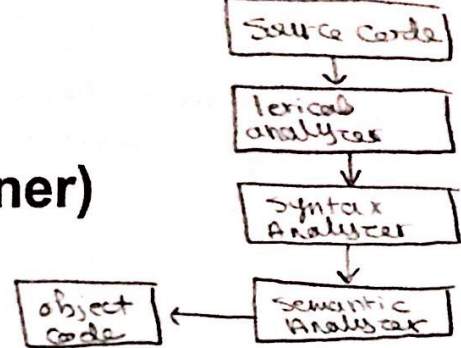


# Chapter 4

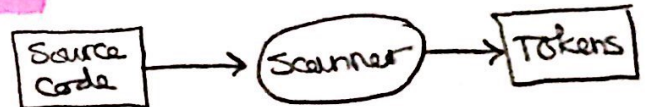
## Lexical Analysis(Scanner)



### What is a Lexical Analyzer?

is an algorithm which recognizes the set of tokens in the source code. A Lexical Analyzer or Scanner is an algorithm that groups the characters of the source code to form the Tokens. Afterwards, it returns the Internal Representation Number of these tokens, which is an ID that matches the reserved word fetched from a keywords table which the implementer of the compiler predefines.

These tokens are divided into 3 kinds:



1. **Names** : Which is any name we have in a program. These in turn are divided into 2 types:
  - a. **Keywords/Reserved**, which are words such as if/else/while. These names can't be used as variable names. They have a specific place and function
  - b. **User Defined Names**, Which are the names declared by the user. *(i.e) variable names, constant names, function names ----*
2. **Values** : such as integers(1, 2, 3, 4) or floating point(1.1, 2.34, 5234.123) etc
3. **Special Symbols/Tokens** : And these are the logical(==, & ||) and arithmetic operations(+, -, \*, /), parenthesis([], {}, ()), or any other tokens that are not from the first or second kind.

!!! Comments aren't part of the source code so they aren't considered tokens

Let us apply the scanner to this short segment of code:

```

while(x>=100)
{
    n +=x;
    x++
}
  
```

This results in this set of tokens :

```

While , ( , x , >= , 100 , ) , { , n , += , x , ; , x , ++ , } .
  
```

Referencing these tokens against a certain keywords table like this one :

Token-ID  
Token

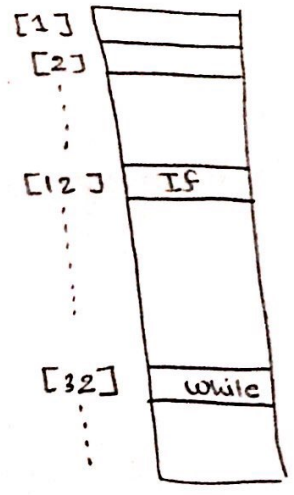
index	Symbol
..	..
33	While
..	..
67	>=
..	..
..	..

Leads us to these ID's :

Token-ID

Token	Internal Representation Number
While	33
(	84
x	100
>=	67
100	200
)	85
{	92
n	100
+=	77
x	100
;	81
x	100
++	75
;	81
}	93

\* the Compiler maintains a table of the reserved words (sorted)



\* All the user defined names are given the same ID(index) because the type of the name of the user defined name doesn't affect the syntax.

note that all user defined names have the same number. This is because to the syntax analyzer, it doesn't matter what the variable is, it just matters that there is a variable there.

**Type Checking implementation**

During this process of analysis, The Compiler builds what is called the **Symbol Table**. The Symbol Table is a table of the name of each user defined name (mostly variables), its type, and its values. The Symbol table for this segment of code would be:

```

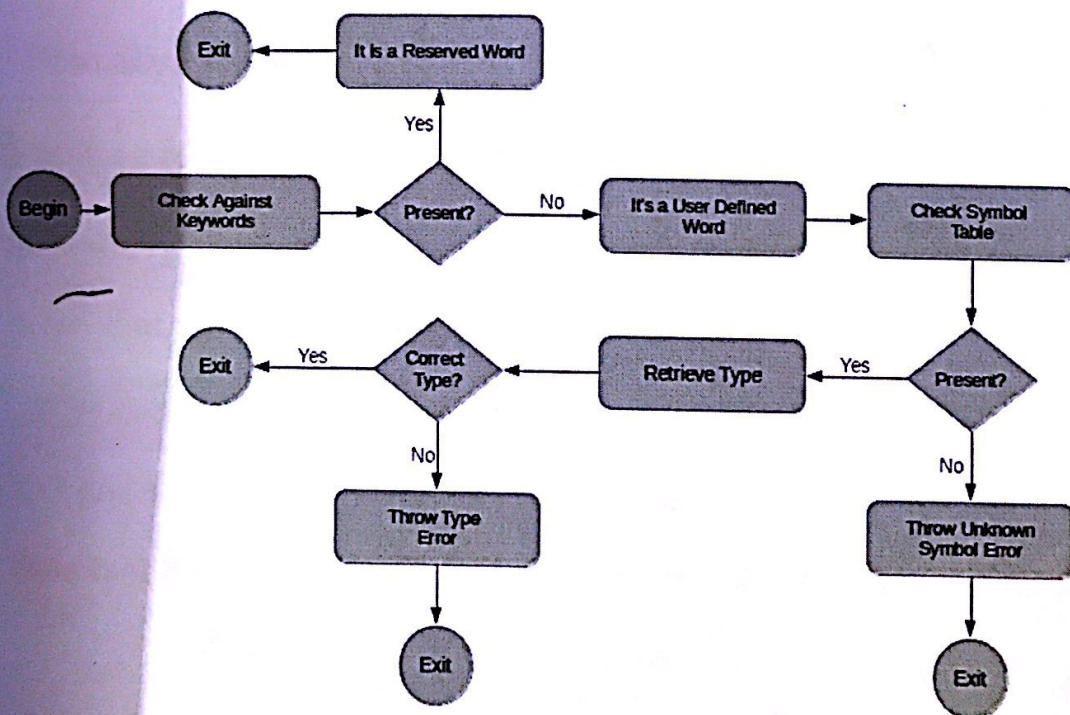
int compute(int,int);
int n; // 0.52 < Float -> Error [Type mismatch]
float x,y;
const int m=10;
compute(int,int) -> Error [Identifier isn't declared]

```

Name	Type	Value
compute	function-name	0
n	Integer	0
x	float	0
y	float	0
m	const-int	10

} user-defined names.

To perform type checking, the compiler takes the name, and checks the keywords table. If it is not in the keywords table, it is a user defined variable. if it is a user defined variable, it then goes to check the symbol table. If it is not defined in the symbol table, it returns that the variable is not defined (unknown symbol/variable deceleration error), if it is in the symbol table, it retrieves its type. If the operation being performed on the variable is not compatible with the type of the variable, it returns that the operation is not compatible( type error ).



# Regular Expressions(Regular Languages)

Regular Languages are a class of language that are important for lexical analysis, since we use them to define and generate tokens. This Class of languages is defined recursively.

## Defining a Language as a Set

We Say that :

**Def:** An Alphabet is a Set of Symbols.

For example, our alphabet is the set  $V = \{a, b, \dots, y, z\}$ .

$V = \{0, 1\}$  Binary alphabet

**Def:** A String is a Sequence of Symbols Taken From the Alphabet.

So if  $V$  is our alphabet, then :

abc  
dsd  
a  
qwe  
az  
asa  
sd

⚠ When we define  $*$  on  $\mathbb{R}$   
then  $*$  is a Binary  
operation defined on  $\mathbb{R}$ .

Are all strings defined on  $V$ .

Note: we can define an infinite number of strings on an alphabet.

Formally, we define:

$S$  is the set of all strings over some alphabet  $V$ .

→ takes 2 operations.

let  $V = \{0, 1\}$  with  $S$  defined over it, let us define a binary operation on  $S$

given that  $x, y \in S$ , then we can define a concatenation

operation on  $x$  and  $y$  as follows:

$XY = \{\text{the set of strings formed by following } x \text{ with } y\}$ .

عَرَبِيَّة

Note that  $XY \neq YX$ . We say that the concatenation operation is not commutative.

**Def:** Given a set of strings  $S$  on  $V$ , then define the Concatenation operation  $(\cdot)$  on  $S$  as follows:

$$X \cdot Y = XY = \{XY \text{ such that } Y \text{ is following } X\}$$

$X, Y \in S$

(4)

Let us Define a special string, called the empty string, which we denote with  $\lambda$ .  
 Notice that  $\lambda X = X\lambda = X$ .

Formally, we can say:

$\lambda$  is the **identity** element of the **concatenation** operation.

Where the Identity element is formally defined as :

Generally: If  $e$  is the Identity element of some operation  $\otimes$  then:  
 $e \otimes x = x \otimes e = x$

Putting all of this together, We can define a language as :

**Def:** Given an alphabet  $V$ , a language  $L$  over  $V$  is a set of strings formed from  $V$ .

By this definition, these sets Are all languages:

- $L_1 = \{a, b, c\}$
- $L_2 = \{asdasd, qwe, asd\}$
- $L_3 = \{abb\}$

- $V = \{a, b, c, \dots, x, y, z\}$
- $L_4 = \{\lambda\} \equiv$  language which has only  $\lambda$  as a string.
- $L_5 = \{aa, abc, bc, \lambda\}$
- $L_6 = \emptyset \equiv \{\}$  empty string  
 a language which has no elements  $\equiv$  empty language

This definition also leads us to the conclusion :

There are an  $\infty$  number of languages defined on an alphabet.

### Set of Operations On Languages

$|S| = \infty$   
 $\hookrightarrow$  cardinality of  $S \equiv \#$  of elements of  $S$

Given an Alphabet  $V$ , assume that :

1.  $L = \{\text{set of all languages defined on } V\}$
2.  $L = \{L_1, L_2, L_3, \dots, L_n\}$

We will define a 3 operations on  $L$ , ie, the operands are languages belonging to  $L$ .

### Concatenation Operation "BINARY OPERATION"

Given that  $L, M$  are languages over an alphabet  $V$ , then  $L, M \in \mathcal{L}$

$LM =$  "L concatenated with M"  $= \{xy \mid x \in L, y \in M\}$ .

For Example let  $L = \{a, b, c\}$  and  $M = \{aa, bb\}$ , then :

$LM = \{aaa, abb, baa, bbb, caa, cbb\}$

$ML = \{aaa, aab, aac, bba, bbb, bbc\}$

Note that :

1.  $LM \neq ML$ . (Concatenation on languages is not commutative).
2.  $L\{\lambda\} = \{\lambda\}L = L$ . ( $\lambda$  is the identity for concatenation).

NOTE  $\triangle$

$\lambda \lambda \lambda$



Example on the or operation:

given  $L = \{ab, bc, bb\}$ ,  $M = \{\lambda, ccc\}$

$L \cup M = \{ab, bc, bb, \lambda, ccc\}$

$M \cup L = \{\lambda, ccc, ab, bc, bb\}$

Example on the closure operation:

given  $L = \{aab, b\}$  over  $V = \{a, b\}$

$L^* = L \cup L^2 \cup L^3 \dots = \bigcup_{i=0}^{\infty} L^i$

$L^* = \{\lambda\} \cup \{aab, b\} \cup \{aabaab, aabb, baab, bb\} \cup \dots$

$L^* = \{\lambda, aab, b, aabaab, aabb, baab, bb, \dots\}$

$L^+ = L^* - \{\lambda\}$

Example on the recursive definition of Regular language

$V = \{a, b, c\}$

①  $\emptyset = \{\}$

②  $\lambda = \{\lambda\}$

③  $a = \{a\}$

$b = \{b\}$

$c = \{c\}$

④  $R = LR$

$S = LS$

⑤

$RS$  a regular expression denoting the languages

$LR \& LS \Rightarrow$  (i.c)  $ab = \{a\} \cdot \{b\} = \{ab\}$

⑥  $RIS$  a regular expression denoting  $LR \cup LS = LR \cup LS$

⑦  $R^* = LR^*$

$\rightarrow (ab)^* = (\{a\} \cup \{b\})^* = (\{a, b\})^*$   
 $= \{\lambda, ab, aa, ab, ba, bb, \dots\}$

(i.d)  $alb = \{a\} \cup \{b\}$   
 $= \{a, b\}$

(i.c)  $a lab = \{a\} \cup \{ab\}$   
 $= \{a, ab\}$

$$3. L\{\} = \{\}L = \{\}.$$

## The OR "|" Operation "BINARY OPERATION"

OR is  $\cup$  operation in set theory

Union

Given that L, M are languages over an alphabet V, then

$$L|M = "L OR M" = \{x \mid x \in L \text{ or } x \in M\} = L \cup M.$$

Note that :

1.  $L|M = M|L$ . (OR on language is commutative)
2.  $L|\{\} = \{\}|L = L$ . (The empty set is the identity element for the OR operation)

NOTE

## The Closure "\*" Operation (A Unary Operation)

→ we use (\*) to denote this operation.

Given that L is a language over an alphabet V then  $L^*$  is:

$$L^* = L^0 \cup L^1 \cup L^2 \cup \dots \cup L^\infty$$

$$L^+ = L^* - \{\lambda\}$$

## The Recursive Definition of Regular Languages

Given an alphabet V then :

1.  $\emptyset = \{\}$  = empty language is a regular language denoting the language  $\{\}$
2.  $\lambda = \{\lambda\}$  is a regular language denoting the language  $\lambda$
3. For every element  $a \in V$ ,  $a = \{a\}$  is a regular language denoting the language  $\{a\}$
4. Given R and S are regular languages denoting the regular languages LR and LS respectively, then :

a) RS is a regular language denoting LR LS

b) R|S is a regular language denoting LR|LS

c)  $R^*$  is a regular language denoting  $LR^*$

## EXAMPLES:

Given  $R=\{a\}$  and  $S=\{b\}$ , then :

-  $RS=\{ab\}$ ,

-  $R|S=\{a,b\}$ ,

-  $R^* = \{a\}^0 \cup \{a\}^1 \cup \dots$

$= \{\lambda, a, aa, aaa, \dots\}$

= A string that consists of any number of a's

- Lets say we took  $(RS)^*$  then

$(RS)^* = \{ab\}^0 \cup \{ab\}^1 \cup \dots$

$= \{\lambda, ab, abab, ababab, abababab, \dots\}$

= A string that consists of any number of "ab"s

- Lets say we took  $(a|b)^*$  , then :

$(a|b)^* = (\{a\}|\{b\})^* = (\{a\} \cup \{b\})^*$

$= (\{a,b\})^* = \text{Any string of a's and b's}$


- Let us say we took  $(0|1)^*00$ , then By the definitions above, this results in any binary string that ends with 00, such as  $\{\lambda, 00, 100, 000, 1100, 0000, \dots\}$  *(!) any string of 0's & 1's that ends with 00.*

- Let us say we took  $(a|b)^*bbb(a|b)^*$ , then By the definitions above, this results in any string of a's and b's that contains at least 3 consecutive b's such as  $\{bbb, abbb, bbba, bbbbbb, \dots\}$



## Example on Exponential Notation

$[+|-] \underline{d} . \underline{d}^+ E \underline{(+|-)} d^+$   
must have only one digit  $\leftarrow$   $\leftarrow$  any digit  $\leftarrow$  there must be  
but  $\lambda$   $\leftarrow$  + or -

 All tokens in the source code are strings (members) in language  $\rightarrow$  type regular expression.

# Defining Tokens Using Regular Languages:

we have 3 types of tokens:

1. Names
2. Values
3. Special Symbols

The scanner must recognize these and be able to distinguish them. ⚠

## Names

In programming languages, names are: letter followed by letters or digits. The regular language for names is :

$$\text{Letter}(\text{Letter}|\text{digit})^* = L(L|d)^*$$

## Values

In programming languages, there are multiple types of values, and they can all be defined using a regular language

1. **Integers :**

$$[+|-]\text{digit}(\text{digit})^* = [+|-]d^* = [+|-]d^+$$

we either take + or don't E. or we

all digits

Note : [x] means we take x zero or one time only. ⚠

2. **Floating Point Numbers :**

$$[+|-]d^+.d^+ \rightarrow \text{Fixed Point}$$

$$[+|-]d^*.d^+ \rightarrow \text{Exponential}$$

0.5, .5

$$[+|-]d.d^+E(+|-)d^+ \rightarrow \text{Exponential}$$

must

## Special Symbols

not these brackets mean that there must be either + or -

The Set of special symbols {+, -, <=, ...} are each given by its own regular languages. For example, the symbol + has its own regular languages given by :

$$+ = \{+\}$$

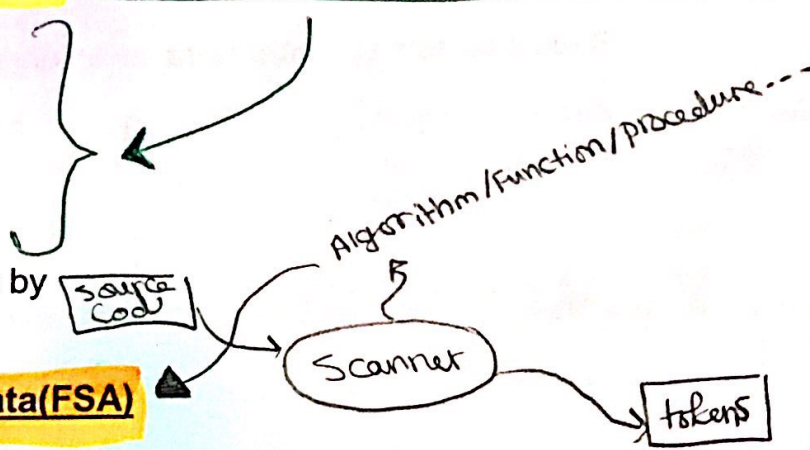
$$* = \{*\}$$

$$<= = \{<\} \{=\} = \{<=\}$$

or ++, which is given by

$$++ = \{+\}\{+\} = \{++\}$$

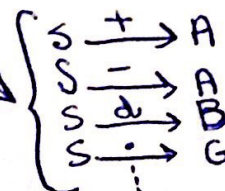
## Finite State Automata(FSA)



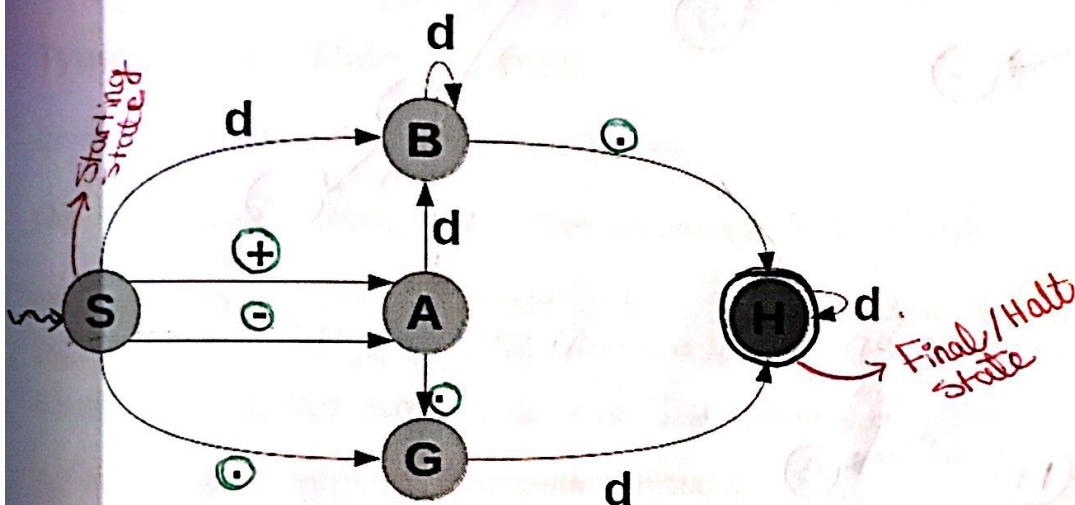
The Question remains: How do we build an algorithm to recognize(accept) strings whose languages are regular languages?

Tokens in source codes are strings of regular languages. The algorithm that recognizes these strings is called the **Finite State Automata** or the **Finite State Machine** or the **Finite State System**. The Finite State Automata contains :

1. A Set of states.  $Q = \{S, A, B, G\}$
2. Transitions between states.
3. Input string to be examined.

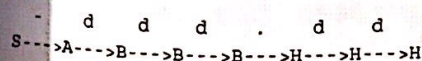


Given that we have this Finite State Automata(FSA) :



- The set of states  $Q = \{S, A, B, G, H\}$   
S is called the **Starting State**.  
H is called the **Final(Halt) State**.
- The transitions between the states are given by  $\{+, -, d, .\}$
- Given the input string is "-ddd.dd" eg "-511.32".

Tracking through the states on this string, we start at state S :

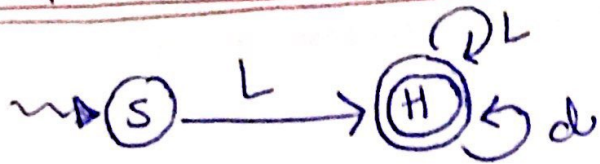


⚠ If after scanning the inputs, the machine ends up in a Final/Halt state we say

Since H is the final state, we say that the string is accepted, or more formally :

**Def:** A string is accepted or recognized if after scanning the

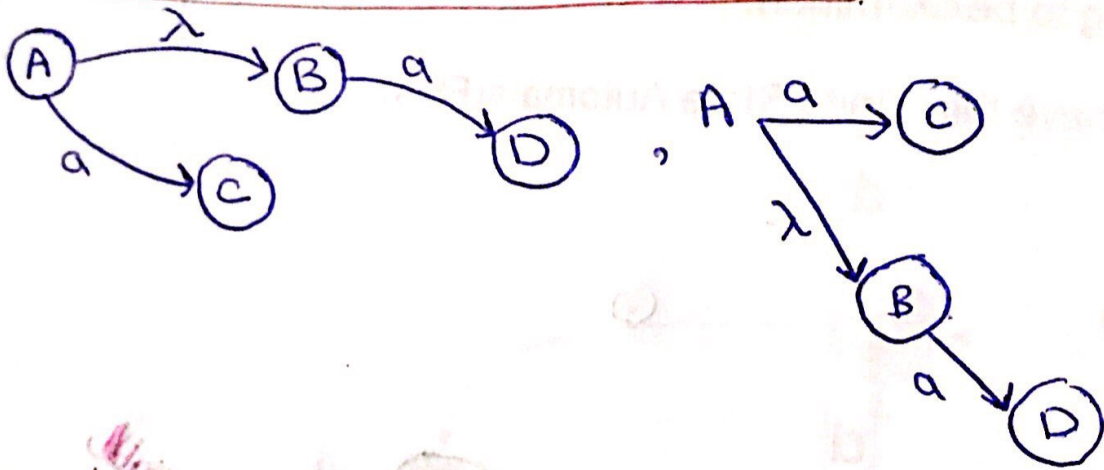
# Example on Finite State Machine For Names



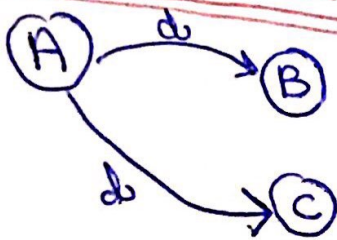
Example:  $\{+ =\}$



(1) If there are  $\lambda$ -transitions in the FSA:



(2) There's more than one transition on the same input from the same state.



string we end up with a final state.

The Regular Language(expression) generated(accepted)

by this machine  $L(M)$  is given by :

$L(M) = [+|-]\{d^+, .d^+, d^+.d^+\}$  → it accepts all floating point numbers in the Fixed point notation.  
+ / - / unsigned.

Other examples of finite state machines are :

1. Names
2. Integers
3.  $\leq$

## Types of Finite State Automata

there are 2 types of Finite State Automata

### (A) Non-Deterministic Finite State Automata [NDFSA]

An algorithm is non-deterministic or fuzzy if there are options(choices) in the algorithm. An example of a non-deterministic algorithm is the solution of the Knight Tour Problem, which is based on Backtracking Techniques.

A Finite State Automata is non-deterministic if :

1. There are  $\lambda$ -transitions(moves) in the FSA :
2. Or There is more than one transition from the same state on the same input :

In Both cases, There is a choices (trial and error) to make. The only way to solve non-deterministic machines is to use backtracking. This is practical in a compiler, because backtracking is a very compute-heavy method and is extremely slow.

Fortunately, There are algorithms to transform any NDFSA to a DFSA. Therefore, we can assume always in the assumption that our machine is deterministic

## (B) Deterministic Finite State Automata

If A Machine is **not** non-deterministic, we call it a **Deterministic Finite State Automata**, ie :

1. There are NO  $\lambda$ -moves(transitions).
2. AND There is NO more than one transition from the same state on the same input.

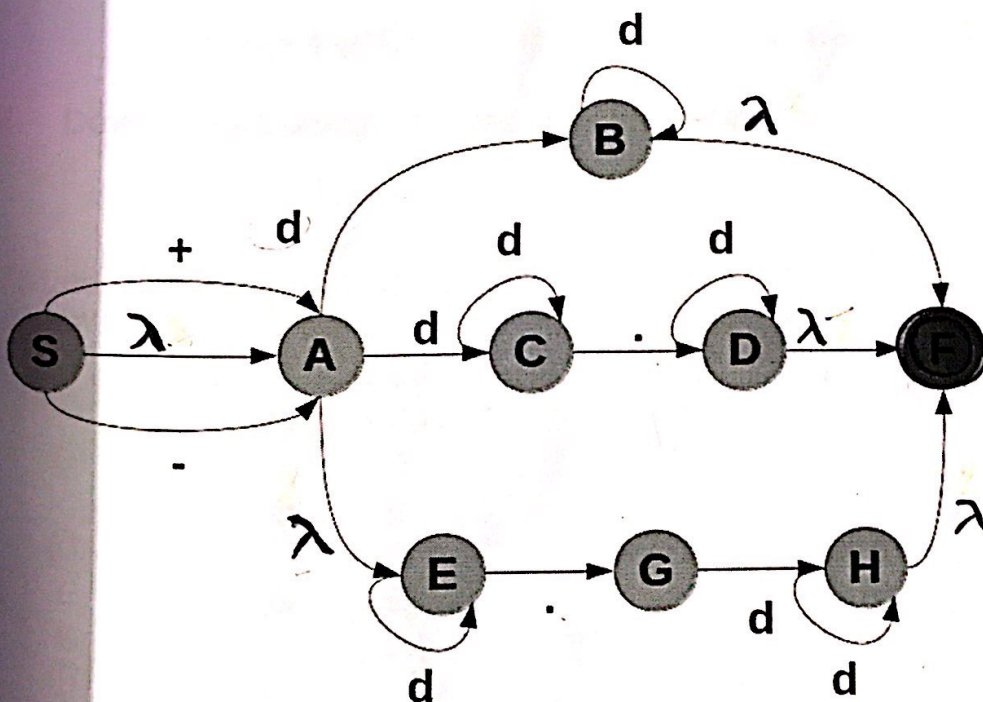
Only Deterministic Finite State Automata are used for compilers.

### Transformation of NDFSA to DFSA

The Algorithm that transforms an NDFSA to a DFSA consists of the following steps :

1. Removal of  $\lambda$  transitions.
2. Removal of non-determinism.
3. Removal of inaccessible states.
4. Merging equivalent states.

Given the following NDFSA :



What is  $L(G) = ?$

$L(G) = [+|-] (d^+ | d^+ \cdot d^+ | d^+ \cdot | \cdot d^+)$

*Signed*  
*unsigned*

*Integers*  
*Floats*

And we want to transform it into a DFSA. Let's follow through the steps :

Let's transform the Finite state machine into a transition table: ①

State \ VT	+	-	.	d	$\lambda$
S	A	A			A
A				B,C	E
B				B	F
C			D	C	
D				D	F
E			G	E	
G				H	
H				H	F

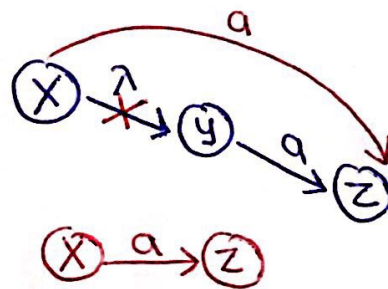
### (1) Removal Lambda Transitions

① Consider  $S \xrightarrow{\lambda} A$

Add all transition in Row A to S.

② Repeat Step(1) for all States with  $\lambda$  Transitions

③ Mark all states from which there is a  $\lambda$  Transition to a final State. Mark it as the final State. *B, D, A are marked as Final states*



4. Delete the  $\lambda$  Column This results in this table :

State \ VT	+	-	.	d
S	A	A	G	B,C,E
A			G	B,C,E
B				B
C			D	C
D				D
E			G	E
G				H
H				H

## (2) Removal Of Non-Determinism

Which mean not having more than 1 transition on 1 input.

- ① Consider [B,C,E]. Lets add this and treat it as a new state in the table.
- ② If at least one of the states [B,C,E] is a final state, then we make it a final state.
- ③ Repeat steps (1) and (2) for all non-deterministic states

④ The Machine is now deterministic

This results in this table :

State \ VT	+	-	.	d
S	A	A	G	B,C,E
A			G	B,C,E
<del>B</del>				B
<del>C</del>			D	C
<del>D</del>				D
<del>E</del>			G	E
<del>G</del>				H
H				H
F				
<b>B,C,E</b>			D,G	B,C,E
<b>D,G</b>				D,H
<b>D,H</b>				D,H

○ → Final State

## (3) Removal of Non-Accessible States

→ the states that can't be reached.

- ① Mark the Initial State
- ② Mark all states for which there is a transition from S
- ③ Repeat step (2) for all marked states.



Do one more extra iteration because some times you mark states above the current state.

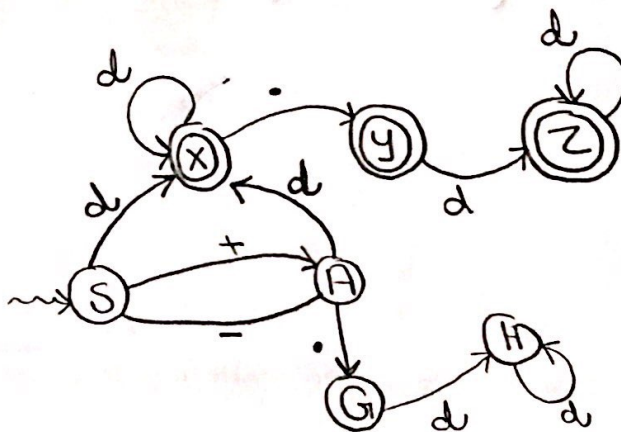
This results in this table :



State \ VT	+	-	.	d
✓S	A	A	G	B,C,E
✓A			G	B,C,E
✗B				B
✗C			D	C
✗D				D
✗E			G	E
✓G				H
✓H				H
✗F				
✓B,C,E			D, G	B,C,E
✗			G	
y ✓D,G				D,H
z ✓D,H				D,H

4. Delete all unmarked states . This results in this simplified Table :

State \ VT	+	-	.	d
✓S	A	A	G	B,C,E
✓A			G	B,C,E
✓G				H
✓H				H
✓B,C,E			D, G	B,C,E
✓D,G				D,H
✓D,H				D,H



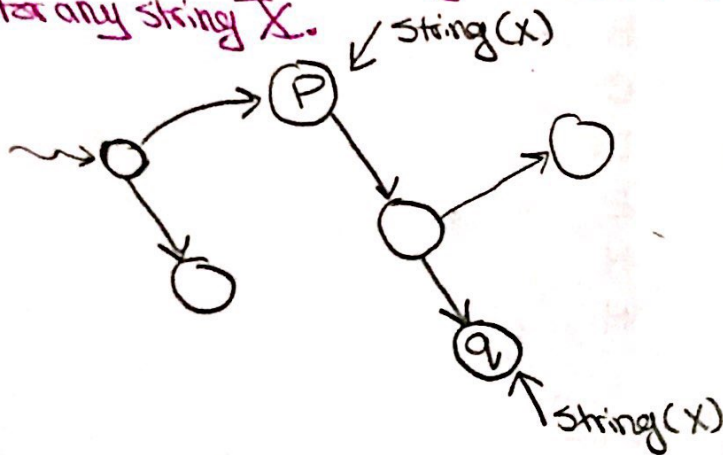
This is now a Deterministic Machine That accepts the same languages as the original NDFSA . For Clarity, Let.s rename [B,C,E] to X, [D,G] to Y, [D,H] to Z.

## (4) Merging equivalent states

### Definition

two states  $p, q$  in the FSA are said to be equivalent if:

$p$  accepts the string  $X$  iff,  $q$  accepts the string  $X$ , for any string  $X$ .



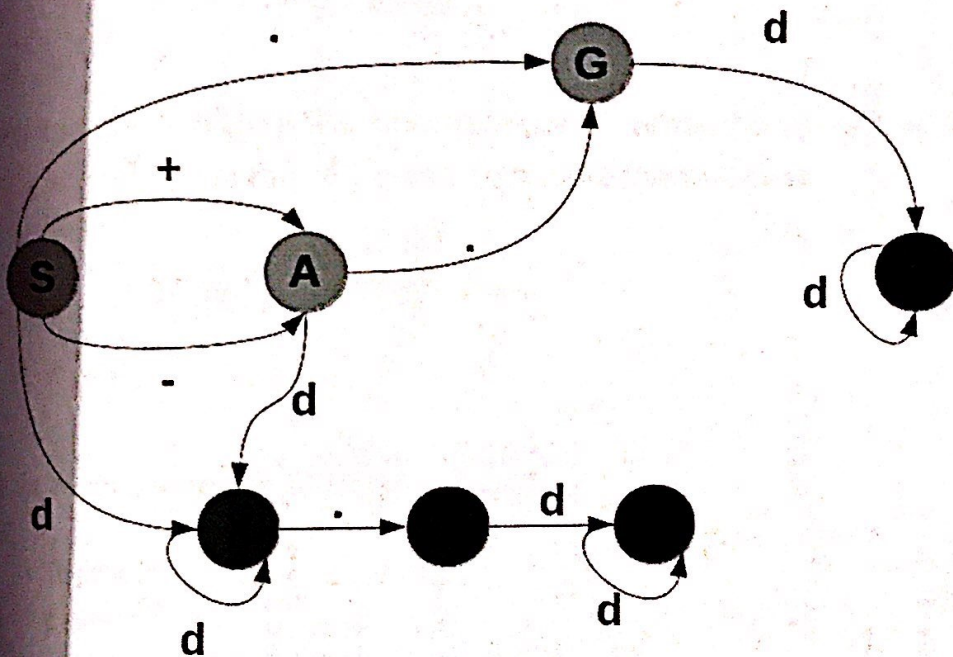
### Revision

A state  $q$  in machine  $M$  & state  $p$  in machine  $M'$  are said to be Equivalent if given a string  $X$ , the machine  $M$  in state  $q$  accepts  $X$  iff machine  $M'$  in state  $p$  accepts  $X$  for every string  $X$ .

$$\{p \equiv q\}$$

State	+	-	.	d
VT				
✓S	A	A	G	X
✓A			G	X
✓G				H
✓H				H
✓X			Y	X
✓Y				Z
✓Z				Z

And the graph now looks like this :



#### (4) Merging Equivalent States

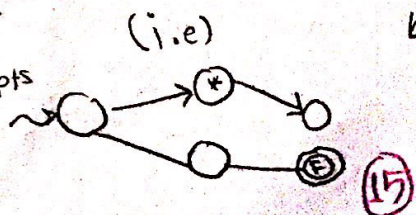
We will use the **feasible-pairs table** method in merging equivalent states.

A state pair  $(p,q)$  is a feasible pair if:

- ①  $\{p,q\} \subset F$  OR  $\{p,q\} \subset Q-F$  ie, either both  $\{p,q\}$  are final states or both  $\{p,q\}$  are not final states.
- ② for every token(symbol)  $a \in VT$ , either both  $\{p,q\}$  have transitions on "a", or both  $\{p,q\}$  don't have transitions on "a".
- ③  $p \neq q$ ;

Note that  $(p,q) \equiv (q,p)$ .

State \* doesn't accept  $\lambda$ , but Final state accepts  $\lambda$  so they're not Feasible



### Definition

A Feasible state pair  $(P, Q)$  is marked if there's a transition from  $(P, Q)$  to a pair  $(R, S)$  such that  $(R, S)$  is either marked or not among the feasible pairs &  $R \neq S$ .

for example, given the following NDFSA, represented by the following transition table :

State \ VT	a	b	c
1	2	5	⊆
2	3	4	⊆
3	5	2	
4	6		1
5	1	4	1
6	4		1
7	3	5	3

(1,2) can't be feasible since there's transition to c(2), but there's no transition to c(1)

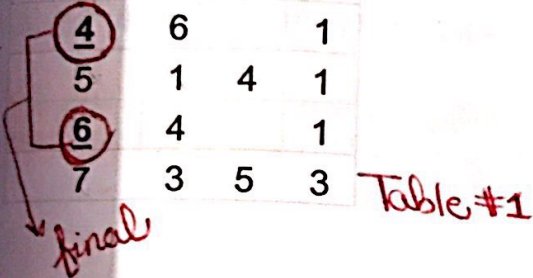


Table #1

To find the feasible pairs, first we must separate the set of final states from the set of non-final states. therefore, we have these 2 sets :

$Q-F = \{1,2,3,5,7\}$  "Non-Final"

$F = \{4,6\}$  "Final"

this results in this feasible-pairs table :

feasible pairs \ VT	a (r,s)	b (r,s)	c (r,s)
(1,3)	2,5 <sub>x</sub>	5,2 <sub>x</sub>	
(2,5)	3,1 <sub>x</sub>	4,4 <sub>x</sub>	1,1 <sub>x</sub>
✓ (2,7)	3,3 <sub>x</sub>	4,5 <sub>✓</sub>	1,3 <sub>x</sub>
✓ (5,7)	1,3 <sub>x</sub>	4,5 <sub>✓</sub>	1,3 <sub>x</sub>
(4,6)	6,4 <sub>x</sub>		1,1 <sub>x</sub>

(1-3) → (2,5) } X  
 ↳ not marked  
 ↳ among feasible pairs X  
 → (5,2) ≡ (2,5) X

Those pairs are feasible to be EQUIVALENT

⚠ repeat the derivation one more time.

We then mark all feasible pairs (p,q) where There is a transition to a pair (r,s) such that :

1.  $r \neq s$ .
2. (r,s) is either marked OR not among the feasible pairs.

This results in this feasible-pairs table :

feasible pairs \ VT	a	b	c
(1,3)	2,5	5,2	
(2,5)	3,1	4,4	1,1
✓(2,7)	3,3	4,5	1,3
✓(5,7)	1,3	4,5	1,3
(4,6)	6,4		1,1

We go through the table once more, in case we marked something later on in the table that would effect the pairs in the top of the table.

if a pair (p,q) remains unmarked, that means that p is equivalent to q. therefore, we merge p and q, choosing one of them :

- ① 1 ≡ 3 ---> 1
- ② 2 ≡ 5 ---> 2
- ③ 4 ≡ 6 ---> 4

We then merge, replacing every 3 with a 1, every 5 with a 2, and every 6 with a 4, resulting in this state table :

State \ VT	a	b	c
1	2	2	
2	1	4	1
④	4		1
7	1	2	1

≡ machine in Table #1.

This is the machine with the minimum number of states.  
Let's go back to our example Last time, we reached this state table :

State \ VT	+	-	.	d
S	A	A	G	X
A			G	X
G				H
⑧				H
⑩			Y	X
⑪				Z
⑫				Z

Lets quickly apply what we learned on this table.

⑬

Separate the final from the non-final states :

$$Q-F = \{S, A, G\}$$

$$F = \{H, X, Y, Z\}$$

Constructing the feasible pairs table :

feasible pairs \ VT	+	-	.	d
(H, Y)				H, Z
(H, Z)				H, Z
(Y, Z)				Z, Z

1. Marking feasible pairs

feasible pairs \ VT	+	-	.	d
(H, Y)				H, Z
(H, Z)				H, Z
(Y, Z)				Z, Z

2. Merge

and Replace

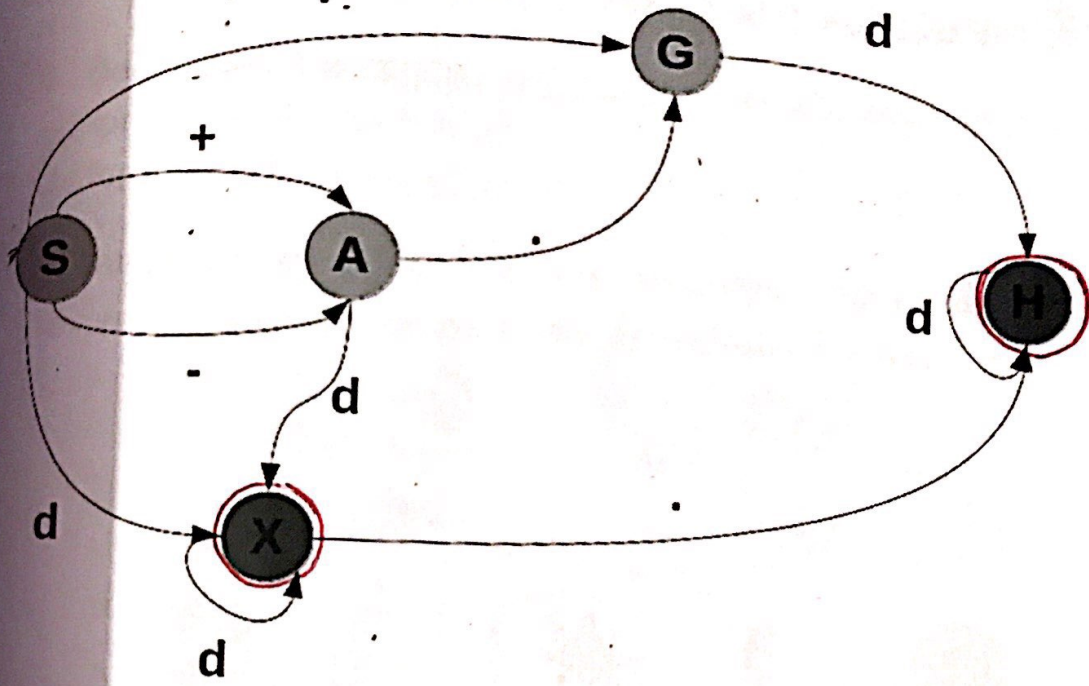
$$H \equiv Y \equiv Z \rightarrow H$$

Resulting in this state table :

State \ VT	+	-	.	d
S	A	A	G	X
A			G	X
G				H
H				H
X			H	X

This is the simplest form of the machine.

Now we must check if the machine accepts the same language as our original machine.



This machine accepts the language.

$$L(G) = [+|-] \{ \underline{d} \underline{d} \underline{d} \underline{d}, \underline{d} \underline{d} \underline{d} \underline{d} \underline{d}, \underline{d} \underline{d} \underline{d} \underline{d} \underline{d} \underline{d}, \underline{d} \underline{d} \underline{d} \underline{d} \underline{d} \underline{d} \}$$

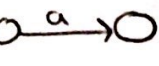
Which is the same language of our original machine.



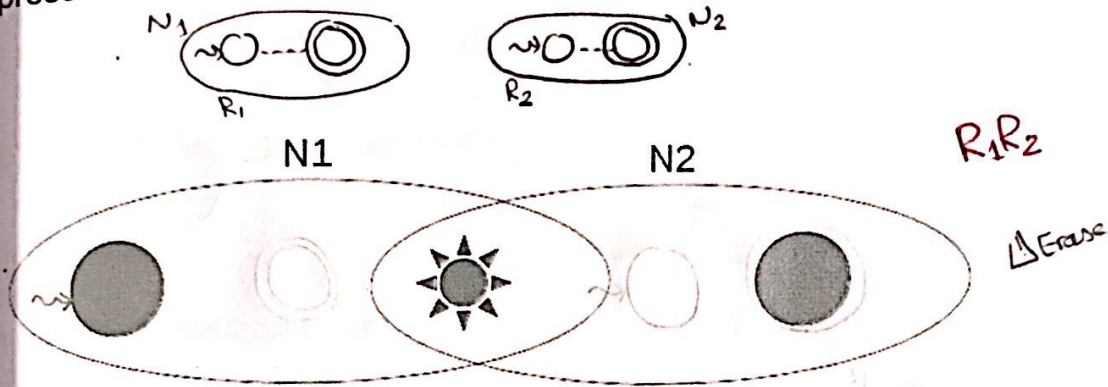
# Creating a NDFSA From a Regular Expression \*

① Decompose the regular expression to its primitive components :

a for  $\lambda$ ,  $X \xrightarrow{\lambda} Y$ . 

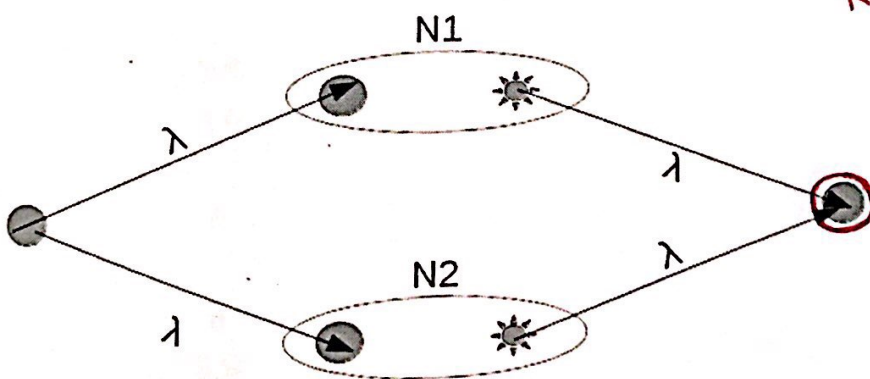
b for  $a$ ,  $X \xrightarrow{a} Y$ . 

② Supposed that  $N_1, N_2$  are transition diagrams for the regular expressions  $R_1, R_2$  respectively,  $N_1$  accepts  $R_1$  &  $N_2$  accepts  $R_2$ , then :

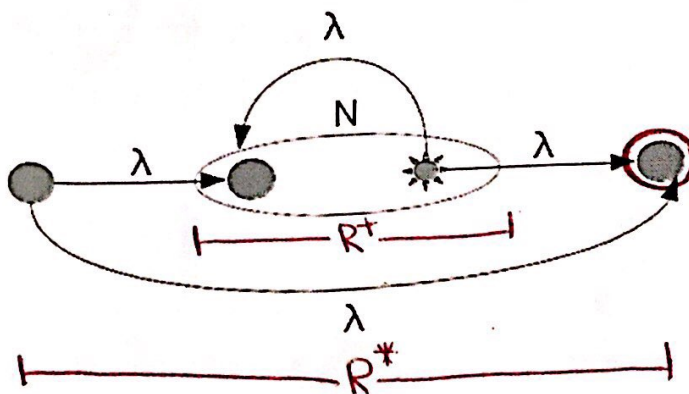


a. which represents  $R_1R_2$  is :

b. which represents  $R_1|R_2$  is :

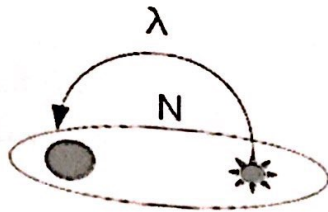


c. which represents  $Rx^*$  is :



d.  $x^+$  N

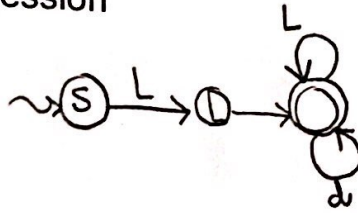
which represents  $R^+$  is :



Note that this is the same as  $Rx^*$  except we removed all the states that result in a  $\lambda$ .

Example: Given the regular expression

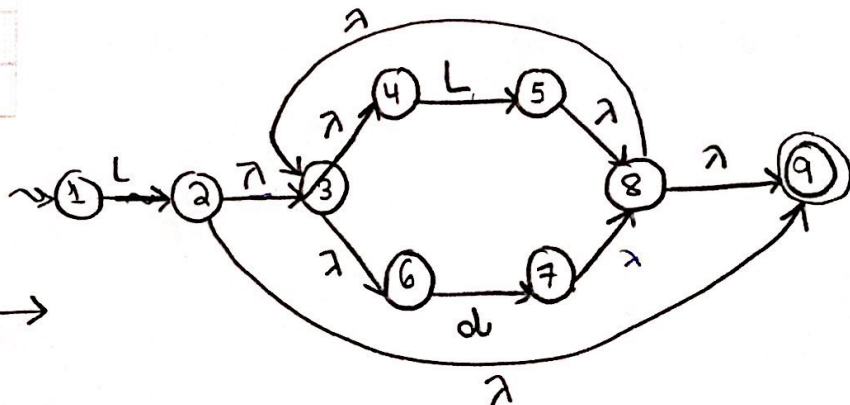
$L(L|d)^*$  "Names"



Which has this transition table

State \ VT	L	d	$\lambda$
1	2		
2			3,9
3			4,6
4	5		
5			8
6		7	
7			8
8			3,9
9			

NDFSA  $L(L|d)^*$  →



State\ VT	L	d	λ
1	2		
2	5	7	3,9,8,6
3	5	7	4,6,8,3,9
4	5		
5	5	7	8,3,9,4,6
6		7	
7	5	7	8,3,9,4,6
8	5	7	3,9,8,4,6
9			

State\ VT	L	d
✓1	2	
✓2	5	7
3	5	7
4	5	
✓5	5	7
6		7
✓7	5	7
8	5	7
9		

State\ VT	L	d
✓1	2	
✓2	5	7
✓5	5	7
✓7	7	7

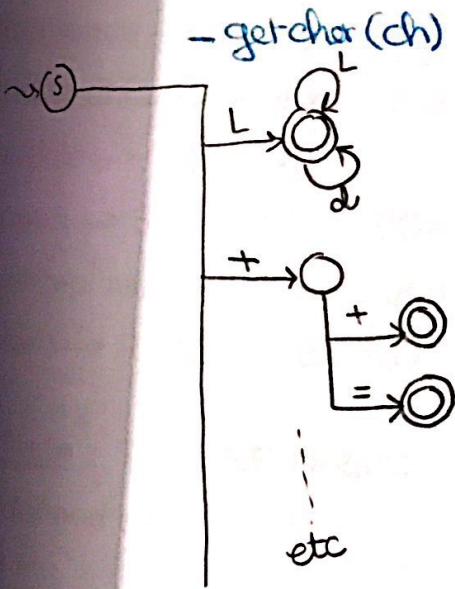
feasible pairs \ VT	L	d
(2,5)	(5,5)	(7,7)
(7,7)	(5,5)	(7,7)
(5,7)	(5,5)	(7,7)

$2 \equiv 5 \equiv 7$



State \ VT	L	d
1	2	-
2	2	2

\*Our Scanner [lexical Analyzer] Algorithm, contains a set of DFSA.



Token .

Name

- +            X = n + 5;
- ++          X ++;
- +=          X + = 5;

# Chapter 5

## Syntax Analysis(Parser)

A **Syntax analyzer** is formally defined as :

An Algorithm that Groups the Set of Tokens Sent by the Scanner to Form **Syntax Structures** Such As Expressions, Statements, Blocks, etc.

Simply, the parser examines if the source code written follows the grammar(production rules) of the language.

The Syntax structure of programming languages and even spoken languages can be expressed in what is called **BNF** notation, which stands for **Bakus Naur Form**.

For example, in spoken English, we can say the following:

sentence --> noun-phrase verb-phrase *↪ First Notation*

noun-phrase --> article noun

article --> THE | A | ...

noun --> STUDENT | BOOK | ...

verb-phrase --> verb noun-phrase

verb --> READS | BUYS | ....

**Note** : The BNF Notation uses different symbols, for example, a sentence is defined as :

< sentence > ::= < noun-phrase > < verb-phrase > *↪ Second notation.*

But this is very cumbersome, so we use the first notation, since it is easier to use. Now, let us derive

a sentence :

sentence --> noun-phrase verb-phrase

--> article noun verb-phrase

--> THE noun verb-phrase

--> THE STUDENT verb-phrase

①

- > THE STUDENT verb noun-phrase
- > THE STUDENT READS noun-phrase
- > THE STUDENT READS article noun
- > THE STUDENT READS A noun
- > THE STUDENT READS A BOOK

! STEP BY STEP

In the same way, the parser tries to **derive** your source program from the starting symbol of the grammar. Let us say we have these

sentences :

THE BOOK BUYS A STUDENT  
THE BOOK WRITES A DISH

*dream*

Syntax-wise, all of these sentences are correct. However, their meaning is not correct, and they are not useful. What differentiates 2 sentences that are grammatically correct is their meaning or their **semantics**. You and I can agree that the meaning of a grammatically correct sentence is not correct, but how does the computer do it?

## Grammar of a programming language:-

A grammar  $G = (V_N, V_T, S, P)$  where:

1.  $V_N$ : A finite set of nonterminals (nonterminals set).
2.  $V_T$ : A finite set of terminals (terminals set).
3.  $S \in V_N$ : The Starting symbol of the grammar. "non-terminal"
4.  $P =$  A set of production rules (productions). <-- Pending  $\Leftrightarrow$  Basically the whole grammar.

Note :

1.  $V_N \cap V_T = \emptyset$ .
2.  $V_N \cup V_T = V$  (the vocabulary of the grammar).

Note : We will use:

1. Uppercase Letters A, B, ..., Z for non-terminals.  $\mathbb{N}$
2. Lowercase Letters a, b, ..., z for terminals.  $\mathbb{N}$
3. Greek letters  $\alpha, \beta, \gamma, \dots$  for strings  $\mathbb{N}$

formed from  $V_N$  OR  $V_T = V$ .

$$V_N \cup V_T$$

eg, if  $V_N = \{S, A, B\}$ ,  $V_T = \{0, 1\}$

then

$$\alpha = A11B$$

$$\beta = S110B$$

$$\gamma = 0010$$

## Productions

Contains at least one non-terminal

1. A Production  $\alpha \rightarrow \beta$  (alpha derives beta) is a rewriting rule such that the occurrence of  $\alpha$  can be substituted by  $\beta$  in any string.

Note that  $\alpha$  must contain at least one

nonterminal from  $V_N$ , i.e.  $\in V_N$ . For

example, Assume we have the string

$\gamma\alpha\sigma$ ,

$$\gamma\alpha\sigma \rightarrow \gamma\beta\sigma$$

(3)

2. A Derivation is a sequence of strings  $\alpha_0, \alpha_1, \alpha_2, \dots, \alpha_n$ , then:

$$\alpha_0 \rightarrow \alpha_1 \rightarrow \alpha_2 \rightarrow \dots \rightarrow \alpha_n$$

Given a grammar G, then:

$L(G)$  = Language Generated By the Grammar.

for example, Given the Grammar,  $G = (\{S, B, C\}, \{a, b, c\}, S, P)$

- P:
- $S \rightarrow aSBC$
  - $S \rightarrow abC$
  - $CB \rightarrow BC$
  - $bB \rightarrow bb$
  - $bC \rightarrow bc$
  - $cC \rightarrow cc$

∴ Contains at least one non-terminal.  
∴ Non-terminal represented by Capital letter.

What is  $L(G)$ ? ← mark it!

Let us do some derivations: "Derive sentences"

$S \rightarrow abC \rightarrow \boxed{abc}$  (all terminals)  $\in L(G)$  ← A sentence

$S \rightarrow aSBC \rightarrow aabCBC \rightarrow aabbcBC \rightarrow$  blocked, so we try another path

$S \rightarrow aSBC \rightarrow aabCBC \rightarrow aabBCC \rightarrow aabbCC \rightarrow aabbcC$

$\rightarrow \boxed{aabbcC} \in L(G)$  ← A sentence

$S \rightarrow aSBC \rightarrow \dots \rightarrow \boxed{aaabbbccc} \in L(G)$  ← A sentence

Therefore,  $L(G) = \{a^n, b^n, c^n \mid n \geq 1\}$

As another Example, we have these productions

$E \rightarrow E+T$  } ← we can write the productions 1 and 2 as a single production  $E \rightarrow E+T \mid T$

$E \rightarrow T$  }

$T \rightarrow T * F$  }

$T \rightarrow F$  }

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

$F \rightarrow (E)$  }

$F \rightarrow n$  }

← we can write the productions 5 and 6 as a single production  $F \rightarrow (E) \mid n$

$$\begin{cases} V_N = \{E, T, F\} \\ V_T = \{+, *, (, ), n\} \\ S = E \end{cases}$$

Let us follow through some derivations

$E \rightarrow T \rightarrow F \rightarrow n \in L(E)$

$E \rightarrow E+T \rightarrow T+T \rightarrow T+F \rightarrow T+n \rightarrow F+n \rightarrow n+n \in L(E)$

$E \rightarrow E+T \rightarrow T+T \rightarrow F+T \rightarrow n+T \rightarrow n+F \rightarrow n+(E) \rightarrow n+(T) \rightarrow n+(T * F) \rightarrow$

$n+(F * F) \rightarrow n+(n * F) \rightarrow n+(n * n) \in L(E)$

$L(G)$  :- The set of all arithmetic operations ~~derived~~ with + & \* operations. (4)



⚠ The Algorithm we used in the Derivation steps in **Backtracking** "Try & Error" and the Language is **Non-Deterministic** because it has options!



Therefore,  $L(G) = \{\text{Any arithmetic expression with } * \text{ and } + \text{ operations}\}$ ,  $n$  is an operand here.

Note that, if we add the productions

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

$T \rightarrow T * F \mid T / F \mid T \% F$

We would have a language to express all arithmetic expressions with  $(*, \backslash, +, -)$  operations.

Let us Take another Example (Tokens between double quotes are terminals)

Program  $\rightarrow$  block "#"

block  $\rightarrow$  "{" stmt-List "}"

stmt-List  $\rightarrow$  statement ";" stmt-List  $\mid \lambda$

statement  $\rightarrow$  if-stmt  $\mid$  while-stmt  $\mid$  read-stmt  $\mid$  write-stmt  $\mid$  assignment-stmt  $\mid$  block if-

if-stmt  $\rightarrow$  "if" condition....

while-stmt  $\rightarrow$  "while" condition.....

....

....

read-stmt  $\rightarrow$  "read"

write-stmt  $\rightarrow$  "write"

$V_N = \{\text{Program, block, stmt-List, statement, if-stmt, while-stmt, read-stmt, write-stmt, assignment-stmt}\}$  **THOSE TOKENS DON'T APPEAR IN THE SOURCE CODE.**

$V_T = \{\text{"{", "}", "\#", ";", "if", "while", "read", "write"}\}$  **THOSE TOKENS APPEAR IN THE SOURCE CODE.**

Let us Follow through some derivations :

Program  $\rightarrow$  block #  $\rightarrow$  { stmt-list } #  $\rightarrow$  {  $\lambda$  } #

Program  $\rightarrow$  block #  $\rightarrow$  { stmt-list } #  $\rightarrow$  { statement ; stmt-list } #  $\rightarrow$

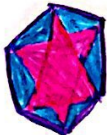
{ statement ; statement ; stmt-list } #  $\rightarrow$

{ statement ; statement ;  $\lambda$  } #  $\rightarrow$  { statement ; statement ; } #  $\rightarrow$  { READ-statement ;

statement ; } #  $\rightarrow$  { READ ; statement ; } #  $\rightarrow$

{ READ ; write-statement ; } #  $\rightarrow$

{ READ ; WRITE ; } #



5

The language of this language is defined as

$L(G) = \{\text{Set of all programs that can be written in this language}\}$   
*Set of all programs derived from this grammar.*

This is only a simple example, of a simple language. For something more complex such as C or Pascal, there are hundreds of productions.

### Algorithms for Derivation

A Leftmost derivation is a derivation in which we replace the leftmost nonterminal in each derivation step.

A Rightmost derivation is a derivation in which we replace the rightmost nonterminal in each derivation step.

DEFINITIONS ←

For example, given

the grammar

$V \rightarrow SR\$$  *Stop Symbol*

$S \rightarrow + \mid - \mid \lambda$

$R \rightarrow \cdot dN \mid dN \cdot N$

$N \rightarrow dN \mid \lambda$

$V_N = \{V, R, S, N\}$

$V_T = \{+, -, \cdot, d, \$\}$

Let us follow through on the leftmost derivation

$V \rightarrow SR\$ \rightarrow -R\$ \rightarrow -dN.N\$ \rightarrow -ddN.N\$ \rightarrow -dddN.N\$ \rightarrow -ddd.N\$ \rightarrow -ddd.d\$ \rightarrow -ddd.d\$ \leftarrow$  A sentence.

Let us follow through on the rightmost derivation

$V \rightarrow SR\$ \rightarrow SdN.N\$ \rightarrow SdN.dN\$ \rightarrow SdN.d\$ \rightarrow sddN.d\$ \rightarrow sdddN.d\$ \rightarrow Sddd.d\$ \rightarrow -ddd.d\$ \leftarrow$  A sentence.

### Derivation Trees

A Derivation Tree is a Tree that displays the derivation of some sentence in the language. For example, let us look at the tree for the previous example

6

## Examples on Leftmost & Rightmost Derivations

$$E \rightarrow E + T \quad | \quad T$$

$$T \rightarrow T * F \quad | \quad F$$

$$F \rightarrow (E) \quad | \quad n$$

\*1

$$\left\{ \begin{array}{l} E \xrightarrow{Lm} E + T \xrightarrow{Lm} T + T \xrightarrow{Lm} F + T \xrightarrow{Lm} n + T \xrightarrow{Lm} n + F \\ \xrightarrow{Lm} n + n \in L(G) \end{array} \right.$$

LEFTMOST

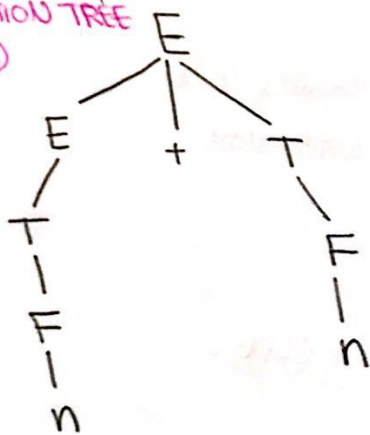
\*2

$$\left\{ \begin{array}{l} F \xrightarrow{RM} E + T \xrightarrow{RM} E + F \xrightarrow{RM} E + n \xrightarrow{RM} T + n \\ \xrightarrow{RM} F + n \xrightarrow{RM} n + n \in L(G) \end{array} \right.$$

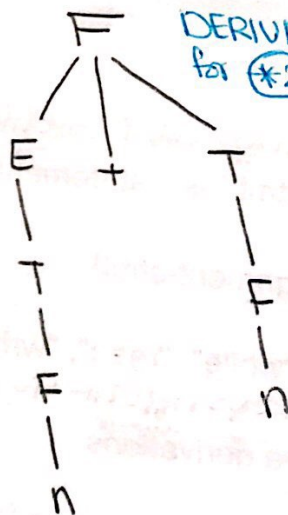
RIGHTMOST

**⚠ Note:** If we traverse the derivation tree **INORDER** considering only the leaves, we obtain the sentence.

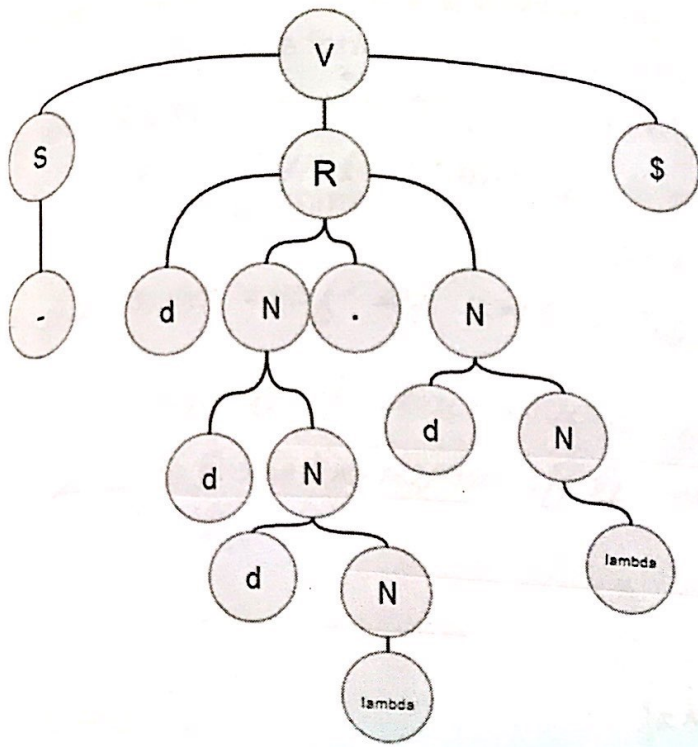
DERIVATION TREE for \*1



DERIVATION TREE for \*2



**⚠** Same Derivation TREE for the Sentence  $n+n$



Note that if we traverse the tree in order, recording **only** the leaves, we obtain the sentence.

## Classes of Grammars

According to Chomsky, There are 4 classes of grammars :

1. **Unrestricted Grammars** : No restrictions whatsoever except the restriction by definition that the left side of the production contains at least one nonterminal from  $V_N$ . This grammar is not practical and we cannot work with it.

*its very hard to work with*

2. **Context-Sensitive Grammars** : For each production  $\alpha \rightarrow \beta$ ,  $|\alpha| \leq |\beta|$ , i.e., the **length of alpha( $\alpha$ )** is less than or equal to the **length of Beta( $\beta$ )**. This means that in this class of grammar, there are no  $\lambda$  productions in the form  $A \rightarrow \lambda$ , since  $|\lambda| = 0$  and  $A \geq 1$ .

$(A \rightarrow \lambda) \times$

3. **Context-Free Grammar(CFG)** : Each production in this grammar class is of the form  $A \rightarrow \alpha$ , where  $A \in V_N$  and  $\alpha \in V^*$

that is to say, the left hand side is **only** one nonterminal.

This is the most important class of grammar. Most programming languages structures are context-free. We will mostly be working with this class of grammar. Most of the examples we have taken are CFG.

(7)

4. **Regular Grammar (Regular Expressions)** : Each production in this grammar class is of the form  $A \rightarrow aB$  or  $A \rightarrow a$ , where  $A, B \in VN$  and  $a \in VT$ , with the exception of  $S \rightarrow \lambda$

↳ Starting Symbol.

5. For example, given the grammar:

$A \rightarrow aA$   
 $A \rightarrow a$

Therefore, we get  $L(G) = a^+$

However, adding the production  $A \rightarrow \lambda$

Results in the grammar

$L(G) = a^*$

⚠ Regular Grammar is a subset of Context Free Grammar

## Parsing Techniques

There are 2 main parsing techniques used by a compiler.

### Top-Down Parsing

① leftmost derivation until it derives the sentence and deadlock.  
 ③ Major question: which production the parser must select for the root derivation step.

In Top-Down Parsing, the parser builds the derivation tree from the root ( $S$ : the starting symbol) down to the leaves (sentence).

In Simple words, the parser tries to derive the sentence using **leftmost**!

derivation. For example, say we have this grammar:  $V \rightarrow SR\$$

$S \rightarrow + | - | \lambda$

$R \rightarrow \cdot dN | dN \cdot N$

$L(G) = [+|-] \left\{ \begin{matrix} d^+ \\ d^+ \cdot d^+ \\ d^+ \cdot d^+ \end{matrix} \right\}$

$N \rightarrow dN | \lambda$

Let us examine if the sentence  $dd.d\$$  if it is derived from this grammar.

$V \rightarrow SR\$ \rightarrow \lambda RS\$ \rightarrow dN.N\$ \rightarrow ddN.N\$ \rightarrow$

$dd.N\$ \rightarrow dd.dN\$ \rightarrow dd.d\$$

**Problem.** The Parser does not know which production it should select in each derivation statement. We will learn how to solve these issues later in the course.

# DEFENITION

## PARSER [Syntax Analyzer]:

is an algorithm which takes your source code and tries to :-

- (1) Build the derivation tree for your source code
- OR  
(2) Derive your source code from the production rules of the grammar. (Using left most Derivation)

### TOP-DOWN PARSING PROBLEM

$E \rightarrow T$   
 $T \rightarrow F$   
 $F \rightarrow a$  the only sentence.

⚠ Not having multiple productions means that we only have one sentence.

## Bottom-Up Parsing

In Bottom-Up Parsing, the parser builds the derivation tree from the leaves (sentence) up to the root (S : Starting Symbol). This type of tree, built from the leaves to the root, is called a B-Tree.

In Simple words, the parser starts with the given sentence, does reduction (opposite of derivation) steps, until the starting symbol is reached.

Note that the string  $\lambda$  is present everywhere in the string,

and we can use it wherever we like. Let us follow the

reduction of the example given above.

$+dd.d\$ \rightarrow +dd\lambda.d\$ \rightarrow +ddN.d\$ \rightarrow +dN.d\$ \rightarrow +dN.d\lambda\$ \rightarrow$

$+dN.dN\$ \rightarrow +dN.N\$ \rightarrow +R\$ \rightarrow SR\$ \rightarrow V$  Which means that the

sentence is in the grammar.

Note that we can run into deadlocks here. say we took this path instead :

$+dd.d\$ \rightarrow +dd\lambda.d\$ \rightarrow +ddN.d\$ \rightarrow +dN.d\$ \rightarrow +dN.d\lambda\$ \rightarrow +dN.dN\$ \rightarrow$

$+dNR\$ \rightarrow +NR\$ \rightarrow SNR\$ \rightarrow$  Deadlock This technique also has a major

**problem** : Which substring should we select to reduce in each reduction

step?

how do we solve this?

**Ambiguity**  $\rightarrow$  Ambiguous  $\equiv$  Not Deterministic  $\equiv$  If there are choices.

$\rightarrow$  non-deterministic Algorithm.

Given the following grammar :

$num \rightarrow num d$

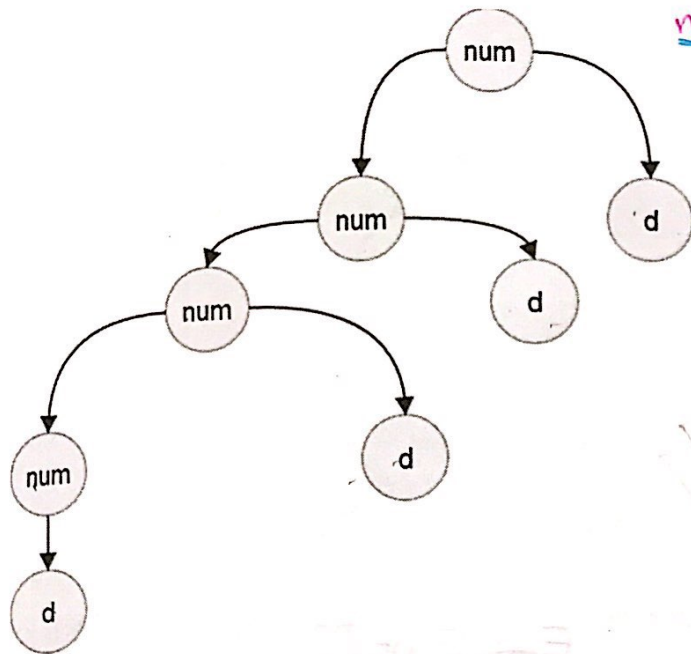
$num \rightarrow d$

} Similar to regular grammar.

Let us draw the derivation tree for the sentence  $dddd$

9





$\text{num} \rightarrow \text{num } d \rightarrow \text{num } d d$   
 $\rightarrow \text{num } d d d \rightarrow d d d d \in L(G)$

⚠ one sentence  $\rightarrow$  one derivation tree  
 $\hookrightarrow$  Unambiguous

Question : is there another derivation tree that represents the sentence? The answer is **no**.

If there is only one derivation tree representing the sentence, this means there is only one way to derive the sentence. Based on this, we can say that

$\hookrightarrow$  Unambiguous.

**Def:** A Grammar G is said to be **ambiguous** if there is <sup>at-least</sup> one sentence with more than one derivation tree. That is, there is more than one way to derive the sentence.

This means that our algorithm is **non-deterministic**.

EX: given the grammar:

$E \rightarrow E + E$

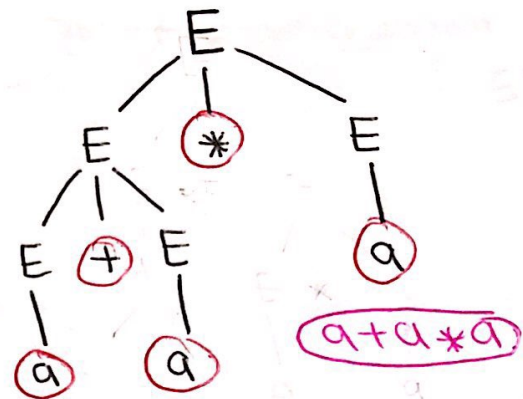
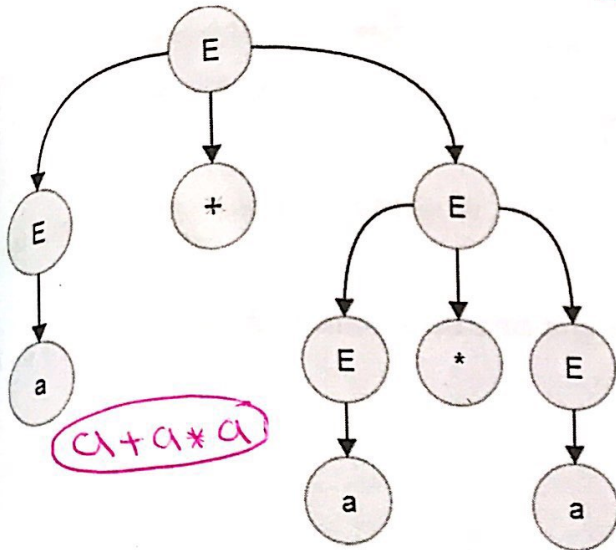
$E \rightarrow E * E$

$E \rightarrow (E) | a$

Take the sentence :  $a + a * a$

Let us draw the derivation tree

$E \rightarrow E + E \rightarrow a + E \Rightarrow a + E * E \rightarrow a + a * E \rightarrow a + a * a \in L(G)$



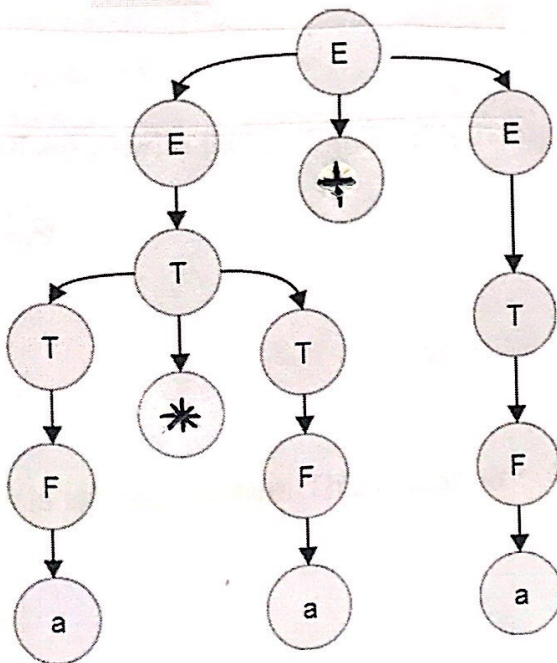
⚠ Same sentence, different derivation trees. "ambiguous"

Due to the fact that we have 2 trees that give the same result, we can say that this grammar is ambiguous. In this case, to enforce the associativity rule, this grammar can be re-written as :

$E \rightarrow E + E \mid T$   
 $T \rightarrow T * T \mid F$   
 $F \rightarrow (E) \mid a$

} - ambiguity  
 - Precedency  
 - associativity

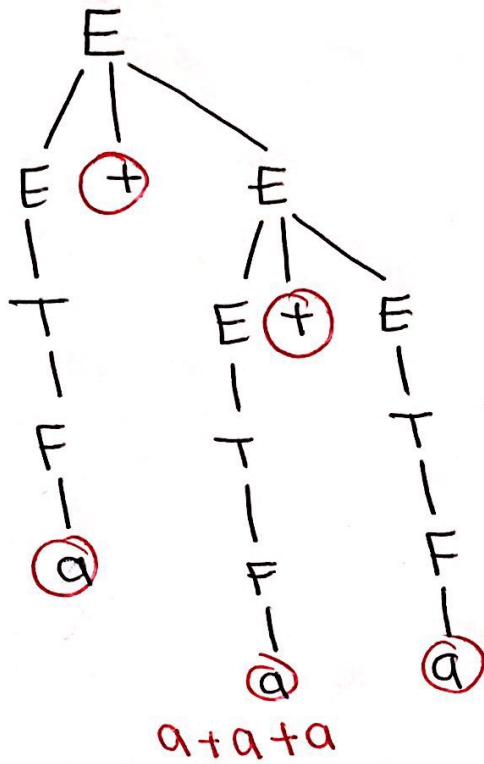
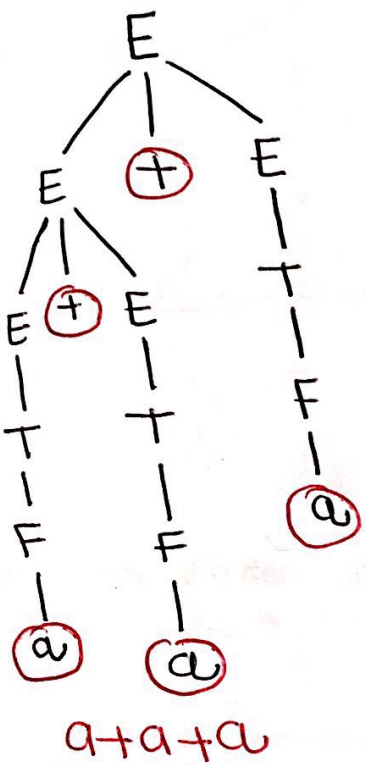
Now, Take the sentence  $a + a * a$  and find the derivation tree now.



There is only possible derivation trees now. This solves the associativity problem with + and \* of the grammar before with the operations.

But Let us say we have the sentence : a + a + a

Let us try to find the derivation tree and any alternative trees.



We can see here that there is more than 1 derivation tree, and the language is still ambiguous.

We can solve this if we rewrite the grammar with the **left-associative rule**

- $E \rightarrow E + T \mid T$
- $T \rightarrow T * F \mid F$
- $F \rightarrow (E) \mid a$

The grammar now is left-associative. This grammar solves the problems of :

- ambiguity.
- precedence.
- associativity.

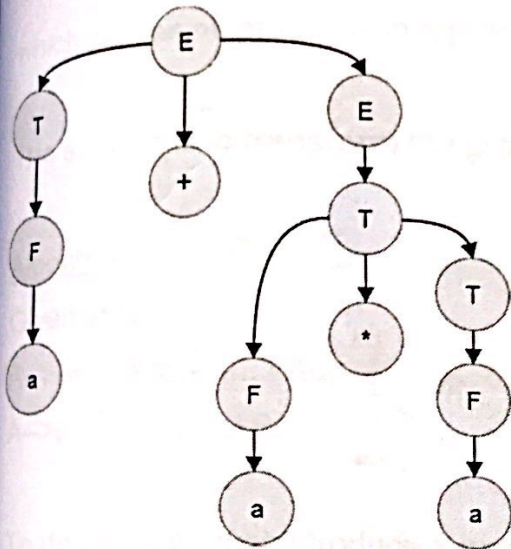
Let us try rewriting it with the right-associative rule

$E \rightarrow T + E \mid T$

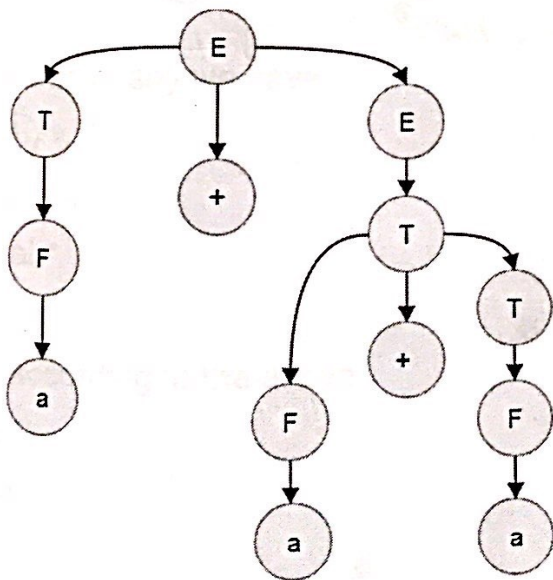
$T \rightarrow F * T \mid F$

$F \rightarrow (E) \mid a$

Let us try creating the derivation tree of  $a + a * a$



Now let us check the derivation tree of  $a + a + a$



This new grammar is not ambiguous, however it does not solve the fact that associativity issue according to our standard.

# left-recursive Grammar

This causes problems when it comes to Top-down parsing techniques (see why later).

A grammar is said to be left recursive if there is a production of the form:

$$A \rightarrow A\alpha$$

Conversely, a grammar is right-recursive if there is a production of the form:

$$A \rightarrow \alpha A$$

which causes no problems in top-down parsing.

The solution is to transform the grammar to a grammar which is not left-recursive.

## Algorithm.

Given that:

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_n$$

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

Transform

To do this, we must introduce a new non-terminal, say  $A'$ .

The grammar now becomes:

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

$$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_n A' \mid \lambda$$

TO IMPORTANT!

- INPUT: Left-Recursive Grammar.
- OUTPUT: An equivalent non Left-Recursive Grammar

For example, say we have

$$A \rightarrow Ab \rightarrow \alpha_1$$

$$A \rightarrow a \rightarrow \beta_1$$

$$L(G) = ab^*$$

Then according to the above

$$A \rightarrow aA'$$

$$A' \rightarrow bA' \mid \lambda$$

which results in the same grammar.

Let us apply this to the grammar:

$$E \rightarrow E + T \mid T \beta_1 \quad \bullet \text{Left-Recursion.}$$

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

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

$$E \rightarrow E + T \rightarrow E \rightarrow E + T \rightarrow \alpha_1$$

$$E \rightarrow T \rightarrow \beta_1$$

Then the new grammar :

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid a$

This grammar is now **perfect**. It solves all our ambiguity issues, and this is a grammar we can use to construct the production rules for our programming language.

Another ambiguity in programming languages is the **if...else** statement.

stmt  $\rightarrow$  if-stmt | while-stmt | ....

if-stmt  $\rightarrow$  IF condition stmt

if-stmt  $\rightarrow$  IF condition stmt **ELSE** stmt

condition  $\rightarrow$  C

stmt  $\rightarrow$  S

RESERVED WORDS

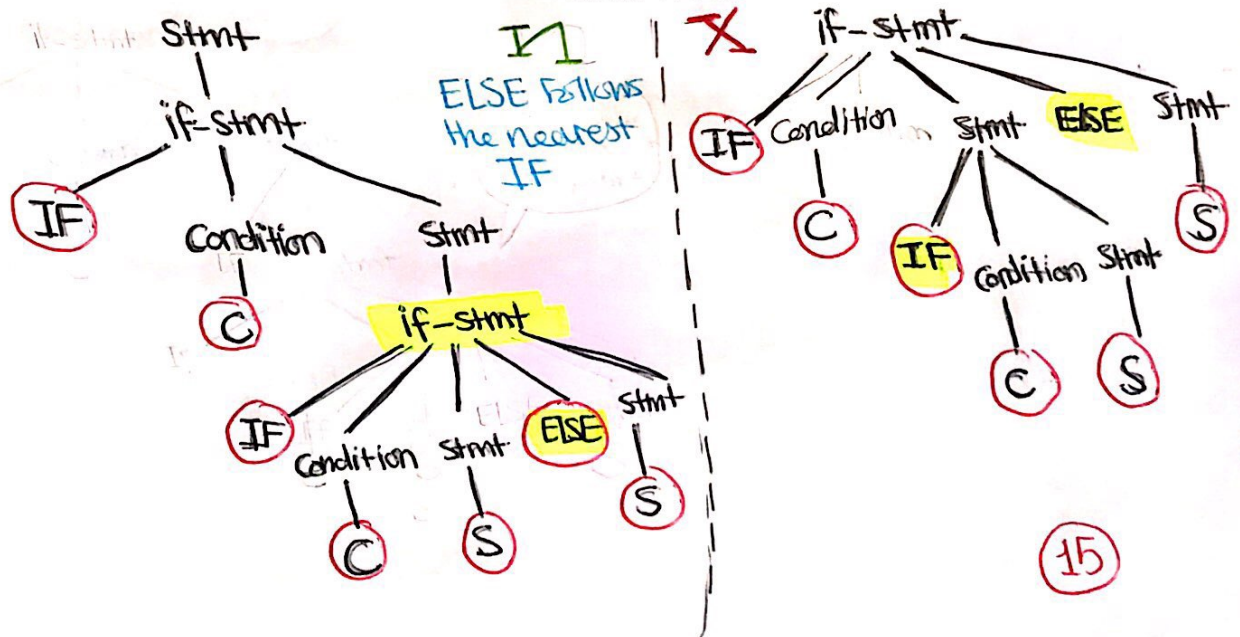
This grammar is **ambiguous**.

Let us take the following nested if...else statement :

```
IF C
  IF C
    S
  ELSE
    S
```

This statement results in 2 derivations trees.

Draw the two derivation Trees



The first results in the ELSE belonging to the first IF , while the second results in the ELSE belonging to the second IF .

The second tree is the correct one since we know that the ELSE statement follows the nearest IF .

But how can the compiler behaves in this case?

There are a bunch of solutions to this problem:

1. Add a delimiter to the IF statement, such as ENDIF or END or FI to the end of the statement, resulting in these productions :

if-stmt --> IF condition stmt **ENDIF**

if-stmt --> IF condition stmt ELSE stmt **ENDIF**

Resulting in this statement :

...

IF C

IF C

S

ELSE

S

ENDIF

ENDIF

...

The grammar is now unambiguous, since we have to clearly state when an "IF" statement ends.

However, this is not a pretty solution, and is extra work for both the programmer and compiler, and result in less readable code.

2. Make the compiler **always** prefers to shift the ELSE when it sees the ELSE in the source code.

## Left Factoring

Consider the productions :

$$\left. \begin{array}{l} A \rightarrow \alpha\beta \\ A \rightarrow \alpha\gamma \end{array} \right\}$$

Note how the **first part** of the productions is the same. This grammar can be transformed by introducing a new **non-terminal B**,

So what happens now is:

$$\left. \begin{array}{l} A \rightarrow \alpha B \\ B \rightarrow \beta \mid \gamma \end{array} \right\}$$

For our grammar, this results in

if-stmt  $\rightarrow$  IF condition stmt

if-stmt  $\rightarrow$  IF condition stmt ELSE stmt

becomes:

if-stmt  $\rightarrow$  IF condition stmt else-part

else-part  $\rightarrow$  ELSE stmt  $\mid \lambda$

Does this solve the ambiguity? No, but it helps in removing choices, since the if-stmt is now one production. If we look at the

statement :

```
IF C
  IF C
    S
  ELSE
    S
```

It still has 2 derivation trees

## Extended BNF Notation

So far, we have been using **BNF Notation**(Production rules) to express languages.

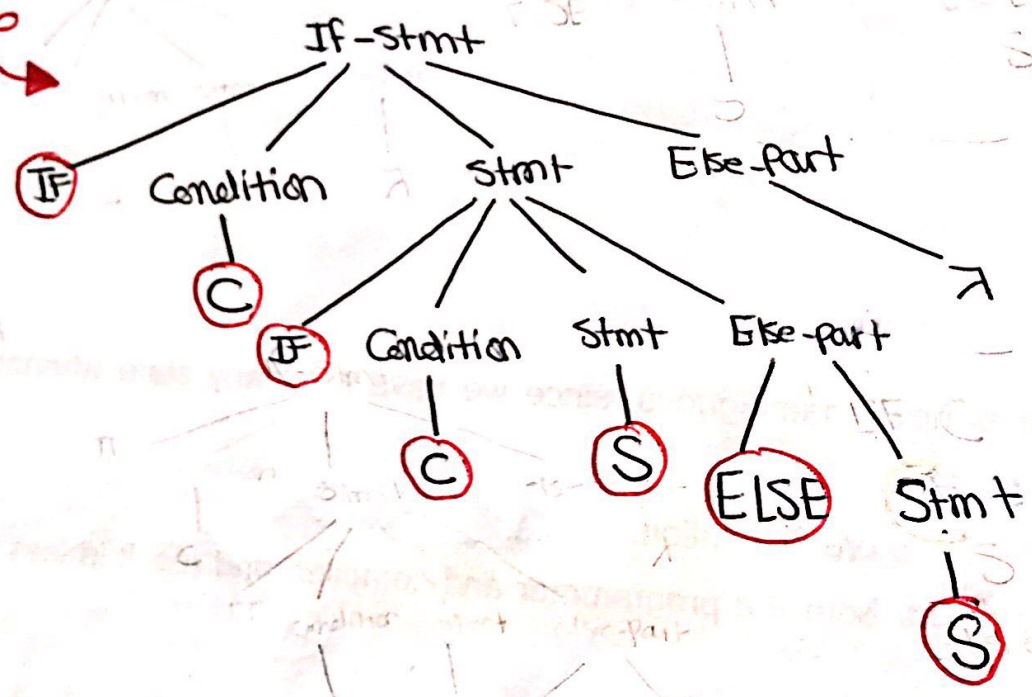
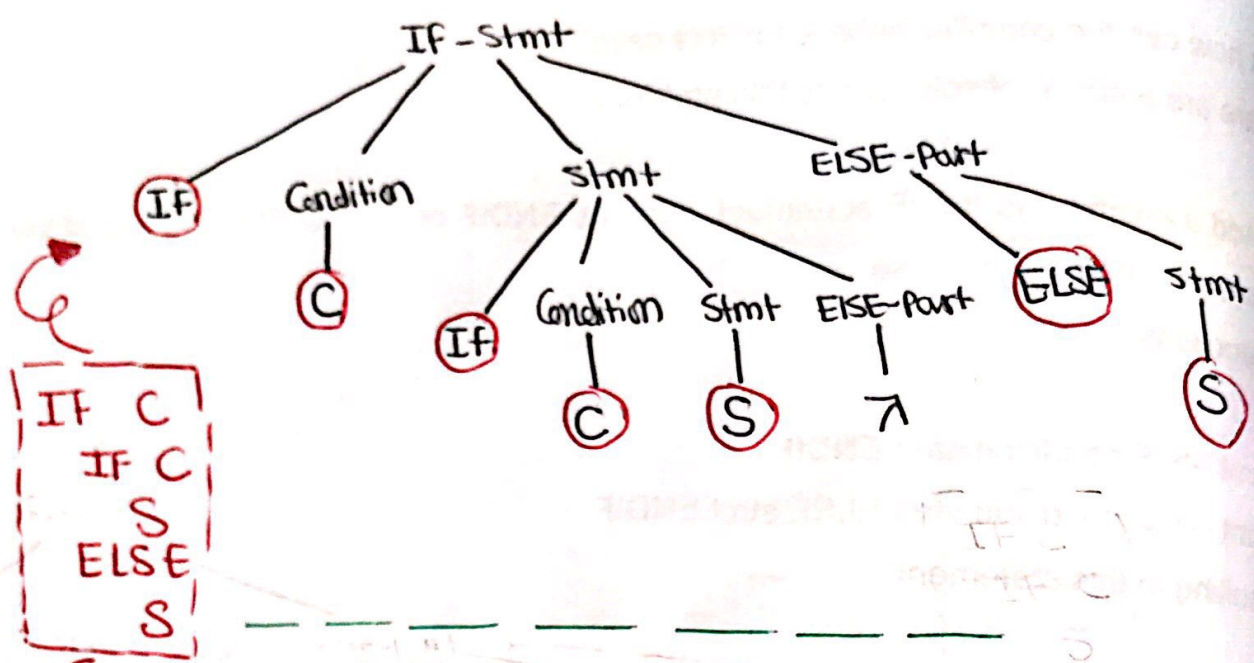
However, there is another form to

Express a language, which is **Extended BNF Notation**

if there is repetition in the grammar, say in the example of the grammar

$$E \rightarrow E + T \mid T$$
$$T \rightarrow T * F \mid F$$
$$F \rightarrow (E) \mid a$$





which can give us a derivation in the form of

$$E \rightarrow E + T \rightarrow E + T + T \rightarrow E + T + T + T \rightarrow T + T + T + T \dots + T$$

or in the same line,

$$T \rightarrow T * F \rightarrow T * F * F \rightarrow T * F * F * F \rightarrow T * F * F * F \dots F$$

We can express this grammar as :

$$E \rightarrow T \{ + T \}$$

$$T \rightarrow F \{ * F \}$$

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

We know that  $[x]$  means that we take  $x$  0 or 1 time only.

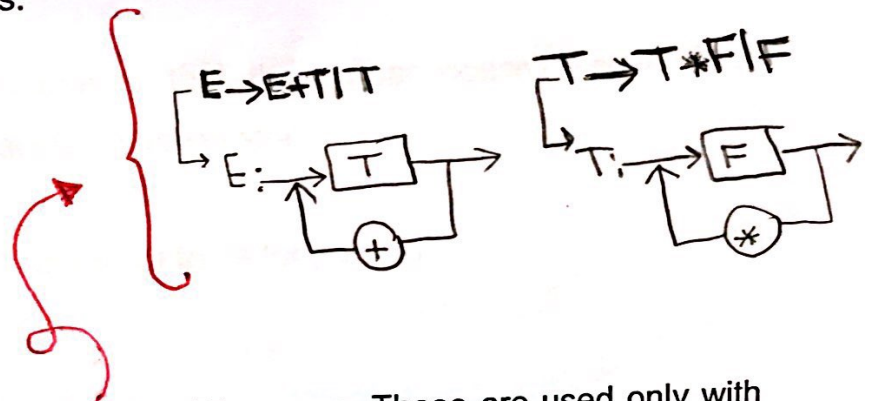
However,  $\{x\}$  means we take  $x$  zero or any number of times. This is equivalent to  $(x)^*$

We can also express this grammar as:

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

$$T \rightarrow F (* F)^*$$

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



### Syntax Diagrams

Another way to express languages are Syntax Diagrams. These are used only with Extended-BNF notation.

A square shape represents a nonterminal and an oval shape represents a terminal.

DRAW

# Parsing Techniques

Recall : The parser is an algorithm which accepts or rejects a sentence in the programming language.

Recall : There are 2 kinds of parsers :

1. **Top-Down Parsers** : In This parsing technique, The parser starts with S using leftmost derivation to derive the sentence.

The **Major problem** with this parsing technique is that the parser doesn't know which production it should select in each derivation step.

2. **Bottom-Up Parsers** : The parser in this parsing technique starts from the sentence, doing reduction steps, until it reaches the starting symbol S of the grammar.

The **Major problem** with this technique is that the parser doesn't know which substring the parser should select in each reduction step.

In Top-Down parsing, we have 2 available algorithms for parsing :

- 1. Recursive Descent Parsing.
- 2. LL(1) Predictive Parsing.

*\* We can build deterministic parser for the grammar.  
\* These subsets are powerful enough to define any programming language.*

In Bottom-Up parsing, we have 2 available algorithms for parsing :

- 1. LR Parsers.
- 2. Operator Precedence Parsers --> Uses matrix manipulation.

*\* we can build deterministic parser for the grammar*

Before we continue, we need to define a few functions

## The FIRST() Function

Given a string  $\alpha \in V^*$ , then

$$\text{FIRST}(\alpha) = \{ a \mid \alpha \xrightarrow{*} aw, a \in V_T, w \in V^* \}$$

in addition, if  $\alpha \xrightarrow{*} \lambda$ , then we add  $\lambda$  to  $\text{FIRST}(\alpha)$ , that is  $\lambda \in \text{FIRST}(\alpha)$ .

That is to say,  $\text{FIRST}(\alpha)$  = Set of all **terminals** that may begin strings derived from  $\alpha$ .

*any string*

*zero or more derivations.*

*$V_T \cup V_N$  including  $\lambda$ .*



## Rules To Compute FIRST() and FOLLOW() Sets

- X)
- $FIRST(\lambda) = \{\lambda\}$ .
  - $FIRST(a) = \{a\}$ .  $a \in V_T$
  - $FIRST(a\alpha) = \{a\}$ .  $a \in V_T, \alpha \in V^*$
  - $FIRST(XY) = FIRST(FIRST(X).FIRST(Y))$  OR  
*preferred*  $FIRST(X.FIRST(Y))$  OR  
 $FIRST(FIRST(X).Y)$ .
  - Given the production  $A \rightarrow \alpha X \beta$ , Then :
    - $FIRST(\beta) \subset FOLLOW(X)$  if  $\beta \neq \lambda$ .
    - $FOLLOW(A) \subset FOLLOW(X)$  if  $\beta = \lambda$ .

Note that the FIRST() and FOLLOW() sets are made of **terminals only**

### Notes :

- $\lambda$  may appear in FIRST() but it doesn't appear in FOLLOW(). We will see this when we define augmented grammars.
- Generally, we start computing the FIRST() from bottom to top, But FOLLOW() from top to bottom.
- When we compute FOLLOW(X), we search for X in the right side of any production.

## Augmented Grammars

### DEFENITION

Given the grammar  $G = (V_N, V_T, S, P)$ , then the augmented grammar  $G' = (V_N', V_T', S', P')$  can be obtained from G as follows:

- $V_N' = V_N \cup \{S'\}$ .
- $V_T' = V_T \cup \{\$ \}$ . *Stopping Symbol*
- $S'$  = new starting point.
- $P' = P \cup \{S' \rightarrow S\$ \}$

For example :

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

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

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

(21)

Becomes :

$$G \rightarrow E\$$$

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

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

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

This is because we want to create a FOLLOW() set for S.

Example 1 :

$$S' \rightarrow S\$$$

$$S \rightarrow AB$$

$$A \rightarrow a \mid \lambda$$

$$B \rightarrow b \mid \lambda$$

Let us compute the FIRST() sets for this grammar :

$$\text{FIRST}(A) = \{a, \lambda\}$$

$$\text{FIRST}(B) = \{b, \lambda\}$$

$$\text{FIRST}(S) = \text{FIRST}(AB) = \text{FIRST}(\text{FIRST}(A).\text{FIRST}(B))$$

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

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

$$= \{a, b, \lambda\}$$

$$\text{FIRST}(S') = \text{FIRST}(S\$) = \text{FIRST}(\text{FIRST}(S).\text{FIRST}(\$))$$

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

$$= \{a, b, \$\}$$

⚠️ we don't calculate Follow(S') ?!

Now Let us compute the FOLLOW() sets for this grammar :

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

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

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

$$\text{Follow}(A) = \text{First}(B) \text{ " } B \neq \lambda \text{ " } \rightarrow \text{Rules Page \#21}$$

$$\text{First}(b) = \{b, \lambda\}$$

λ doesn't appear in Follow(C)

$$\text{Follow}(A) = \text{First}(B) \cup \text{Follow}(S) \text{ " } \text{First}(B) \text{ contains } \lambda \text{ "}$$

$$= \{b\} \cup \{\$\}$$

$$= \{b, \$\}$$

Rules page #21

22

Example 2 :

$$S' \rightarrow S\$$$

$$S \rightarrow aAc b$$

$$S \rightarrow \underline{A}bc$$

$$A \rightarrow b \mid c \mid \lambda$$

Let us take the FIRST() for this grammar :

$$\text{FIRST}(A) = \{b, c, \lambda\}$$

$$\text{FIRST}(S) = \text{FIRST}(aAc b) \cup \text{FIRST}(Abc) = \{a, \lambda\} \cup \{b, c\}$$

$$= \{a, b, c\}$$

$$\text{FIRST}(S') = \text{FIRST}(SS) = \text{FIRST}(\text{FIRST}(S) \cdot \text{FIRST}(S))$$

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

$$= \{a, b, c\}$$

Now let us take the FOLLOW() :

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

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

Example 3:

$$G \rightarrow E\$$$

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

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

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

Let us calculate FIRST() :

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

$$\text{FIRST}(T) = \text{FIRST}(T * F) \cup \text{FIRST}(F) = \text{FIRST}(T * F) \cup \{(, a\}$$

$$= \{(, a\} \text{ (Because every } T \text{ will eventually become an } F)$$

$$\text{FIRST}(E) = \text{FIRST}(E + T) \cup \text{FIRST}(T) = \{(, a\} \cup \{(, a\}$$

$$= \{(, a\}$$

$$\text{FIRST}(G) = \text{FIRST}(E\$) = \{(, a\}$$

Now let us Calculate FOLLOW() :

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

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

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

But what makes all this so important?

Well, All of the parsing techniques we are going to depend will heavily on FIRST() and FOLLOW().

23

$$\text{FIRST}(\{a|b|c\} \cup \{\$\}) = \{a|b|c\}$$

↳ since this doesn't contain  $\epsilon$   
we ignore  $\$$

$$\text{FIRST}(\{a|b|c|\epsilon\} \cup \{\$\}) = \{a|b|c|\$\}$$

↳  $\epsilon$  exists here so we  
take  $\$$  in consideration



## Recursive Descent Parsing

This parsing technique is better with Extended BNF

Recursive Descent Parsing is very simple. It works like this :

Divide the grammar into primitive/simple components

1- For the token "a" :

```
if(token == "a")
    get-next()
else
    report-error()
```

⚠ Disadvantage:  
Slow because it depends on function calling.

2- For  $X = \alpha_1 \alpha_2 \dots \alpha_n$  :  
 $X \rightarrow \alpha_1 \alpha_2 \dots \alpha_n$

Code(X):

```
{
    Code( $\alpha_1$ );
    Code( $\alpha_2$ );
    .
    .
    Code( $\alpha_n$ );
}
```

3- For  $X = \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$ , If none of the  $\alpha_i$ 's =  $\lambda$   
 $X \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$

Code(X):

```
{
    If (token  $\in$  FIRST( $\alpha_1$ ))
        Code( $\alpha_1$ );
    Else
        If (token  $\in$  FIRST( $\alpha_2$ ))
            Code( $\alpha_2$ );
        Else
            .
            .
            Else
                If (token  $\in$  FIRST( $\alpha_n$ ))
                    Code( $\alpha_n$ );
```

TOKEN DOESN'T EXIST

Else

Report-error();

SYNTAX ERROR!

(24)

4- For  $X = \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n = \lambda$ , If one of the  $\alpha_i$ 's =  $\lambda$ , say  $\alpha_n = \lambda$

Code(X):

```
{ If (token  $\in$  FIRST( $\alpha_1$ ))
  Code( $\alpha_1$ );
  Else
    If (token  $\in$  FIRST( $\alpha_2$ ))
      Code( $\alpha_2$ );
    Else
      .
      .
    Else
      If (token  $\in$  FIRST( $\alpha_{n-1}$ ))
        Code( $\alpha_{n-1}$ );
      Else
        If (token is not  $\in$  FOLLOW(X))
          Report-error();
}
```

5- For  $X = \alpha^*$

Code(X):

```
While (token  $\in$  FIRST( $\alpha$ ))
  Code( $\alpha$ );
```

Notes :

1. Every nonterminal has a code(a function).
2.  $S'$  in augmented grammar is represented by the function "main".
3. We only start with calling "get-token" in function "main".

Example:

```
G  $\rightarrow$  E$
E  $\rightarrow$  T( + T ) *
T  $\rightarrow$  F( * F ) *
F  $\rightarrow$  ( E ) | a
```

(25)

```

main() { //represents G
  get-token; only in the main.
  call E();
  if(token!="$")
    Error;
  else
}

```

```

function E(){
  call T();
  while(token == "+"){
    get-token();
    call T()
  }
}

```

// E → T (+ T)\*

⚠ Here we don't need get-token because we already have a token.



While(token ∈ First(α))  
 $\alpha = +T$  // First( $\alpha$ ) = +

```

function T(){ //T--> F (* F)*
  call F();
  while(token == "*"){
    get-token();
    call F();
  }
}

```

```

function F(){ //F--> ( E ) | a
  if(token == "(")
  {
    get-token();
    call E();
    if(token == ")")
      get-token();
    else
      ERROR;
  }
  Else
  if(token=="a")
    get-token();
  else
    ERROR;
}

```

Note that ERROR is a function we should write.

Example:

Given the grammar :

Program  $\rightarrow$  body .

body  $\rightarrow$  Begin stmt ( ; stmt)\* End

stmt  $\rightarrow$  Read | Write | body |  $\lambda$

and

VN = { Program, body, stmt, block}

VT = { ., Begin, ;, End, Read, Write}

examples of programs of this language would be:

Begin

Read;

Write;

Read;

Write;

End.

Or

Begin

Read;

End.

Or

Begin

Read;

Begin

Read;

Write;

End.

Write;

End.

Or

Begin;

;

;

;

;

;

End.

*\$/Stop Symbol.*

Let us write the recursive descent code for this programming language.

```
main(){  
  get-token();  
  call body();  
  if(token != "."){  
    ERROR;  
  }  
}
```

```
else {  
  SUCCESS;  
}
```

```
function body()  
{  
  if(token == "Begin")  
  {  
    get-token();  
    call stmt();  
    while(token == ";")  
    {  
      get-token();  
      call stmt();  
    }  
  }  
  if(token == "End")  
  get-token();  
  else  
  ERROR;  
}
```

```
else  
  ERROR;  
}
```

```
function stmt()  
{
```

```
  if(token == "Read")  
  get-token();
```

```
  else if (token == "Write")  
  get-token();
```

```
  else if(token == "Begin")  
  call body();
```

```
  else if ( token != ;" token != "End" )
```

```
  ERROR();  
}
```

## LL(1) Parsing

This Parsing method is a **table-driven** parsing method. The LL(1) parsing table selects which production to choose for the next derivation step.

### Formal Definition of LL(1)

The **Formal definition of LL(1)** grammars is given by :

Given the Productions :

$A \rightarrow \alpha_1$

$A \rightarrow \alpha_2$

$A \rightarrow \alpha_3$

$A \rightarrow \alpha_n$

then the **grammar is LL(1)** if :

1.  $FIRST(\alpha_i) \cap FIRST(\alpha_j) = \emptyset$  for all  $i, j$

2. if one of  $\alpha_i$  is  $\lambda$ ,  $\alpha_n = \lambda$ , in addition to 1

$FIRST(\alpha_i) \cap follow(A) = \emptyset$ , for  $\forall i < n$

For example, Given the grammar :

$S' \rightarrow S\$$

$S \rightarrow aABC$

$A \rightarrow a \mid bbD$

$B \rightarrow a \mid \lambda$

$C \rightarrow b \mid \lambda$

$D \rightarrow \epsilon \mid \lambda$

let us see if it is LL(1)

$FIRST(a) \cap FIRST(bbD) = \emptyset$

$FIRST(a) \cap FOLLOW(B) = \emptyset$

$FIRST(b) \cap FOLLOW(C) = \emptyset$

$FIRST(\epsilon) \cap FOLLOW(D) = \emptyset$

Then this grammar is LL(1).

... another Grammar :

$\rightarrow SS$

$\rightarrow aAa \mid \lambda$

$\rightarrow abS \mid \lambda$

$FIRST(aAa) \cap FOLLOW(S) = \{a\} \cap \{\$, a\} = \{a\} \neq \emptyset$

... grammar is not LL(1).

### LL(1) Parsing Table Building Algorithm

Let us assume that we have a grammar that is LL(1). How do we build the LL(1) parsing table?

- For each production  $A \rightarrow \alpha$  in the grammar  $G$ ,  
 Add to the table entry  $T[A, a]$  the production  $A \rightarrow \alpha$ , where  $a \in FIRST(\alpha)$   
 If  $\lambda \in FIRST(\alpha)$ , Add to the table entry  $T[A, b]$  the production  $A \rightarrow \alpha$ ,  
 $\forall b \in FOLLOW(A)$ .
- All Remaining Entries are Error Entries.

For example, given the grammar :

- CGF -

$S \rightarrow SR \$_1$

$S \rightarrow +_2 \mid -_3 \mid \lambda_4$

$R \rightarrow dN.N_5 \mid .dN_6$

$N \rightarrow dN_7 \mid \lambda_8$

$$L(G) = [+|-] \left\{ \begin{array}{l} ddd \cdot ddd \\ ddd \cdot \\ \cdot ddd \end{array} \right\}$$

Note that the superscript denotes the production number.

$FIRST(SR \$) = \{+, -, d, .\}$

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

$FOLLOW(S) = \{d, .\}$

$FIRST(R) = \{d, .\}$

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

$FOLLOW(N) = \{., \$\}$

$$\begin{aligned} \Rightarrow FIRST(SR \$) &= FIRST(FIRST(S) \cdot FIRST(R) \cdot \{ \$ \}) \\ &= FIRST(\{+, -, \lambda\} \cdot \{d, .\} \cdot \{ \$ \}) \\ &= (\{+d, -d, d, +., -., \cdot\} \cdot \{ \$ \}) \\ &= \{+, -, d, ., \} \end{aligned}$$

VN	+	-	d	.	\$
IVT					
V	1	-1	1	1	
S	2	3	4	4	
R			5	6	
N			7	8	8

There should be no conflict (multiple entries) in the LL(1) table.  
 (G) of this grammar = all floating point numbers.

The parser works like this

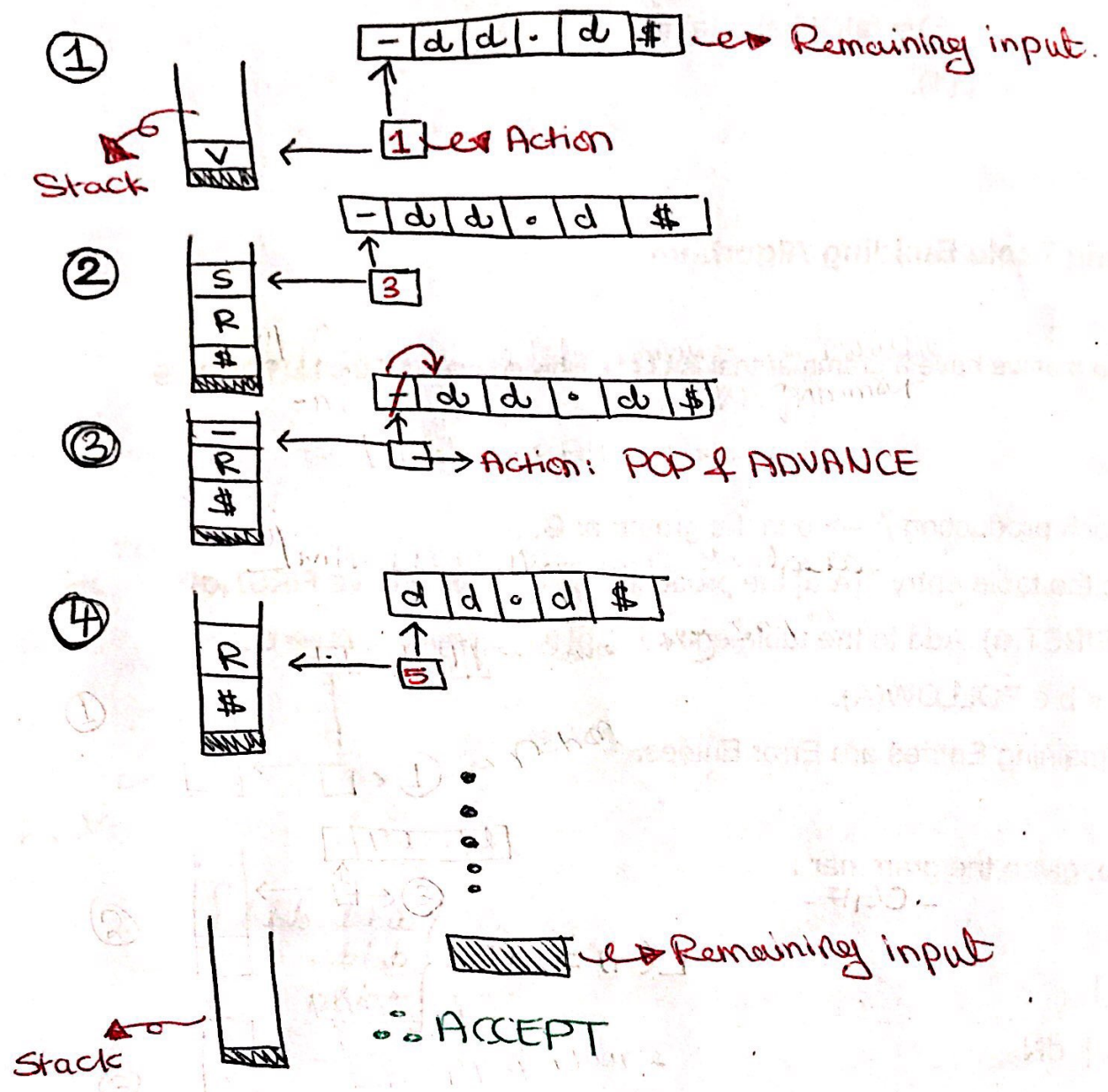
Stack	Remaining Input	Action
V	-dd.d\$	Production 1
SR\$	-dd.d\$	Production 3
-R\$	-dd.d\$	Pop & advance input
R\$	dd.d\$	Production 5
dN.N\$	dd.d\$	Pop & advance input
N.N\$	d.d\$	Production 7
dN.N\$	d.d\$	Pop & advance input
N.N\$	.d\$	Production 8
.N\$	.d\$	Pop & advance input
N\$	d\$	Production 7
dN\$	d\$	Pop & advance
N\$	\$	Production 8
\$	\$	Pop and Advance
$\lambda$	$\lambda$	Accept

If at any point the parser reaches a place where the input and the stack have 2 different terminal symbols, it throws a syntax error.



⚠️ If the parsing table contains no multiple entries (Conflicts) then the grammar is **LL(1) Grammar**.

⚠️ With LL(1) there are **NO** deadlocks.



Let us Take another example. Let the Grammar be :

program  $\rightarrow$  block \$<sub>1</sub>

block  $\rightarrow$  { decls<sub>5</sub> stmts }<sub>2</sub>

decls  $\rightarrow$  D ; decls<sub>3</sub> |  $\lambda$ <sub>4</sub>

stmts  $\rightarrow$  statement ; stmts<sub>5</sub> |  $\lambda$ <sub>6</sub>

statement  $\rightarrow$  if<sub>7</sub> | while<sub>8</sub> | ass<sub>9</sub> | scan<sub>10</sub> | print<sub>11</sub> | block<sub>12</sub> |  $\lambda$ <sub>13</sub>

VNVT = { \$, {, }, D, ,, if, while, ass, scan, print }


$\left. \begin{array}{l} \text{First(block)} = \{ \{ \} \} \\ \text{First(program)} = \{ \{ \} \} \end{array} \right\}$

VNVT	if	while	ass	scan	print	{	}	D	;	\$
Program						1				
block						2				
decls	4	4	4	4	4	4	4	3	4	
stmts	5	5	5	5	5	5	6		5	
statement	7	8	9	10	11	2			13	

Another example is the If..else statement with a delimiter. the grammar looks like this :

$S' \rightarrow S\$$   
 $S \rightarrow iCSE$   
 $E \rightarrow eS | \lambda$   
 $S \rightarrow a \rightarrow \text{ELSE}$   
 $C \rightarrow c$

VNVT	i	a	e	c	\$
S'	1	1			
S	2	5			
E			3,4		4
C				6	

 OUR IF/ELSE GRAMMAR ISN'T LL(1)

There is a conflict. To solve this, we can add a delimiter.

- $S' \rightarrow S\$$
- $S \rightarrow iCSEd$
- $E \rightarrow eS \mid \lambda$
- $S \rightarrow a$
- $C \rightarrow c$

V N   V T	i	a	e	c	d	\$
S'	1	1				
S	2	5				
E			3		4	
C				6		

The grammar is now unambiguous.

Alternatively, we can just remove out the production 4 in the conflict entry from the LL(1) table.

The New table is:

V N   V T	i	a	e	c	d	\$
S'	1	1				
S	2	5				
E			3		4	
C				6		

**Note:**

- if a grammar is LL(1), then it is unambiguous. However, the opposite is not necessarily true.
- Another thing to note is that in Top-Down parsing, we should avoid a grammar that is not LL(1).

## Bottom-Up Parsing

Recall that in Bottom-Up parsing, the parser starts from the given sentence, applying reductions until it reaches the starting symbol of the grammar or a deadlock.

The major problem with Bottom-Up parsing is which substring we should select in each reduction step.

The answer to the above question is :

In each reduction step, we select what is called **the handle**.

→ Solution for the Problem of bottom-up parsing.

The Handle is obtained by a **rightmost derivation in reverse**.

For example, Given the grammar :

$V \rightarrow S R \$$

$S \rightarrow +|-|\lambda$

$R \rightarrow .dN | dN.N$

$N \rightarrow dN | \lambda$

and the sentence

$-dd.d\$$

First, we derive the sentence rightmost.

$V \xrightarrow{rm} SR\$ \xrightarrow{rm} SdN.N\$ \xrightarrow{rm} SdN.dN\$ \xrightarrow{rm} SdN.dd\$ \xrightarrow{rm} SddN.d\$ \xrightarrow{rm} Sdd.d\$ \xrightarrow{rm} -dd.d\$$

So our handles would be :

$V \leftarrow SR\$ \leftarrow SdN.N\$ \leftarrow SdN.dN\$ \leftarrow SdN.d\lambda\$ \leftarrow SddN.d\$ \leftarrow Sdd\lambda.d\$ \leftarrow -dd.d\$$

But Compilers do not work like this. We already derived the sentence, why would we go back and do it again?

We could not build a Bottom-Up parser for every Context-Free Grammar. However, we are fortunate enough that there exist subsets of the Context-Free Grammar for which we can build a deterministic Bottom-Up parser i.e. the parser can determine/decide precisely where the handle is in each reduction step. Some of these subsets are:

LR Parsers :

- SLR(Simple-LR).
- LALR(Look-Ahead LR).
- LR.

Operator Precedence. "depends on matrix manipulation"

We will only be talking about the LR parsers, just to get an idea of how Bottom-Up parsing works.

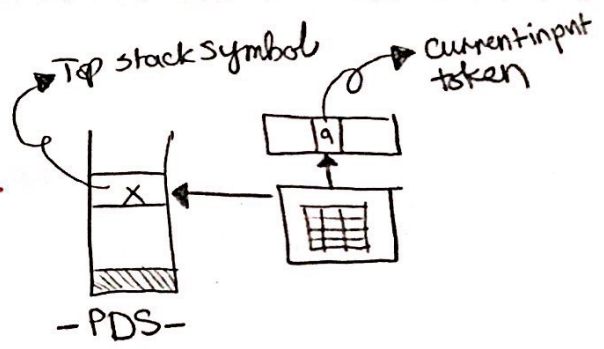
SLR Parsing

SLR parsing, and LR parsing in general, is a table driven parsing method.

All LL(1) grammars are a subset of SLR grammars.

All LR parsers contains :

1. A parsing table.
2. A stack. *e → Push Down Stack.*
3. The input string.



As a reminder, the LL(1) parser contains :

1. A parsing table.
2. A stack.
3. The input string.

**! Question:**  
How we build the LR parsing table?!

However, the way we build it is different.

## Building the SLR Parsing Table

Def: An LR(0) item of a grammar G is a production in G with a dot(.) at some position in the right side.

For example, the production

$A \rightarrow aBY$

This production generates the following LR(0) items :

$A \rightarrow \cdot aBY$

$A \rightarrow a \cdot BY$

$A \rightarrow aB \cdot Y$

$A \rightarrow aBY \cdot \rightarrow$  complete item

already had seen a string derived from aB

$A \rightarrow aB \cdot Y$

expect to see on input a string derived from Y.

Note that for  $A \rightarrow \lambda$ , this generates only  $A \rightarrow \lambda$ .

Generally speaking, if the right side of the production is of length n, then there are  $n+1$  LR(0) items.

For  $A \rightarrow \alpha, \alpha \neq \lambda$

$|\alpha| = n,$

LR(0) generates  $n+1$  items.

The LR(0) item

$A \rightarrow aB \cdot Y$

Means that the parser has scanned on the input a string derived from aB and expects to see a string derived from Y.

We need to define the following 2 functions.

### The CLOSURE function

function CLOSURE(I) // I is a set of LR(0) items

Repeat

For (every LR(0) item  $A \rightarrow \alpha \cdot B\beta$  in I, and for every production  $B \rightarrow \delta$  in G,

Add the LR(0) item  $B \rightarrow \cdot \delta$  to I)

Until no more items to be added;

Let us apply this to our grammar :

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

! This grammar isn't LL(1) !  
because there will be conflict.

This grammar is not LL(1) because  
 $FIRST(T*F) \cap FIRST(F) = \{ (, a \} \neq \emptyset$

We will need to build the LR(0) sets of items.

we start with :  $I_0: E' \rightarrow .E$

CLOSURE( $I_0$ )

- $I_0: E' \rightarrow .E$
- $E \rightarrow .E+T$
- $E \rightarrow .T$
- $T \rightarrow .T*F$
- $T \rightarrow .F$
- $F \rightarrow .(E)$
- $F \rightarrow .a$

The GOTO function

function GOTO( $I, X$ ) =

CLOSURE(all items  $A \rightarrow \alpha X \beta$  Where  $A \rightarrow \alpha \underline{X} \beta$  in  $I$ )

Let us apply this to the grammar above. groups:

- $I_1: E \rightarrow .E, E \rightarrow .E+T$
- $I_2: E \rightarrow .T, T \rightarrow .T*F$
- $I_3: T \rightarrow .F$
- $I_4: F \rightarrow .(E)$

	a	(	+	x	)	#
E'	0	0				
E	1, 2	1, 2				
T	3, 4	3, 4				
F	6	5				

⚠ It's Not LL(1) because there are conflicts



15:  $F \rightarrow \cdot a$

and take the CLOSURE for all these sets. The resultant is :

- 16:  $E \rightarrow \cdot E$
  - 17:  $E \rightarrow \cdot E + T$
  - 18:  $E \rightarrow \cdot T$
  - 19:  $T \rightarrow \cdot T * F$
  - 20:  $T \rightarrow \cdot F$
  - 21:  $F \rightarrow \cdot (E)$
  - 22:  $F \rightarrow \cdot a$
- since they are the same they go to a new state.*

16:  $E \rightarrow \cdot E$   
 $E \rightarrow \cdot E + T$  } complete

17:  $E \rightarrow \cdot T$   
 $T \rightarrow \cdot T * F$  } complete

18:  $T \rightarrow \cdot F$   
Complete

19:  $F \rightarrow \cdot (E)$   
 $E \rightarrow \cdot E + T$

20:  $E \rightarrow \cdot T$   
 $E \rightarrow \cdot T * F$

21:  $T \rightarrow \cdot F$   
 $F \rightarrow \cdot (E)$   
 $F \rightarrow \cdot a$

22:  $F \rightarrow \cdot a$

23:  $E \rightarrow \cdot E + T$   
 $E \rightarrow \cdot T * F$   
 $T \rightarrow \cdot F$   
 $F \rightarrow \cdot (E)$   
 $F \rightarrow \cdot a$

24:  $E \rightarrow \cdot T * F$   
 $F \rightarrow \cdot (E)$   
 $F \rightarrow \cdot a$

25:  $F \rightarrow \cdot (E)$   
 $E \rightarrow \cdot E + T$

26:  $E \rightarrow \cdot E + T$   
 $T \rightarrow \cdot T * F$

27:  $T \rightarrow \cdot T * F$

28:  $F \rightarrow \cdot (E)$

CRC(0) sets of items  
 $\{I_0, I_1, I_2, \dots, I_{11}\}$

- $I_0: E' \rightarrow \cdot E \quad I_1$
- ①  $E \rightarrow \cdot E + T \quad I_1$
- ②  $E \rightarrow \cdot T \quad I_2$
- ③  $T \rightarrow \cdot T * F \quad I_2$
- ④  $T \rightarrow \cdot F \quad I_3$
- ⑤  $F \rightarrow \cdot (E) \quad I_4$
- ⑥  $F \rightarrow \cdot a \quad I_5$

$I_1: E' \rightarrow E \cdot \text{ Complete}$   
 $E \rightarrow E \cdot + T \quad I_6$

$I_2: E \rightarrow T \cdot \text{ Complete}$   
 $T \rightarrow T \cdot * F \quad I_7$

$I_3: T \rightarrow F \cdot \text{ Complete}$

- $I_4: F \rightarrow (\cdot E) \quad I_8$
- $E \rightarrow \cdot E + T \quad I_8$
- $E \rightarrow \cdot T \quad I_2$
- $T \rightarrow \cdot T * F \quad I_2$
- $T \rightarrow \cdot F \quad I_3$
- $F \rightarrow \cdot (E) \quad I_4$
- $F \rightarrow \cdot (a) \quad I_5$

$I_5: F \rightarrow a \cdot \text{ Complete}$

- $I_6: E \rightarrow E + \cdot T \quad I_9$
- $T \rightarrow \cdot T * F \quad I_9$
- $T \rightarrow \cdot F \quad I_3$
- $F \rightarrow \cdot (E) \quad I_4$
- $F \rightarrow \cdot a \quad I_5$

- $I_7: T \rightarrow T * \cdot F \quad I_{10}$
- $F \rightarrow \cdot (E) \quad I_4$
- $F \rightarrow \cdot a \quad I_5$

- $I_8: F \rightarrow (E \cdot) \quad I_{11}$
- $F \rightarrow E \cdot + T \quad I_6$

$I_9: E \rightarrow E + T \cdot \text{ Complete}$   
 $T \rightarrow T \cdot * F \quad I_7$

$I_{10}: T \rightarrow T * F \cdot \text{ Complete}$

$I_{11}: F \rightarrow (E) \cdot \text{ Complete}$

LR(0) sets of items  
 $\{I_0, I_1, I_2, \dots, I_{11}\}$

# Constructing the SLR table

Input : LR(0) sets of items

Output : SLR(1) parsing table

- For every item  $A \rightarrow \alpha \cdot aB$  in  $I_i$ ,  $a \in V_T$ , and  $GOTO(I_i, a) = I_j$ , then set:  $ACTION[i, a] = S_j$  (shift and push  $j$  on top stack).
- For item  $A \rightarrow \alpha \cdot$  (complete item) in  $I_i$ ,  $ACTION[i, b] = \text{Reduce by } A \rightarrow \alpha$  FOR ALL  $b \in FOLLOW(A)$ .
- For  $S' \rightarrow S$  in  $I_i$ ,  $ACTION[i, \$] = \text{Accept}$ .
- If  $GOTO(I_i, A) = I_j$  then set,  $GOTO(i, A) = j$ .
- All remaining entries are error entries.

Terminals

Let us apply this to the example above and generate the table

Follow (E) = {+, }, \$  
 Follow (T) = {+, \*, }, \$  
 Follow (F) = {+, \*, }, \$

State/ token	ACTION [ $V_T$ ]					GOTO [ $V_N$ ]			
	a	+	*	(	)	\$	E	T	F
0	S5						1	2	3
1		S6				A accept			
2		R2	S7		R2	R2			
3		R4	R4		R4	R4			
4	S5			S4			8	2	3
5		R6	R6		R6	R6			
6	S5			S4				9	3
7	S5			S4					10
8		S6			S11				
9		R1		S7	R1	R1			
10		R3	R3		R3	R3			
11		R5	R5		R5	R5			

Action section

→ a → shift & go to

Go-to section

→ A → go to

→ product in number

No conflict --> SLR(1) grammar

# Parsing The SLR Table

- E+T
- T
- T\*F
- F
- (E)
- a

Let us examine the sentence a + a \$

Stack	Remaining	Action
0	a + a \$	S5
0 a 5	+ a \$	R6
0 F ③	+ a \$	R4
0 T 2	+ a \$	R2
0 E 1	+ a \$	S6
0 E 1 + 6	a \$	S5
0 E 1 + 6 a 5	\$	R6
0 E 1 + 6 F 3	\$	R4
0 E 1 + 6 T 9	\$	R1
0 E 1	\$	Accept

R: F → a

R: T → F

H: Handles

# LR Parsing Techniques

## → Definition

The main difference between LR and SLR is the CLOSURE function,

function  $CLOSURE(I) // I$  is a set of LR(1) items

Repeat

for every LR(1) item  $[A \rightarrow \alpha \cdot B \beta, a]$  in  $I$ , and for every production  $B \rightarrow \delta$  in  $G$ ,

Add the LR(1) item  $[B \rightarrow \cdot \delta, b]$  where  $b$  belongs to  $FIRST(\beta a)$  to  $I$

Until no more items to be added;

Where An LR(1) item is an LR(0) item with a Look-ahead Symbol.

For example  $[A \rightarrow \alpha \cdot \beta, a]$  where  $a$  is the look-ahead. The look-ahead symbol " $a$ " has no effect whatsoever on an item  $[A \rightarrow \alpha \cdot \beta, a] \beta \neq \lambda$ .

(not complete item) However, if the item is a complete  $[A \rightarrow \alpha \cdot, a]$ , this means we reduce by the production  $A \rightarrow \alpha$  on token " $a$ ".

For example, Give the grammar :

$$S' \rightarrow S$$

$$(1) S \rightarrow CC$$

$$(2) C \rightarrow cC$$

$$(3) C \rightarrow d$$

10:

$S' \rightarrow \cdot S, \$$	$l_1$
$S \rightarrow \cdot CC, \$$	$l_2$
$C \rightarrow \cdot cC, c, d$	$l_3$
$C \rightarrow \cdot d, c, d$	$l_4$

$\xrightarrow{e} First(C\$)$

$$\begin{array}{l}
 A \rightarrow \alpha \cdot B \beta \\
 B \rightarrow \cdot \delta \\
 \beta = \lambda
 \end{array}
 \xrightarrow{a}
 \begin{array}{l}
 S' \rightarrow \cdot S \$ \\
 S \rightarrow \cdot CC
 \end{array}$$

$$First(B\$) = First(\lambda \$) = \{ \$ \}$$

11:

$$S' \rightarrow S \cdot \$ \quad \text{Accept}$$

12:

$S \rightarrow C \cdot C, \beta, \$$	$l_5$
$C \rightarrow \cdot cC, \$$	$l_6$
$C \rightarrow \cdot d, \$$	$l_7$

13:

$C \rightarrow c \cdot C$	$c, d$	$l_8$
$C \rightarrow \cdot cC$	$c, d$	$l_3$
$C \rightarrow \cdot d$	$c, d$	$l_4$

→ Definition:

An LR(1) item is an LR(0) item with look-ahead symbol (token). That's,  $A \rightarrow \alpha \cdot \beta$  is an LR(0) item, then

LR(1) item is:  $[A \rightarrow \alpha \cdot \beta, a], a \in V_T$

\* We start building the LR(1) sets of items,

We start with LR(1) item  $I_0: [S' \rightarrow \cdot S, \$]$

↪ look-ahead

$[A \rightarrow \alpha \cdot B \beta, a]$

\* Add

$[B \rightarrow \cdot \gamma, b] \quad b \in \text{First}(\beta a)$



- \* Look-ahead has no effect if the item isn't complete.
- \* Look-ahead affects the parser if the item is complete.

14:  $C \rightarrow d.$  c,d Complete

15:  $S \rightarrow CC.$  \$ Complete

16:  $C \rightarrow c.C$  \$  $I_9$   
 $C \rightarrow .cC$  \$  $I_6$   
 $C \rightarrow .d$  \$  $I_7$

17:  $C \rightarrow d.$  \$ Complete

18:  $C \rightarrow cC.$  c,d Complete

19:  $C \rightarrow cC.$  \$ Complete

V	ACTION $V_T$			GOTO $V_N$	
	c	d	\$	S	C
0	S3	S4		1	2
1			A		
2	S6	S7			5
3	S3	S4			8
4	R3	R3			
5			R1		
6	S6	S7			9
7			R3		
8	R2	R2			
9			R2		

No conflict, the grammar is an LR grammar. (1)

Multiple entries

If we look at the above example, we can see that some sets of items have the same core items (LR(0) items), but the look-ahead is different. For example (I7, I4), (I3, I6), (I8, I9). Let us say we merge the states.

V	ACTION			GOTO	
	c	d	\$	S	C
0	S3	S4		1	2
1			A		
2	S3	S4			5
3	S3	S4			8
4	R3	R3	R3		
5			R1		
8	R2	R2	R2		

⚠ Some of the LR(1) sets of items, the core items are identical.

- I3, I6
- I4, I7
- I8, I9

→ We merge these sets of items.

This is now a simplified table. if the parsing table after merging has no conflicts (like in the above example), then the grammar is an LALR(1) Grammar.



# Conflict Types:

given S(shift) - R(Reduce)

→ Possibilities of Conflicts :-

(1) SS (Shift-Shift)

In one of the LR(0) sets of items,

$I_1$ :  
 $\vdots$   
 $A \rightarrow \alpha \cdot a \beta$   $I_6$   
 $B \rightarrow \delta \cdot a \delta$   $I_5$

⚠ Impossible to happen

Since 'a' both should go to the same State.

∴ S.S Conflict Impossible.

(2) SR (Shift-Reduce)

$I_i$ :  
 $\vdots$   
 $A \rightarrow \alpha \cdot a \beta$   $I_j$   
 $B \rightarrow \delta \cdot a$

	a
$\vdots$	
$i$	$S_j$ $R_B \rightarrow \delta$

(3) RR (Reduce-Reduce)

$I_i$ :  
 $\vdots$   
 $A \rightarrow \alpha \cdot$  a C  
 $B \rightarrow \delta \cdot$  a C

	a
$\vdots$	
$i$	$R_A \rightarrow \alpha$ $R_B \rightarrow \delta$



(1)

⚠ In LALR parsing table there's only Reduce-Reduce Conflict.

- Assume the grammar is LR(1) & there's a Shift-Reduce Conflict in the LALR parsing table.

Since there's a S-R conflict in the LALR(1), this means, in one of the LALR(1) sets of items there are:

$I_i :$   
 $\vdots$   
 $A \rightarrow \alpha \cdot a \beta \quad c \quad I_j$   
 $B \rightarrow \delta \cdot \quad a \quad C$

		a	
	$\vdots$		
i		Sj R <sub>B→δ</sub>	

But the items  $A \rightarrow \alpha \cdot a \beta \quad c$   
 $B \rightarrow \delta \cdot \quad a$

Came from the LR(1) parsing table so it's not LR(1)

∴ this contradicts our assumption.