# Lexical Analysis

Rakib Mahmud

# Grammar

G is a grammar, which consists of a set of production rules. It is used to generate the strings of a language.

G= { V, T, P, S }

V= Variable (Represented with capital letter)

T= Terminal (Represented with small letter)

P= Production Rule

S= Start symbol

$$S \to aSb \,/\, \varepsilon$$

$\varepsilon$, aSb, aaSbb, aaaSbbb

$\varepsilon$, ab, aabb, aaabbb

$\varepsilon$, ab, $a^2b^2$, .......$a^nb^n$    n>= 0

Rakib Mahmud

# DFA

- In DFA, for each input symbol, one can determine the state to which the machine will move. Hence, it is called Deterministic Automaton.

- As it has a finite number of states, the machine is called Deterministic Finite Machine or Deterministic Finite Automaton.
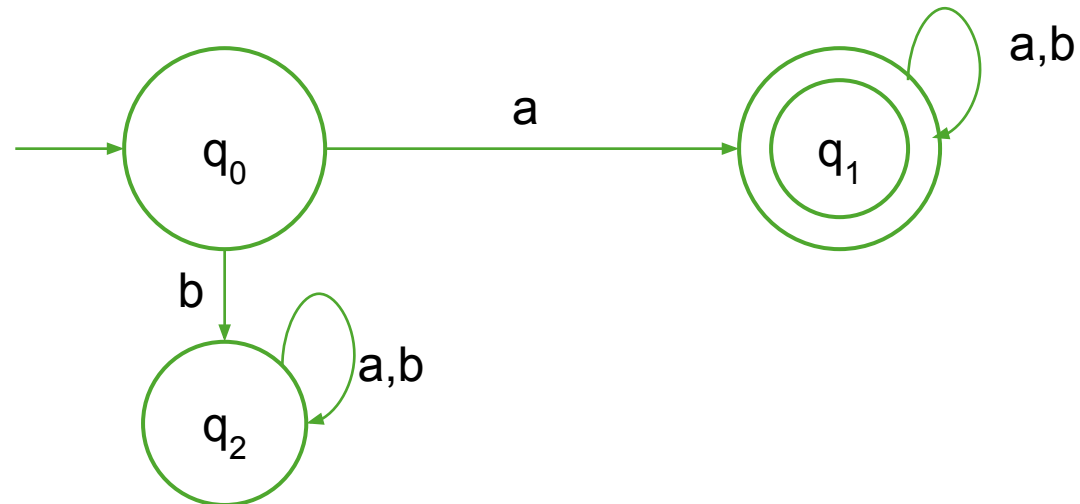
# Formal Definition of a DFA

- A DFA can be represented by a 5-tuple $(Q, \sum, \delta, q_0, F)$ where −
  - Q is a finite set of states.
  - $\sum$ is a finite set of symbols called the alphabet.
  - $\delta$ is the transition function where $\delta: Q \times \sum \to Q$
  - $q_0$ is the initial state from where any input is processed ($q_0 \in Q$).
  - F is a set of final state/states of Q ($F \subseteq Q$).

Rakib Mahmud

# Graphical Representation of a DFA

- A DFA is represented by digraphs called state diagram.
- The vertices represent the states.
- The arcs labeled with an input alphabet show the transitions.
- The initial state is denoted by an empty single incoming arc.
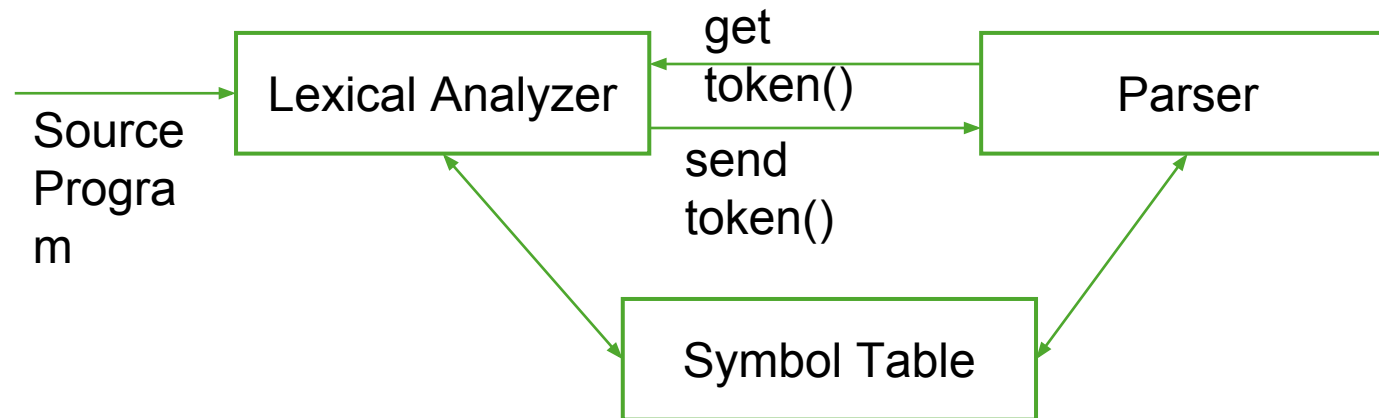- The final state is indicated by double circles.

# Example

- **Make a DFA where all strings start with a.**
- You have two alphabets (a,b)
  - { a, aa, aaa, ab, ba }
    ⌘

Rakib Mahmud

# Lexical Analyzer

- Lexical analyzer divides the given program into meaningful words which is known as tokens.

- Tokens are normally identifiers, separator, keywords, operators, constant and special symbol.

- The program which did this is called lexer, tokenizer, scanner.

- Lexical analyzer eliminates comments, white space character. (blank, Tab)

- It helps in giving error message such as exceeding length by providing row no and column no.

# Lexical Analyzer

# Lexical Analyzer

- Lexical analyzer uses DFA to do tokenization.
- While doing tokenization lexical analyzer always gives importance to longest matching.

# Counting no of tokens (Ex: 1)

```
1
2 ▾ int main() { //5
3        // Write C code here. //0
4        int a=20,b=30; //9
5        if(a<b). // 6
6            {return b;} //5
7        else //1
8            {return a;} //5
9
10 } //1
```

Total tokens = 32

# Counting no of tokens (Ex: 2)

```
1
2 ▾ int main() { //5
3       // Write C code here. //0
4       int i=10; //5
5       printf("i=%d, &i=%x",i,&i); //10
6
7 } //1
```

Total tokens = 21

# Counting no of tokens (Ex: 3)

```
1
2   int main() { //5
3       // Write C code here. //0
4       int i; //3
5       for(i=0;i<5;i++)//13
6       { //1
7           int i=102; //5
8           printf("%d the value of i: ",i); //7
9           i++;//3
10      }//1
11      return 0;//3
12
13  } //1
```

Total tokens = 42

# Counting no of tokens (Ex: 4)

```
1
2  main )( } //4
3      // Write C code here. //0
4      x=a+b*c; //8
5      int x,a,b,c;//9
6      y=x+a;//6
7
8  } //1
```

Total tokens = 28

# Counting no of tokens (Ex: 5)

```
1  main() //3
2  { //1
3      int x==10,y<=20; //9
4      printf("%d %d %d",x); //7
5  } //1
```

Total tokens = 21

# Counting no of tokens (Ex: 6)

```
1   main() //3
2 ▾ { //1
3        char *c="string"; //6
4        float b=100.74; //5
5        char d='e'; //5
6        int f=200; //5
7        in t f=200; //6
8        in /* comment */t f=200; //6
9        ch ar d='e'; //6
10       ch/*comment*/ar d='e';//6
11  } //1
```

Total tokens = 50

# Counting no of tokens (Ex: 7)



```
main.c

1  main() //3
2  { //1
3      a=b+++----+++==; //11
4      x===y /*abcd***/*abcd*/; //9
5      int **p; //5
6      printf("%d %d",a,b); //9
7  } //1
```

Total tokens = 39

# Errors Detected in Lexical analysis

- Numeric literals that are too long. (int a=12345678910111213141516171 81920)

- Long identifiers

- Ill-formed numeric literals. (int a= $123)

- Input characters that are not in the source language.

# Error Recovery

- **Delete:** Unknown characters are deleted. Also known as panic mode recovery. Example:
  - 'Charr' corrected as 'char' deleting r.

- **Insert:** An extraa or missing charactered is inserted to form a meaningful token.Example:
  - 'cha' corrected as 'char' inserting 'r'

- **Transpose:** Based on certain rules we can transpose two characters. Example:
  - 'Whiel' can be corrected to 'while'

# Error Recovery

- **Replace**: Based on replacing one character by another. Example:
  - 'chrr' can be corrected as 'char' by replacing 'r' with 'a'.