

# Lexical Analysis(Scanner)

## What is a Lexical Analyzer?

A Lexical Analyzer or Scanner is an algorithm that groups the characters of the source code to form the **Tokens**, moreover, it returns the **Internal Representation Number** of these tokens, which is a kind of an **ID** that assigned for each token.

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.
2. **Values** : such as integers(1, 2, 3, 4) or floating point(1.1, 2.34, 5234.123) et
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.

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 :



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

Leads us to these ID's :

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

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 the table which contains the user defined name(mostly variables), its type, data type, and its values. The Symbol table for this segment of code would be:

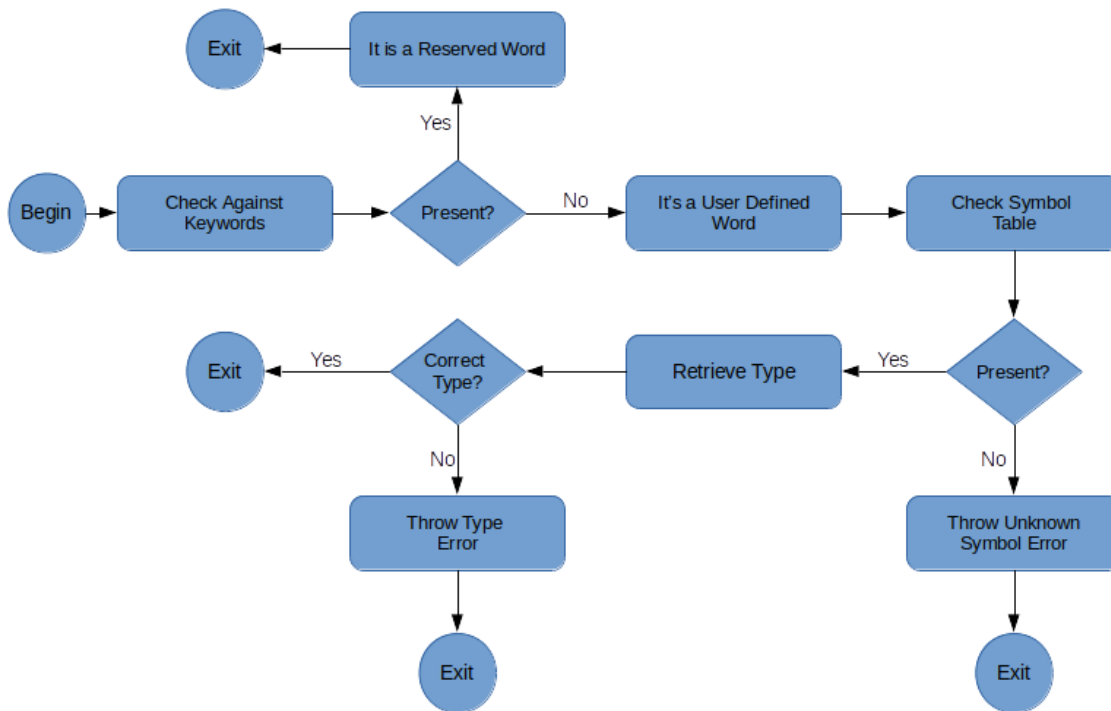
```

.
.
.
int compute(int,int);
int n;
float x=4.5 , y;
y=2*x;
const int m=10;
n=compute(10,m);
.

```

Name	Type	Data type	Value
compute	function-name	integer	0
n	variable	integer	0-20
x	variable	float	4.5
y	variable	float	0-9.0
m	constant	integer	10

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 name, 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( mismatch - type error ).



## Regular Expressions(Regular Languages)

Regular Expressions (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  $V$  is a Set of Symbols.

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

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

Are all strings defined on  $V$ .

Formally, we define:

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

**Def:** given strings  $x, y \in S$ , then define the

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

$XY$  = The string formed by following  $x$  with  $y$

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

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.

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:

Given  $V=\{a, b, , \dots ,z\}$

$L_1=\{aa, gb, c\}$

$L_2=\{asdasd, qwe, asd\}$

$L_3=\{abb\}$

This definition also leads us to the conclusion :

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

## Set of Operations On Languages

Given an Alphabet  $V$ , assume that :

1.  $L = \{\text{set of all languages defined on } V\}$

2.  $L = \{L_1, L_2, L_3, \dots, L_n, \dots, \infty\}$

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

### Concatenation Operation

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

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

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

For Example let  $L=\{ab, bba\}$  and  $M=\{aa,bb\}$ , then :

$LM=\{abaa, abbb, bbaaa, bbabb \}$

$ML=\{aaab, aabba, bbab, bbbba\}$

Note that :

1.  $LM \neq ML$ . (Concatenation on languages is not commutative).

2.  $L\{\lambda\} = \{\lambda\}L = L$ . ( $\lambda$  is the identity for concatenation).

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

### **The OR "|" Operation**

OR is  $\cup$  operation in set theory

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

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

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)

### **The Closure "\*" Operation (A Unary 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  $L_R$  and  $L_S$  respectively, then :
  - a)  $RS$  is a regular language denoting  $L_R L_S$
  - b)  $R|S$  is a regular language denoting  $L_R | L_S$
  - c)  $R^*$  is a regular language denoting  $L_R^*$

### 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  $\{100,000,1100,0000,\dots\}$
- 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  $\{bbb,abbb,bbba,bbbbbb,\dots\}$

### **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

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

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

## 2. Floating Point Numbers :

$$[+|-] d^+.d^+$$

$$[+|-] d.d^+E(+|-)d^+$$

### **Special Symbols**

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

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

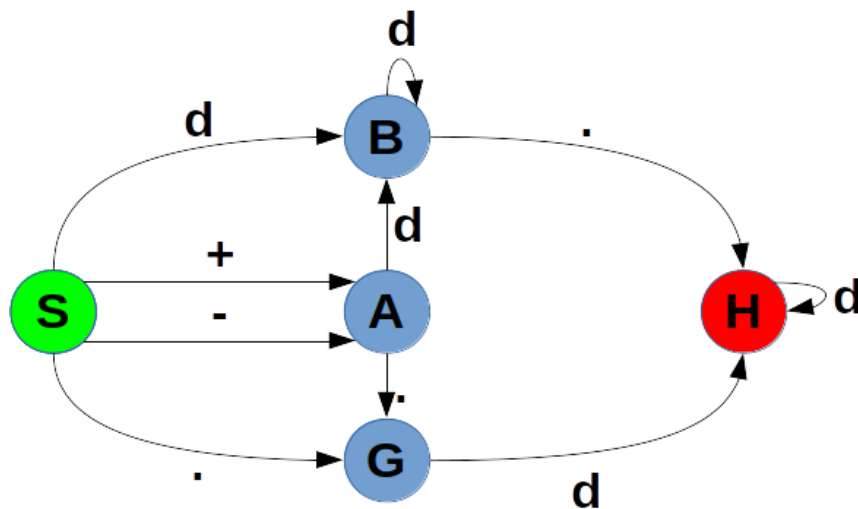
### **Finite State Automata(FSA)**

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.
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 :

```

-   d   d   d   .   d   d
S--->A--->B--->B--->B--->H--->H--->H

```

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 whole string we end up with a final state.

The Regular Language(expression) generated(accepted)

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

$$L(M)=[+|-]\{d^+., .d^+, d^+.d^+\}$$

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) None-Deterministic Finite State Automata (NDFSA)

An algorithm is none-deterministic or fuzzy if there are options(choices) in the algorithm. An example of a none-deterministic algorithm is the solution of the [Knight Tour Problem](#), which is based on [Backtracking Techniques](#).

A Finite State Automata is none-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 are a choices (trial and error) to implement. The only way to solve none-deterministic machines is to use backtracking. This is not practical in a compiler, because backtracking is a very time consuming algorithm 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(DFSA)

If A Machine is **not** none-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.

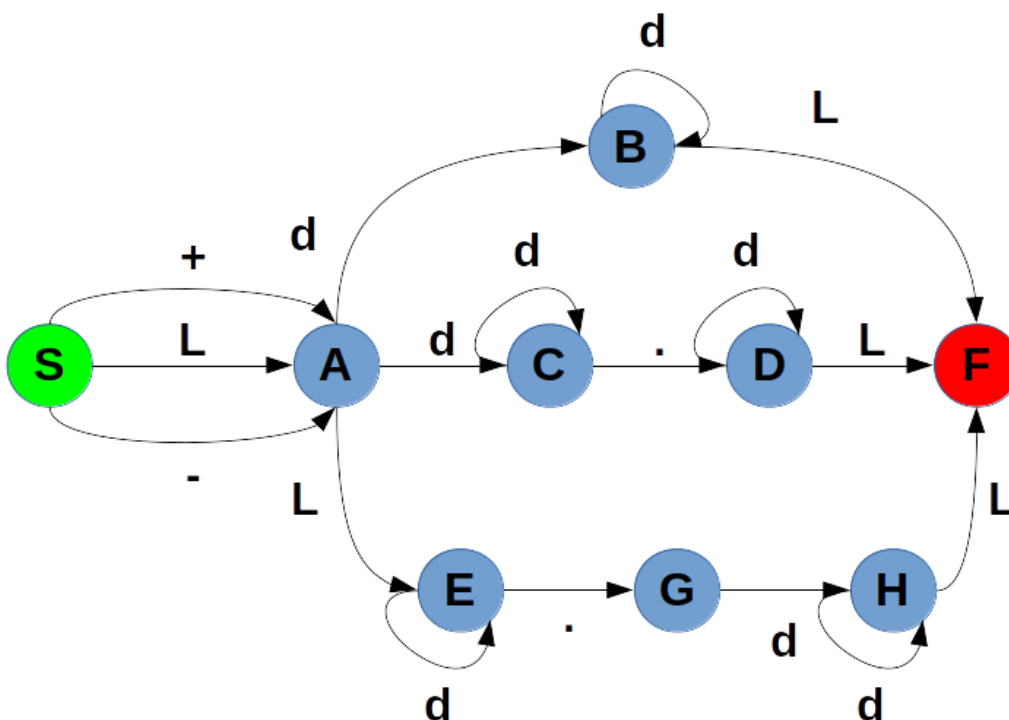
Only Deterministic Finite State Automata are used for compilers.

### 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 none-determinism.
3. Removal of inaccessible states.
4. Merging equivalent states.

Given the following NDFSA :



What is  $L(G) = ?$

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

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:

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
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.

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</b>				B
C			D	C
<b>D</b>				D
E			G	E
G				H
<b>H</b>				H
<b>F</b>				

## (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
B				B
C			D	C
D				D
E			G	E
G				H
H				H
F				
B,C,E			D,G	B,C,E
D,G				D,H
D,H				D,H

## (3) Removal of Non-Accessible States

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.

This results in this table :

State \ VT	+	-	.	d
$\mathcal{S}$	A	A	G	B,C,E
$\checkmark$ A			G	B,C,E
<b>B</b>				B
<b>C</b>			D	C
<b>D</b>				D
<b>E</b>			G	E
$\mathcal{G}$				H
$\mathcal{H}$				H
<b>F</b>				
$\mathcal{B,C,E}$			D, G	B,C,E
$\mathcal{D,G}$				D,H
$\mathcal{D,H}$				D,H

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

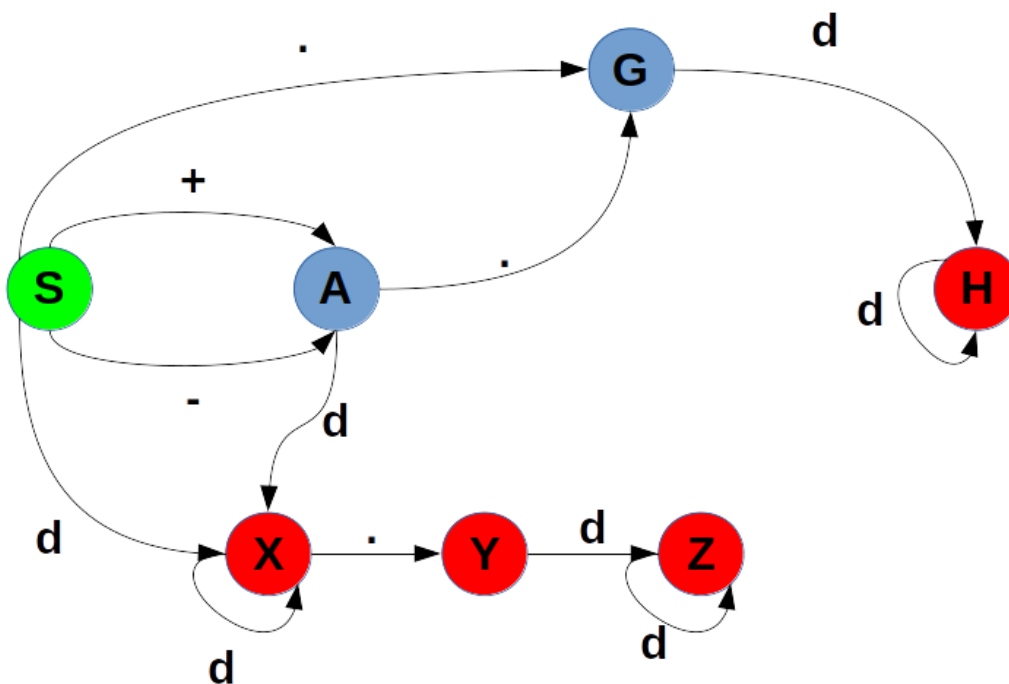
State \ VT	+	-	.	d
$\mathcal{S}$	A	A	G	B,C,E
$\checkmark$ A			G	B,C,E
$\mathcal{G}$				H
$\mathcal{H}$				H
$\mathcal{B,C,E}$			D, G	B,C,E
$\mathcal{D,G}$				D,H
$\mathcal{D,H}$				D,H

This is now a Deterministic Machine That accepts the same languages as the



State \ + - . d				
V T				
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

**DEF:** State  $p$  and state  $q$  in a FSA are said to be equivalent if:

For any string  $X$ , machine  $M$  in state  $p$  accepts  $X$  iff machine  $M$  in state  $q$  accepts  $X$ .

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)  $a \in V_T$ , either both  $\{p,q\}$  have transitions on "a", or both  $\{p,q\}$  don't have transitions on "a".

3.  $p \neq q$  ;

**Note:**  $(p,q) \equiv (q,p)$ .

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

State \ $V_T$	a	b	c
1	2	5	
2	3	4	1
3	5	2	
<u>4</u>	6		1
5	1	4	1
<u>6</u>	4		1
7	3	5	3

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

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

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

this results in this feasible-pairs table :

feasible pairs \ $V_T$	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

**DEF:** A feasible pairs  $(p,q)$  is marked if there is a transition to a pair  $(r,s)$  such that :

1.  $r \neq s$ . **AND**
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.  $1 \equiv 3 \rightarrow 1$
2.  $2 \equiv 5 \rightarrow 2$
3.  $4 \equiv 6 \rightarrow 4$

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

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
✓ H				H
✓ X			Y	X
✓ Y				Z
✓ Z				Z

Lets quickly apply what we learned on this table.

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

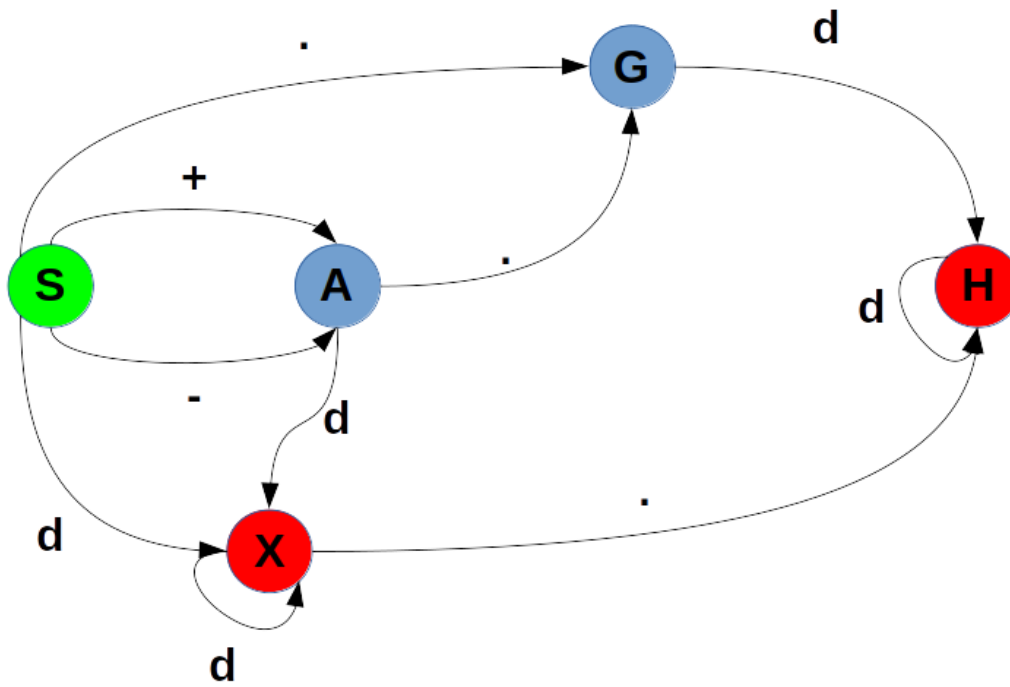
$H \equiv Y \equiv Z \rightarrow H$

Resulting in this state table :

State VT	+	-	.	d
<del>S</del>	A	A	G	X
✓ A			G	X
<del>G</del>				H
<del>H</del>				H
<del>X</del>			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.



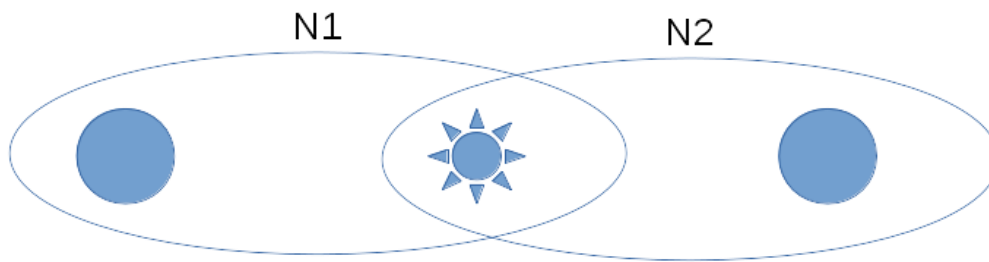
This machine accepts the language.

$L(G) = [+|-] \{ dddd, ddd.dd, dddd., .dddd \}$

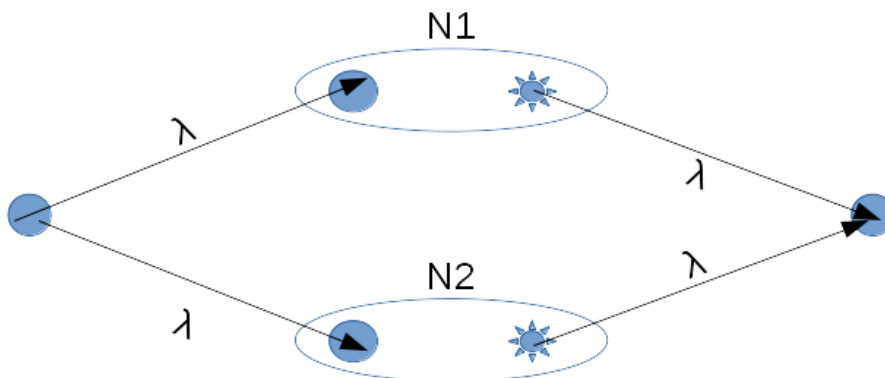
Which is the same language of our original machine.

## Creating a NDFSA From a Regular Expression

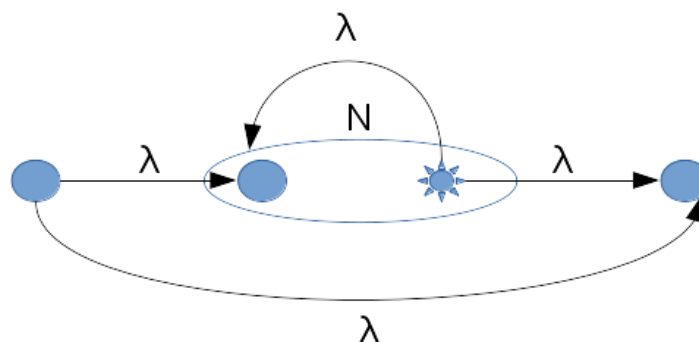
1. Decompose the regular expression to its primitive components :
  - a. for  $\lambda$ ,  $X \xrightarrow{\lambda} Y$ .
  - b. for  $a$ ,  $X \xrightarrow{a} Y$ .
2. Supposed that  $N_1$ ,  $N_2$  are transition diagrams for the regular expressions  $R_1$ ,  $R_2$  respectively,  $N_1$  accepts  $R_1$  &  $N_2$  accepts  $R_2$ , then :



- a. which represents  $R_1 R_2$  is :
- b. which represents  $R_1 | R_2$  is :

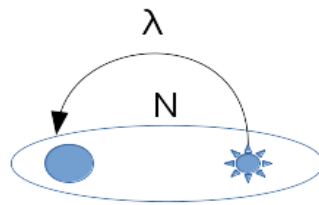


- c. which represents  $R^*$  is :



d.  $x \quad N$   
+

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

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

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			

State\ $V_T$	L	d	$\lambda$
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\ $V_T$	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 \ $V_T$	L	d
1	2	
2	5	7
5	5	7
7	5	7



feasible pairs \ $V_T$	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\ $V_T$	L	d
1	2	-
2	2	2