \mid \beta$  are two distinct productions of  $G$ , the following conditions hold:
  - ✓ **Condition 1:** For no terminal  $a$ , do both  $\alpha$  and  $\beta$  derive strings beginning with  $a$ .  
( $\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \phi$ )
  - ✓ **Condition 2:** At most one of  $\alpha$  and  $\beta$  can derive empty string.
  - ✓ **Condition 3:** If  $\beta \xrightarrow{*} \epsilon$  then  $\alpha$  does not derive any string beginning with a terminal in  $\text{FOLLOW}(A)$ .  
( $\text{FIRST}(\alpha) \cap \text{FOLLOW}(A) = \phi$ )



**INPUT:** Contains string to be parsed with \$ as it's end marker

**STACK:** Contains sequence of grammar symbols with \$ as it's bottom marker.  
Initially stack contains only \$

**PARSING TABLE:** A two dimensional array  $M[A,a]$ , where A is a non-terminal and a is a Terminal



S->AA --- (1) first() - S-{a, b}

A->aA --- (2) A-{a, b}

A->b --- (3) follow() S-{\$}

A-{\$, a, b }

<b>a</b>	<b>b</b>	<b>a</b>	<b>b</b>	<b>\$</b>
----------	----------	----------	----------	-----------

Input Buffer

Lookahead

	<b>a</b>	<b>b</b>	<b>\$</b>
<b>S</b>	S->AA	S->AA	
<b>A</b>	A->aA	A->b	

Parsing Table

S->AA->aAA->abA->abaA->abab

b-match (3)
a- match (1)
A
b -match (3)
a- match (2)
A
A - (1)
A
S (1)
\$

Stack

Check following grammar for LL(1):

This grammar is  
LL(1).

$S \rightarrow CC$   
 $C \rightarrow cC \mid d$

$S \rightarrow CC$   
 $C \rightarrow cC \mid d$

$\text{FIRST}(S) = \{c, d\}$   
 $\text{FIRST}(C) = \{c, d\}$

$\text{FOLLOW}(S) = \{\$ \}$   
 $\text{FOLLOW}(C) = \{c, d, \$ \}$

Non-terminal	Input Symbol		
	c	d	\$
S	$S \rightarrow CC$	$S \rightarrow CC$	
C	$C \rightarrow cC$	$C \rightarrow d$	

This grammar is  
LL(1).

$S \rightarrow AaAb \mid BbBa$

$A \rightarrow \epsilon$

$B \rightarrow \epsilon$

$\text{FIRST}(S) = \{a, b\}$

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

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

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

$\text{FOLLOW}(A) = \{a, b\}$

$\text{FOLLOW}(B) = \{a, b\}$

Non-terminal	Input Symbol		
	a	b	\$
S	$S \rightarrow AaAb$	$S \rightarrow BbBa$	
A	$A \rightarrow \epsilon$	$A \rightarrow \epsilon$	
B	$B \rightarrow \epsilon$	$B \rightarrow \epsilon$	

---

Find if the given grammar is LL1

$S \rightarrow aSbS / bSaS / \epsilon$

First()- {a, b,  $\epsilon$  }

Follow()- {b, a, \$}

	<b>a</b>	<b>b</b>	<b>\$</b>
<b>S</b>	$S \rightarrow aSbS$ $S \rightarrow \epsilon$	$S \rightarrow bSaS$ $S \rightarrow \epsilon$	$S \rightarrow \epsilon$

It shows multiple entries hence it is not LL1 grammar

This grammar is not LL(1).

$S \rightarrow iEtSS' \mid a$   
 $S' \rightarrow eS \mid \epsilon$   
 $E \rightarrow b$

$\text{FIRST}(S) = \{i, a\}$

$\text{FIRST}(S') = \{e, \epsilon\}$

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

$\text{FOLLOW}(S) = \{e, \$\}$

$\text{FOLLOW}(S') = \{e, \$\}$

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

Multiply-defined entries

Non-terminal	Input Symbol					
	a	b	e	i	t	\$
S	$S \rightarrow a$			$S \rightarrow iEtSS'$		
S'			$S' \rightarrow \epsilon$ $S' \rightarrow eS$			$S' \rightarrow \epsilon$
E		$E \rightarrow b$				

# Left recursive grammar

---

A production of grammar is said to have left recursion if the leftmost variable of its RHS is same as variable of its LHS.

A grammar containing a production having left recursion is called as Left Recursive Grammar. Elimination of Left Recursion

Left recursion is eliminated by converting the grammar into a right recursive grammar.

If we have the left-recursive pair of productions-

$$A \rightarrow A\alpha / \beta$$

(Left Recursive Grammar)

where  $\beta$  does not begin with an A.

Then, we can eliminate left recursion by replacing the pair of productions with-

For  $A \rightarrow A\alpha / \beta$

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' / \epsilon$$

(Right Recursive Grammar)

This right recursive grammar functions same as left recursive grammar.

# Right recursive Grammar

---

A production of grammar is said to have right recursion if the rightmost variable of its RHS is same as variable of its LHS.

A grammar containing a production having right recursion is called as Right Recursive Grammar.

Example-

$$S \rightarrow aS / \epsilon$$

(Right Recursive Grammar)

Right recursion does not create any problem for the Top down parsers.

Therefore, there is no need of eliminating right recursion from the grammar.

# Eliminate Left recursion

---

Consider the following grammar and eliminate left recursion-

$$E \rightarrow E + T / T$$

$$T \rightarrow T \times F / F$$

$$F \rightarrow \text{id}$$

Solution-

For $A \rightarrow A\alpha / \beta$
$A \rightarrow \beta A'$
$A' \rightarrow \alpha A' / \epsilon$

The grammar after eliminating left recursion is-

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

$$T' \rightarrow \times FT' / \epsilon$$

$$F \rightarrow \text{id}$$



# Eliminate Left recursion

---

➤ Consider the following grammar and eliminate left recursion-

$A \rightarrow ABd / Aa / a$

$B \rightarrow Be / b$

**Solution:**

$A \rightarrow aA'$

$A' \rightarrow BdA' / aA' / \epsilon$

$B \rightarrow bB'$

$B' \rightarrow eB' / \epsilon$

For $A \rightarrow A\alpha / \beta$
$A \rightarrow \beta A'$
$A' \rightarrow \alpha A' / \epsilon$

➤ Consider the following grammar and eliminate left recursion-

$S \rightarrow S0S1S / 01$

**Solution:**

$S \rightarrow 01A$

$A \rightarrow 0S1SA / \epsilon$

---

If RHS of more than one production starts with the same symbol, then such a grammar is called as Grammar With Common Prefixes.

$$A \rightarrow \alpha\beta1 / \alpha\beta2 / \alpha\beta3$$

This kind of grammar creates a problematic situation for Top down parsers. Top down parsers can not decide which production must be chosen to parse the string in right hand.

To remove this confusion we use left factoring,

**\*It converts nondeterministic CFG into deterministic CFG**

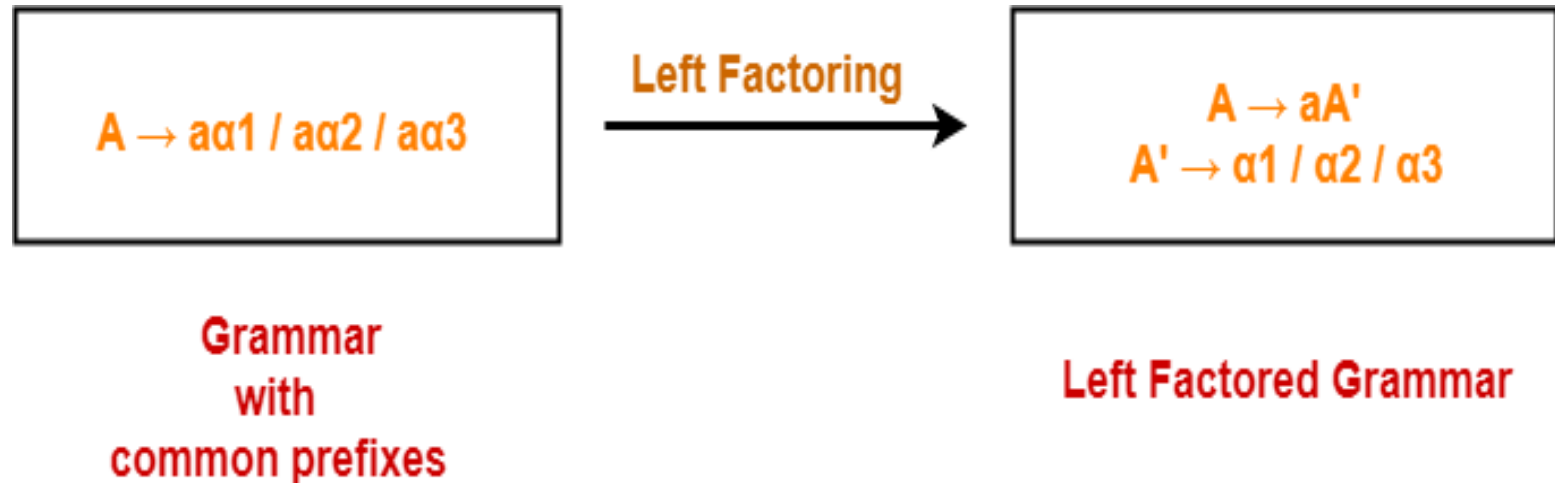
We make one production for each common prefixes.

The common prefix may be a terminal or a non-terminal or a combination of both.

Rest of the derivation is added by new productions.

---

The grammar obtained after the process of left factoring is called as Left Factored Grammar. n, we use left factoring.



Do left factoring in the following grammar-

$A \rightarrow aAB/aA$

$B \rightarrow bB/b$

Solution:

$A \rightarrow aA'$

$A' \rightarrow AB/A$

$B \rightarrow bB'$

$B' \rightarrow B/\epsilon$

$A \rightarrow aA'$   
 $A' \rightarrow \alpha1 / \alpha2 / \alpha3$

Do left factoring in the following grammar-

$A \rightarrow aAB / aBc / aAc$

Solution-

Step1:  $A \rightarrow aA'$

$A' \rightarrow AB / Bc / Ac$

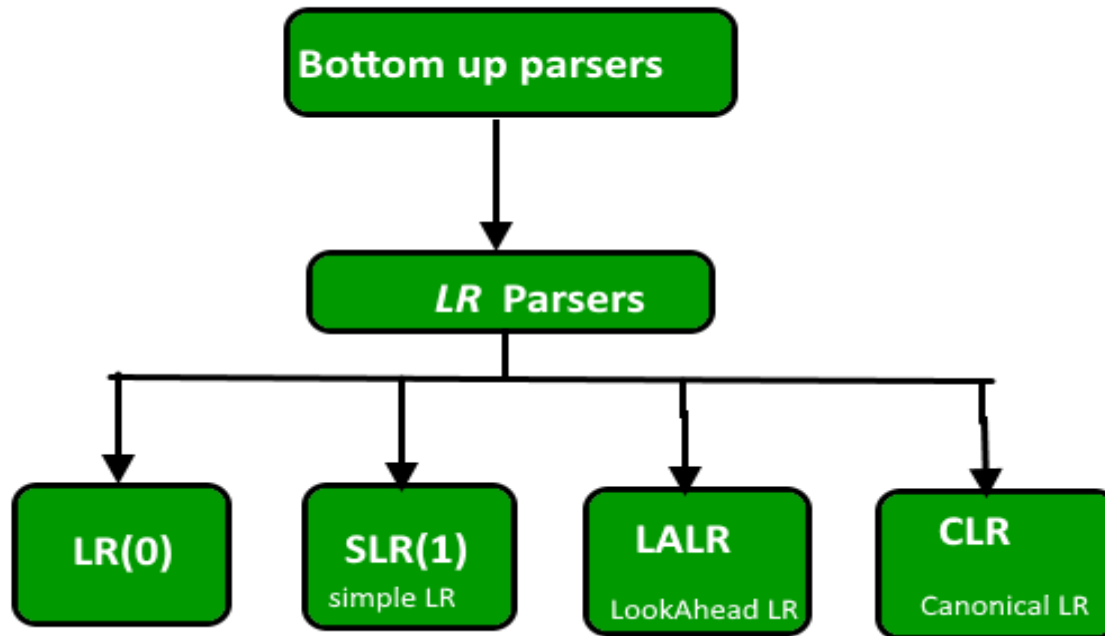
Step2:  $A \rightarrow aA'$

$A' \rightarrow AD / Bc$

$D \rightarrow B / c$

# Bottom up Parsing

---



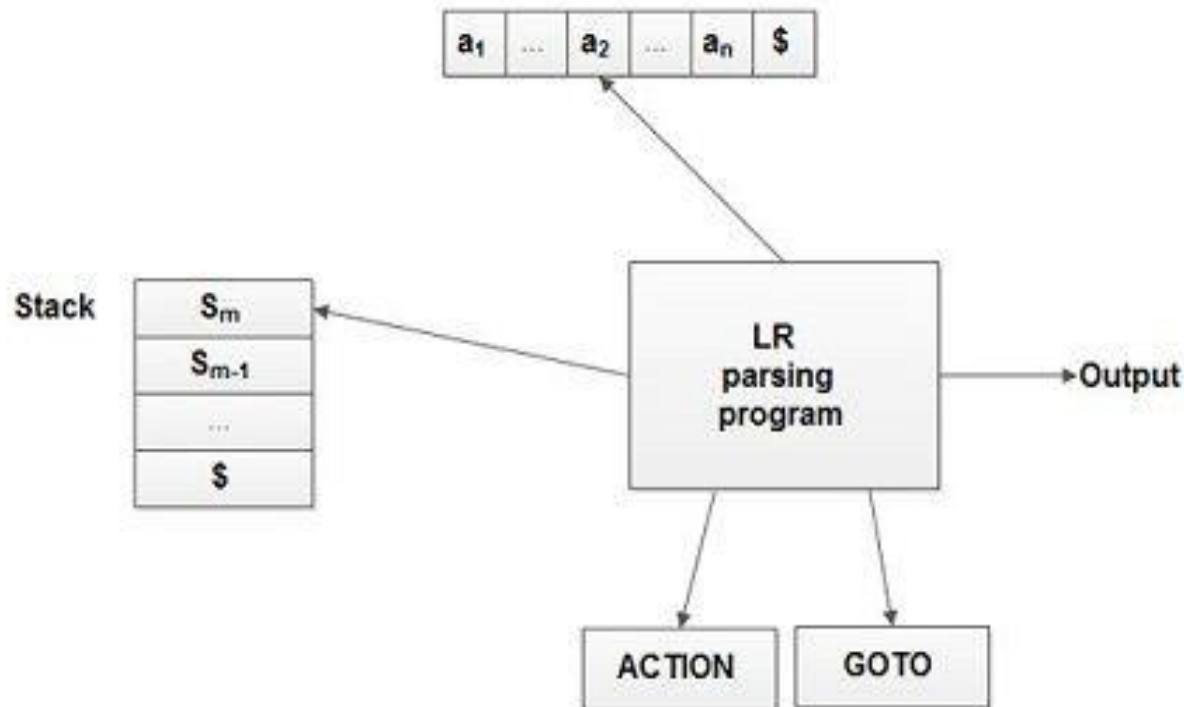
# LR parser

It is a non-recursive shift reduce parser – LR(k)

L- Left to right scanning input stream

R- Construction of rightmost derivation which is in reverse manner

k- Denotes lookahead symbol to make decision



---

An LR-Parser uses

- states to memorize information during the parsing process,
- an action table to make decision (such as shift or reduce) and to compute states
- a goto table to compute states

For a state  $s$  and a terminal  $a$ , the entry  $\text{action}[s, a]$  has one of the following forms

- shift  $s'$  where  $s'$  is a state,
- reduce  $A \rightarrow \beta$ ,
- accept,
- error.

# Example : LR (0) grammar

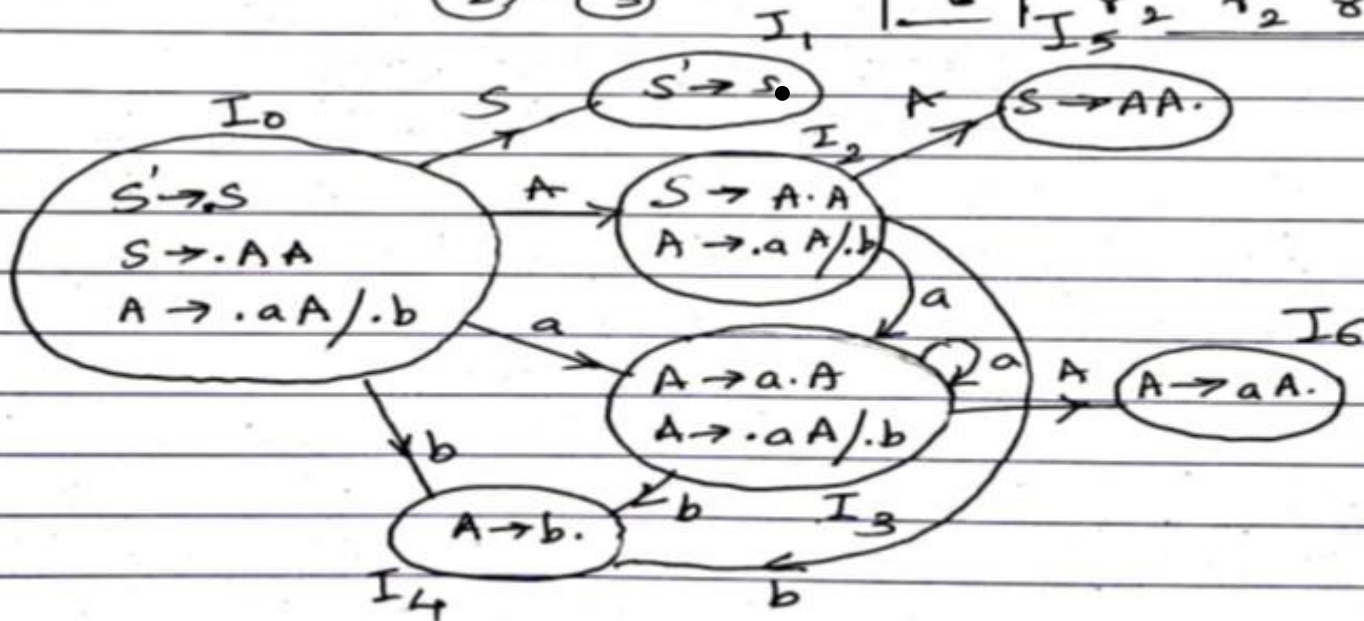
Given :  $S \rightarrow AA$   
 $A \rightarrow aA/b$

Consider the grammar -  
 $S' \rightarrow S$  Augmentation

$S \rightarrow AA$  ①

$A \rightarrow aA/b$   
 ② ③

	action			Goto	
	a	b	\$	A	S
0	<del><math>S_3</math></del>	<del><math>S_4</math></del>		2	1
1			accept		
2	$S_3$	$S_4$		5	
3	$S_3$	$S_4$		6	
4	$r_3$	$r_3$	$r_3$		
5	$r_1$	$r_1$	$r_1$		
6	$r_2$	$r_2$	$r_2$		



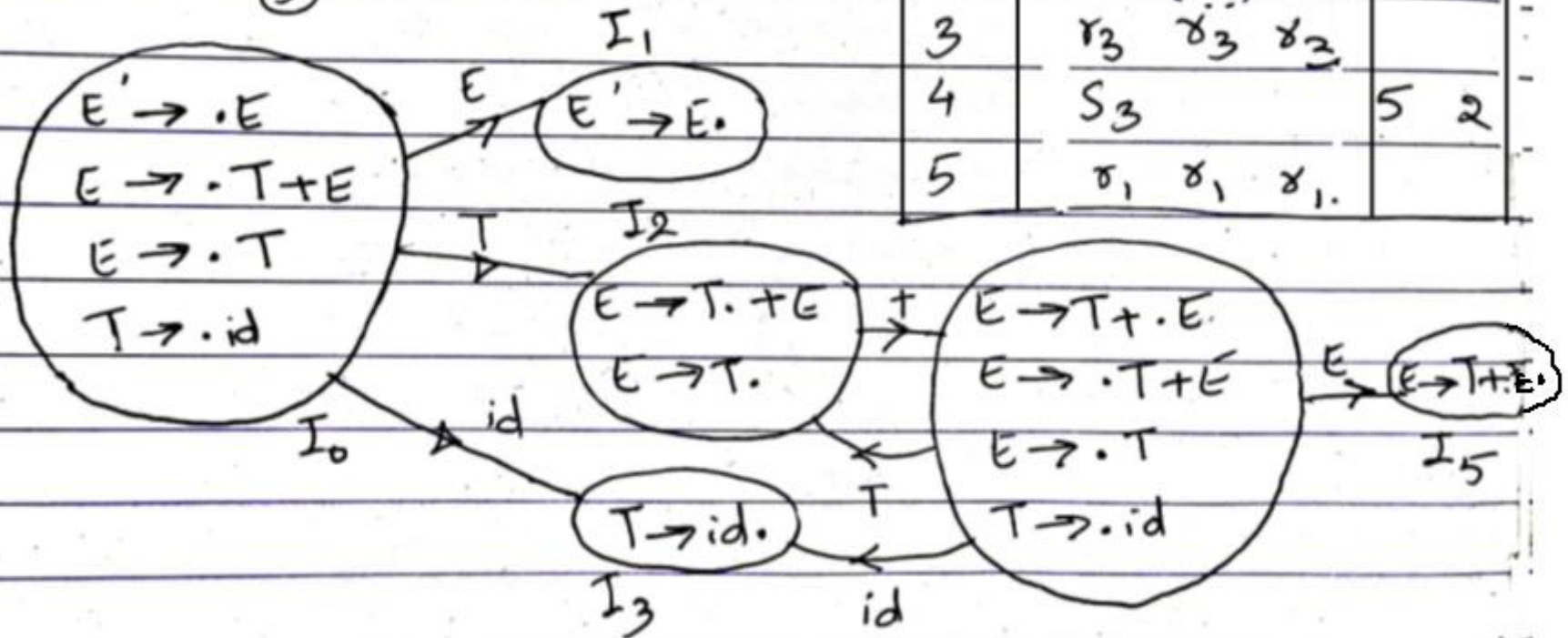
PARSING  
TABLE



# Example: Non LR(0) grammar

Given :  $E \rightarrow T + E / T$   
 $T \rightarrow id$

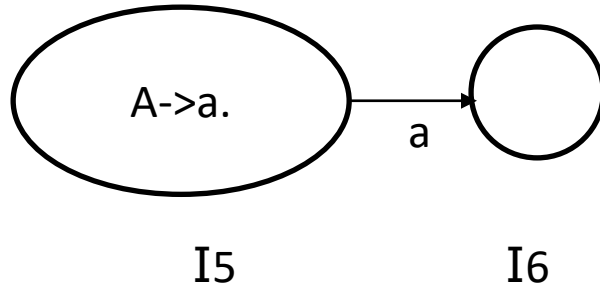
①  
 $E \rightarrow T + E / T$   
 ②  
 $T \rightarrow id$   
 ③



	id	+	\$	E	T
0	$S_3$			1	2
1	-	-	Accept		
2	$r_2$	$S_1/r_2$	$r_2$		
3	$r_3$	$r_3$	$r_3$		
4	$S_3$			5	2
5	$r_1$	$r_1$	$r_1$		

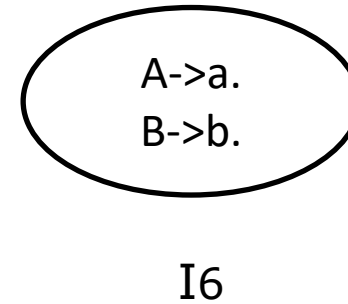
# S-R and R-R conflicts

S-R conflicts



	<b>a</b>	<b>b</b>	
5	S6/r1	r1	

R-R conflicts



	<b>a</b>	<b>b</b>	<b>c</b>	
5	r1/r2	r1/r2	r1/r2	

---

## **Reasons for attractiveness of LR parser**

- LR parsers can handle a large class of context-free grammars.
- The LR parsing method is a most general non-back tracking shift-reduce parsing method.
- An LR parser can detect the syntax errors as soon as they can occur.
- LR grammars can describe more languages than LL grammars.

## **Drawbacks of LR parsers**

- It is too much work to construct LR parser by hand. It needs an automated parser generator.
- If the grammar contains ambiguities or other constructs then it is difficult to parse in a left-to-right scan of the input.

# LL parsing vs LR parsing

---

LL parsing	LR parsing
Left most derivative	Right most derivative
Starts with the root nonterminal on the stack	Ends with the root nonterminal on the stack
Ends when the stack is empty	Starts when the stack is empty
Uses the stack for designating what is still to be expected	Uses the stack for designating what is already seen
Builds the parse tree top-down	Builds the parse tree bottom-up
continuously pops a nonterminal off the stack, and pushes a corresponding right hand side	tries to recognize a right hand side on the stack, pops it, and pushes the corresponding nonterminal
expands nonterminals	reduces them
reads terminal when it pops one off the stack	reads terminals while it pushes them on the stack
uses grammar rules in an order which corresponds to pre-order traversal of the parse tree	does a post-order traversal.

# Operator Grammars

---

Operator grammars have the property that no production on right side is empty (no null productions) or has two adjacent nonterminals. This property enables the implementation of efficient operator-precedence parsers. These parsers rely on the following three precedence relations:

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

Ex:  $E \rightarrow E + E / E * E / id$  ---- is operator grammar (but is ambiguous)

$E \rightarrow E A E / id$  ----- is not an operator grammar

$A \rightarrow + / *$

---

### **Rule-01:**

If precedence of b is higher than precedence of a, then we define  $a < b$

If precedence of b is same as precedence of a, then we define  $a = b$

If precedence of b is lower than precedence of a, then we define  $a > b$

### **Rule-02:**

An identifier is always given the higher precedence than any other symbol.

\$ symbol is always given the lowest precedence.

### **Rule-03:**

If two operators have the same precedence, then we go by checking their associativity.

$E \rightarrow E + E / E \times E / id$

Operation relation (precedence) table is constructed as follows

	Id (Highest prec.)	+	x	\$ (Lowest prec.)
Id	--	.>	.>	.>
+	<.	.> (left associative)	<.	.>
x	<.	.> (left associative)	.>	.>
\$	<.	<.	<.	--

If TOS < input string char, push (shift)  
else pop (Reduce)

\$ < id > + < id > x < id > \$

\$ E + < id > x < id > \$

\$ E + E x < id > \$

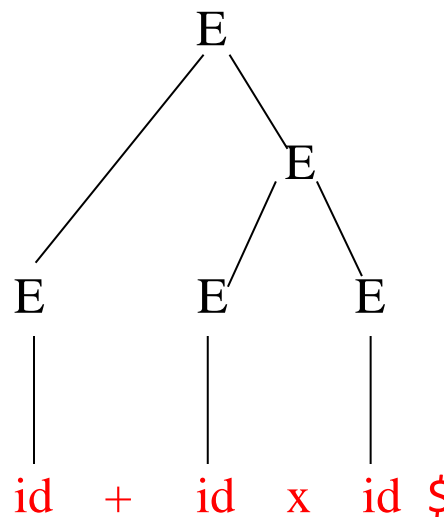
\$ E + E x E \$

\$ + x \$

\$ < + < x > \$

\$ < + > \$

\$ \$



\$	id	+	id	x	id	\$
----	----	---	----	---	----	----

# Shift –Reduce Parsing

$E \rightarrow E + E / E * E / id$

Operation relation (precedence) table is constructed as follows :  $id * id + id$

	Id (Highest prec.)	+	*	\$ (Lowest prec.)
Id	--	.>	.>	.>
+	<.	.> (left associa tive)	<.	.>
*	<.	.>	.> (left associ ative)	.>
\$	<.	<.	<.	--

Stack	Input	Action
\$	id*id+id \$	Shift Id
\$ Id	*id+id \$	Reduce E -> Id
\$ E	*id+id \$	Shift *
\$ E *	id+id \$	Shift Id
\$ E * Id	+id \$	Reduce E -> Id
\$ E * E	+id \$	Reduce E -> E * E
\$ E	+id \$	Shift +
\$ E +	id \$	Shift Id
\$ E + Id	\$	Reduce E -> Id
\$ E + E	\$	Reduce E -> E + E
\$ E	\$	Accept



Consider the following grammar-

$S \rightarrow ( L ) \mid a$

$L \rightarrow L , S \mid S$

Construct the operator precedence parser and parse the string  $( a , ( a , a ) )$ .

	a	(	)	,	\$
a		>	>	>	>
(	<	>	>	>	>
)	<	>	>	>	>
,	<	<	>	>	>
\$	<	<	<	<	

$\$ < ( \underline{\leq a} > , < ( < a > , < a > ) > ) > \$$

$\$ < ( S , < ( \underline{\leq a} > , < a > ) > ) > \$$

$\$ < ( S , < ( S , \underline{\leq a} > ) > ) > \$$

$\$ < ( S , \underline{\leq ( \underline{S} , \underline{S} ) \geq} ) > \$$

$\$ < ( S , \underline{\leq ( \underline{L} , \underline{S} ) \geq} ) > \$$

$\$ < ( S , \underline{\leq ( \underline{L} ) \geq} ) > \$$

$\$ \underline{\leq ( \underline{S} , \underline{S} ) \geq} \$$

$\$ \underline{\leq ( \underline{L} , \underline{S} ) \geq} \$$

$\$ \underline{\leq ( \underline{L} ) \geq} \$$

$\$ \underline{\leq S} > \$$

$\$ \$$

# SLR Parsing

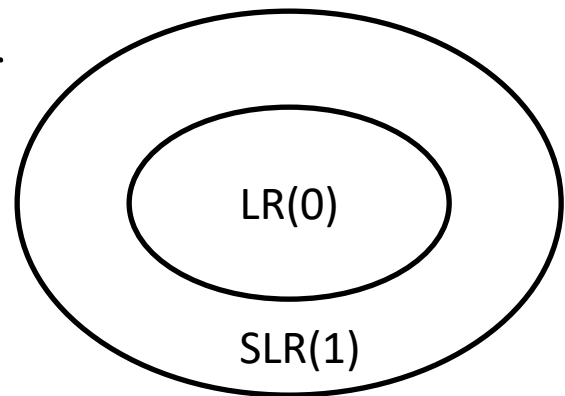
---

SLR (1) refers to simple LR Parsing. It is same as LR(0) parsing. The only difference is in the parsing table. To construct SLR (1) parsing table, we use canonical collection of LR (0) item.

In the SLR (1) parsing, we place the reduce move only in the follow of left hand side.

Various steps involved in the SLR (1) Parsing:

- For the given input string write a context free grammar
- Check the ambiguity of the grammar
- Add Augment production in the given grammar
- Create Canonical collection of LR (0) items
- Draw a data flow diagram (DFA)
- Construct a SLR (1) parsing table



# SLR (<https://www.javatpoint.com/slr-1-parsing>)

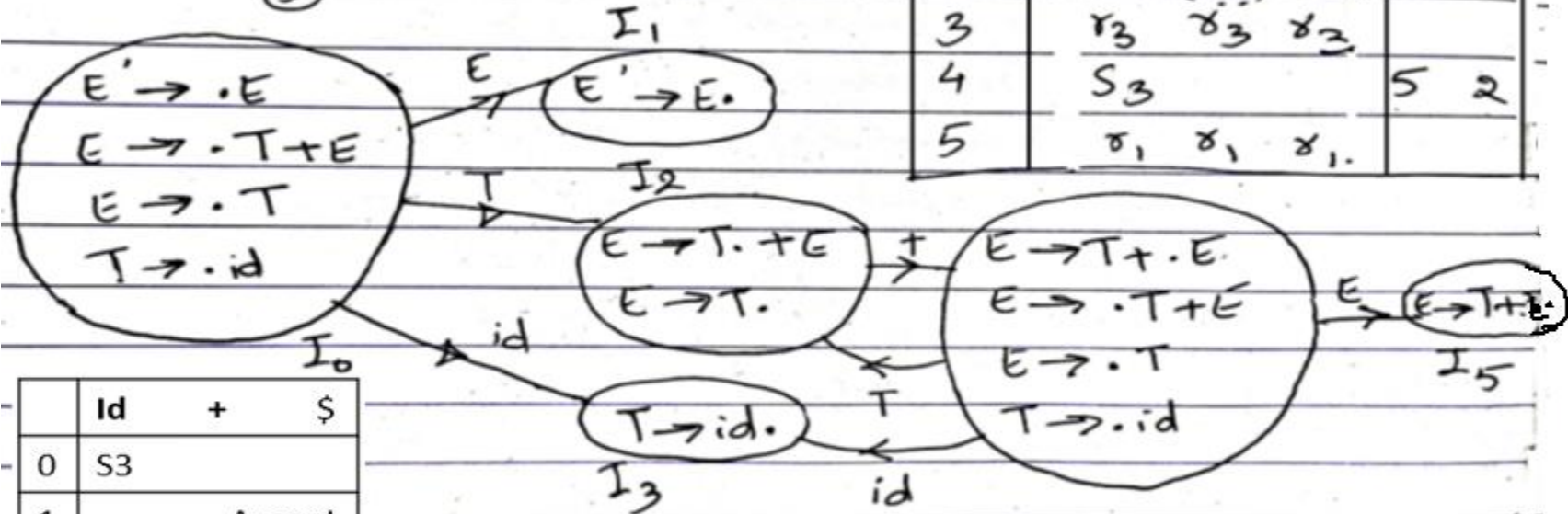
Given :  $E \rightarrow T + E / T$

$T \rightarrow id$

$E \rightarrow T + E / T$

$T \rightarrow id$

	id	+	\$	E	T
0	S <sub>3</sub>			1	2
1	-	-	Accept		
2	r <sub>2</sub>	S <sub>4</sub> /r <sub>2</sub>	r <sub>2</sub>		
3	r <sub>3</sub>	r <sub>3</sub>	r <sub>3</sub>		
4	S <sub>3</sub>			5	2
5	r <sub>1</sub>	r <sub>1</sub>	r <sub>1</sub>		



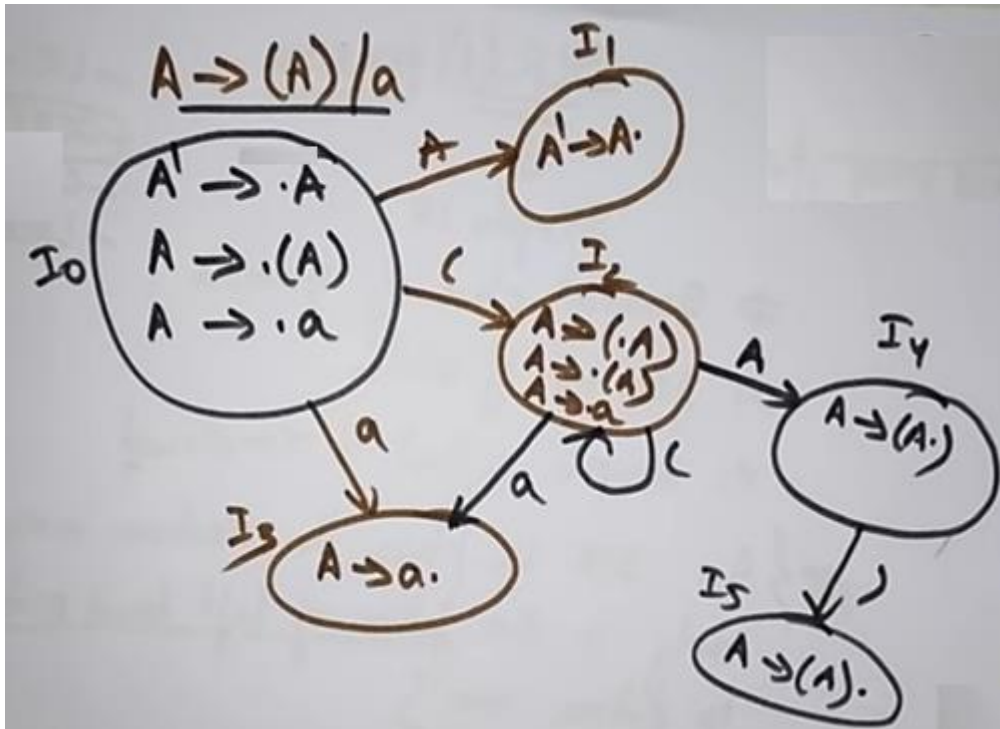
	id	+	\$
0	S <sub>3</sub>		
1		Accept	
2		S <sub>4</sub>	r <sub>2</sub>
3		r <sub>3</sub>	r <sub>3</sub>
4	S <sub>3</sub>		
5			r <sub>1</sub>

Reduced move is written only in the follow of the LHS variable of corresponding production. E.g. r<sub>2</sub> written in only follow of E, r<sub>3</sub> written in only follow of T and r<sub>1</sub> written in follow of E

# Example-SLR

Given:  $A \rightarrow (A) \mid a$  ....  $A \rightarrow (A)$  –rule1 &  $A \rightarrow a$  –rule 2

Augmented rule-  $A' \rightarrow A$



	Action				Goto
	a	(	)	\$	A
0	S3	S2			1
1				Ac c	
2	S3	S2			4
3			r2	r2	
4		S5			
5			r1	r1	

Reduce move is written only in the follow of LHS of the corresponding rule.

Follow(A) = { ), \$ }

# LR(0) vs SLR

---

- LR(0)
  - One LR(0) state mustn't have both shift and reduce items, or two reduce items.
  - So any complete item (dot at end) must be in its own state; parser will always reduce when in this state.
- SLR
  - Peek ahead at input to see if reduction is appropriate.
  - Before reducing by rule  $A \rightarrow XYZ$ , see if the next token is in  $\text{Follow}(A)$ . Reduce *only* in that case. Otherwise, shift.

# Semantic Analysis

---

Semantic analysis checks the semantic consistency of the code. It uses the syntax tree of the previous phase along with the symbol table to verify that the given source code is semantically consistent. It also checks whether the code is conveying an appropriate meaning.

CFG + semantic rules = Syntax Directed Definitions

The semantic analyzer is expected to recognize:

- Type mismatch
- Undeclared variable
- Reserved identifier misuse.
- Multiple declaration of variable in a scope.
- Accessing an out of scope variable.
- Actual and formal parameter mismatch.

# Functions of Semantic analyses phase

---

- Helps you to store type information gathered and save it in symbol table or syntax tree
- Allows you to perform type checking
- In the case of type mismatch, where there are no exact type correction rules which satisfy the desired operation a semantic error is shown
- Collects type information and checks for type compatibility
- Checks if the source language permits the operands or not

```
float x = 20.2;
```

```
float y = x*30;
```

In the above code, the semantic analyzer will typecast the integer 30 to float 30.0 before multiplication

# Syntax directed translation

---

In syntax directed translation, along with the grammar we associate some informal notations and these notations are called as semantic rules.

Grammar + semantic rule = SDT (syntax directed translation)

In syntax directed translation, every non-terminal can get one or more than one attribute or sometimes 0 attribute depending on the type of the attribute. The value of these attributes is evaluated by the semantic rules associated with the production rule.

In the semantic rule, attribute is VAL and an attribute may hold anything like a string, a number, a memory location and a complex record.

In Syntax directed translation, whenever a construct encounters in the programming language then it is translated according to the semantic rules define in that particular programming language.



---

---

Production	Semantic Rules
$E \rightarrow E + T$	$E.val := E.val + T.val$
$E \rightarrow T$	$E.val := T.val$
$T \rightarrow T * F$	$T.val := T.val * F.val$
$T \rightarrow F$	$T.val := F.val$
$F \rightarrow (F)$	$F.val := F.val$
$F \rightarrow num$	$F.val := num.lexval$

$E.val$  is one of the attributes of  $E$ .

$num.lexval$  is the attribute returned by the lexical analyzer.

# Syntax directed translation scheme

---

The Syntax directed translation scheme is a context -free grammar.

The syntax directed translation scheme is used to evaluate the order of semantic rules.

In translation scheme, the semantic rules are embedded within the right side of the productions.

The position at which an action is to be executed is shown by enclosed between braces. It is written within the right side of the production.

**Annotated Parse Tree** – The parse tree containing the values of attributes at each node for given input string is called annotated or decorated parse tree.

**Features** –

High level specification

Hides implementation details

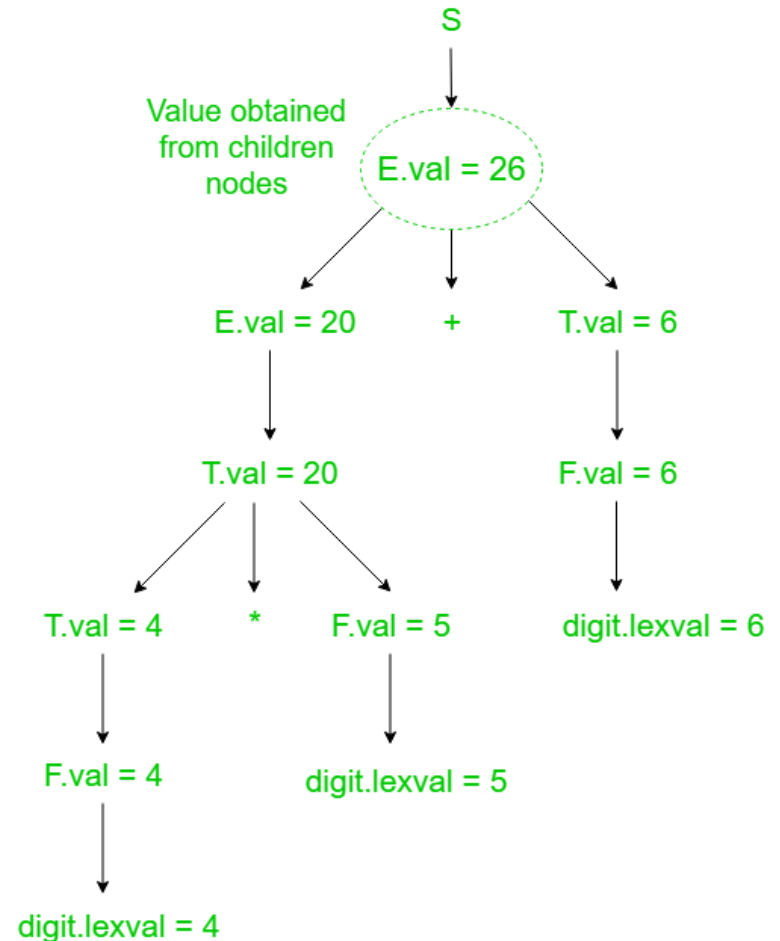
Explicit order of evaluation is not specified

# Implementation of SDT

Consider a grammar: an input string 4\*5+6

```
S --> E
E --> E1 + T
E --> T
T --> T1 * F
T --> F
F --> digit
```

Production	Semantic Actions
$S \rightarrow E$	$\text{Print}(E.\text{val})$
$E \rightarrow E_1 + T$	$E.\text{val} = E_1.\text{val} + T.\text{val}$
$E \rightarrow T$	$E.\text{val} = T.\text{val}$
$T \rightarrow T_1 * F$	$T.\text{val} = T_1.\text{val} * F.\text{val}$
$T \rightarrow F$	$T.\text{val} = F.\text{val}$
$F \rightarrow \text{digit}$	$F.\text{val} = \text{digit}.\text{lexval}$



Annotated Parse Tree

# Implementation of SDT

Syntax direct translation is implemented by constructing a parse tree and performing the actions in a left to right depth first order.

SDT is implementing by parsing the input and produces a parse tree as a result.

Production	Semantic Rules
$S \rightarrow E \$$	{ printE.VAL }
$E \rightarrow E + E$	{ E.VAL := E.VAL + E.VAL }
$E \rightarrow E * E$	{ E.VAL := E.VAL * E.VAL }
$E \rightarrow (E)$	{ E.VAL := E.VAL }
$E \rightarrow I$	{ E.VAL := I.VAL }
$I \rightarrow I \text{ digit}$	{ I.VAL := 10 * I.VAL + LEXVAL }
$I \rightarrow \text{digit}$	{ I.VAL := LEXVAL }

23\*5+4

