

CS & IT ENGINEERING

compiler Design

Lexical Analysis & Syntax Analysis

Lecture No. **2**



By- DEVA Sir



01 Lexical Analysis

02 Symbol Table

03 Types of Errors

04

05

Group-1 (phase)

1. Lexical Analysis → A
2. Syntax Analysis → B
3. Semantic Analysis → C
4. Code Generator → D

Group-2 (Input to phase)



- A. character stream
- B. Token stream
- C. Parse Tree
- D. Intermediate Code
- E. Semantic Tree
- F. Assembly Code

Group - 1 (phase)

1. Lexical Analysis → B
2. Syntax Analysis → C
3. Semantic Analysis → E
4. Code Generator → F

Group - 2 (output of phase)



- A. character stream
- B. Token stream
- C. Parse Tree
- D. Intermediate Code
- E. Semantic Tree
- F. Assembly Code

Group-1 [Compiler]

1. Lexical phase → D

2. Syntax phase → C

3. Semantic phase → B

4. Compiler → A

Group-2 [TOC]

A. Halting TM

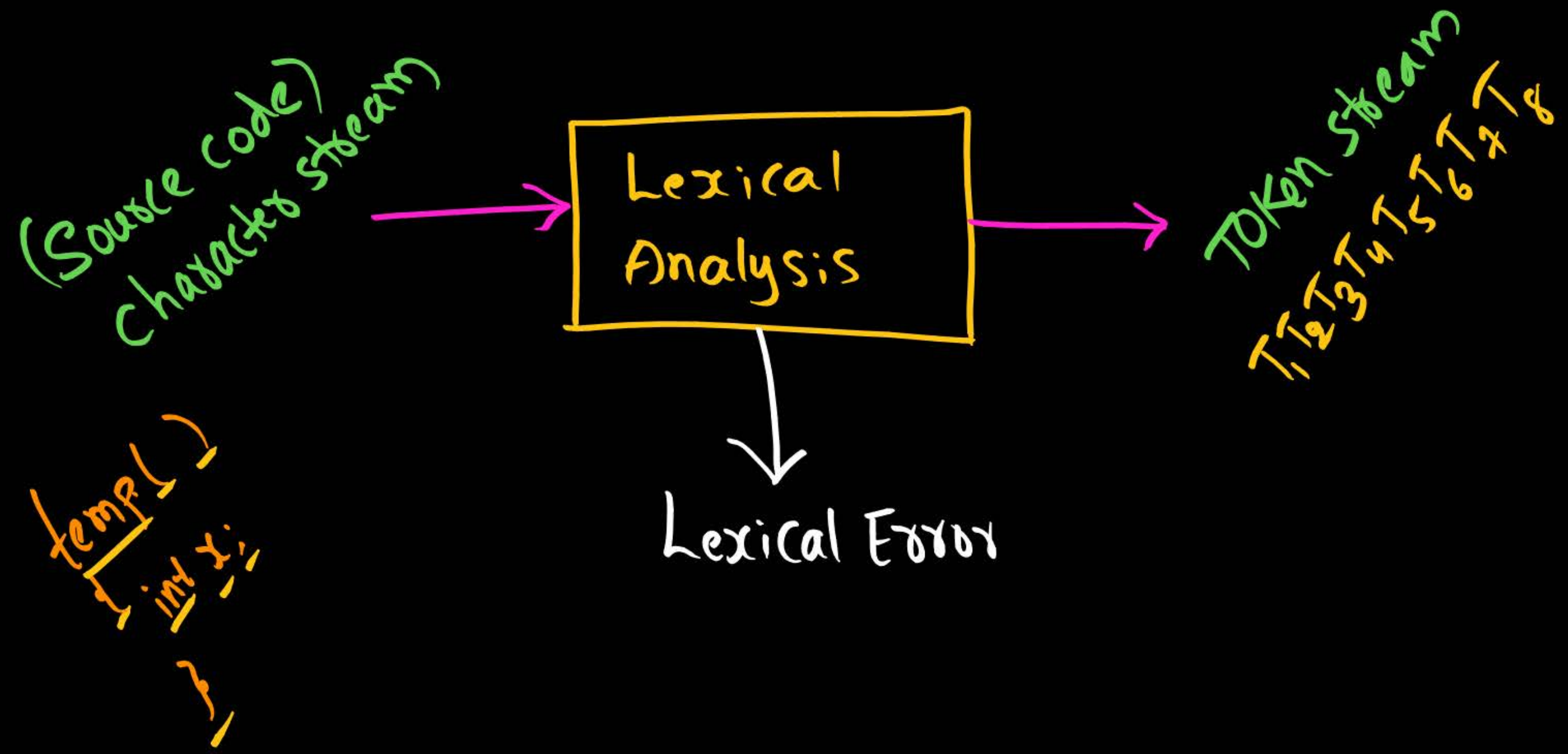
B. LBA

C. PDA

D. FA



Lexical Analysis :



Lexical Analysis



- Token Recognizer
- Token Generator/producer
- Scanner
- Linear phase
- Lexical phase
- First phase of compiler

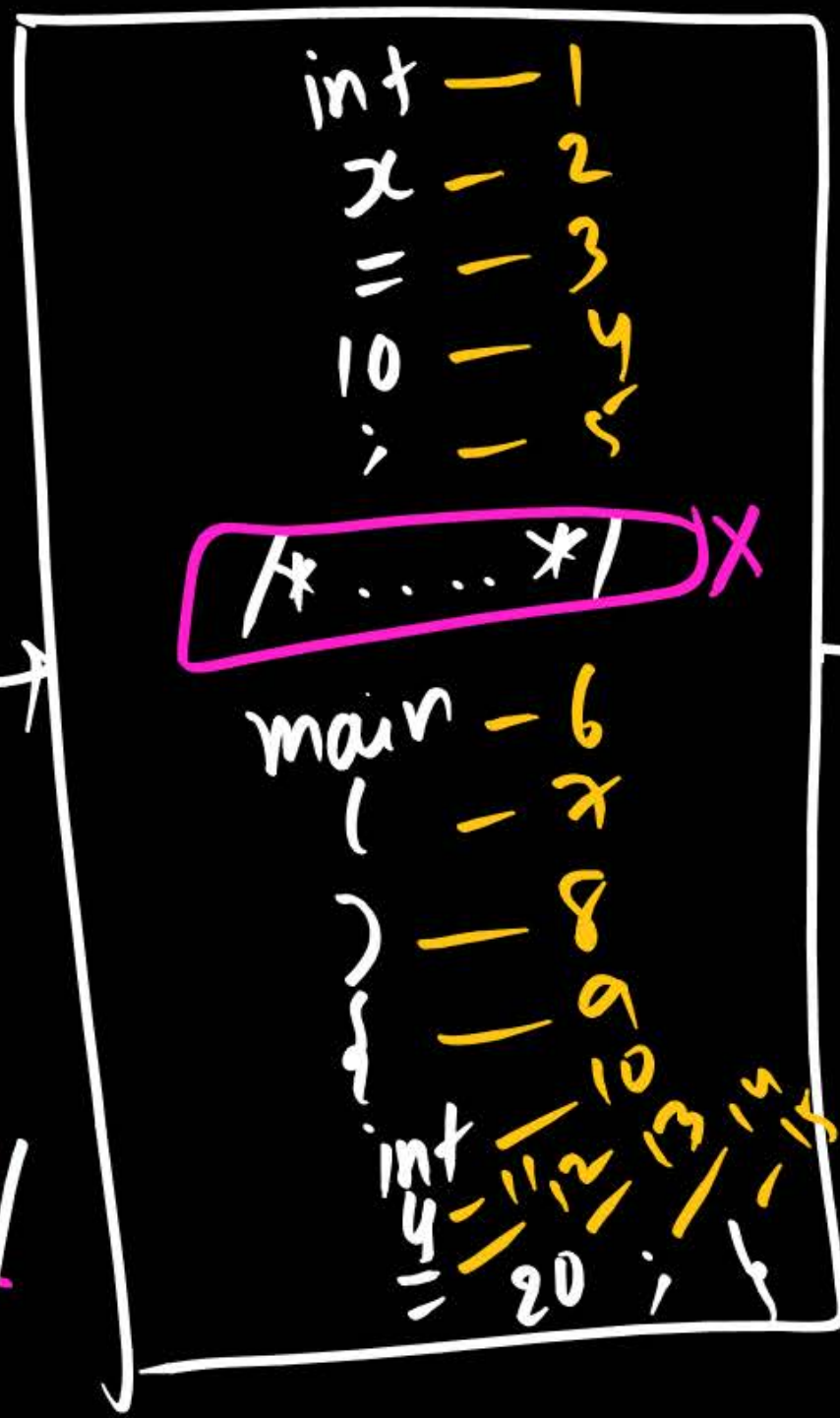
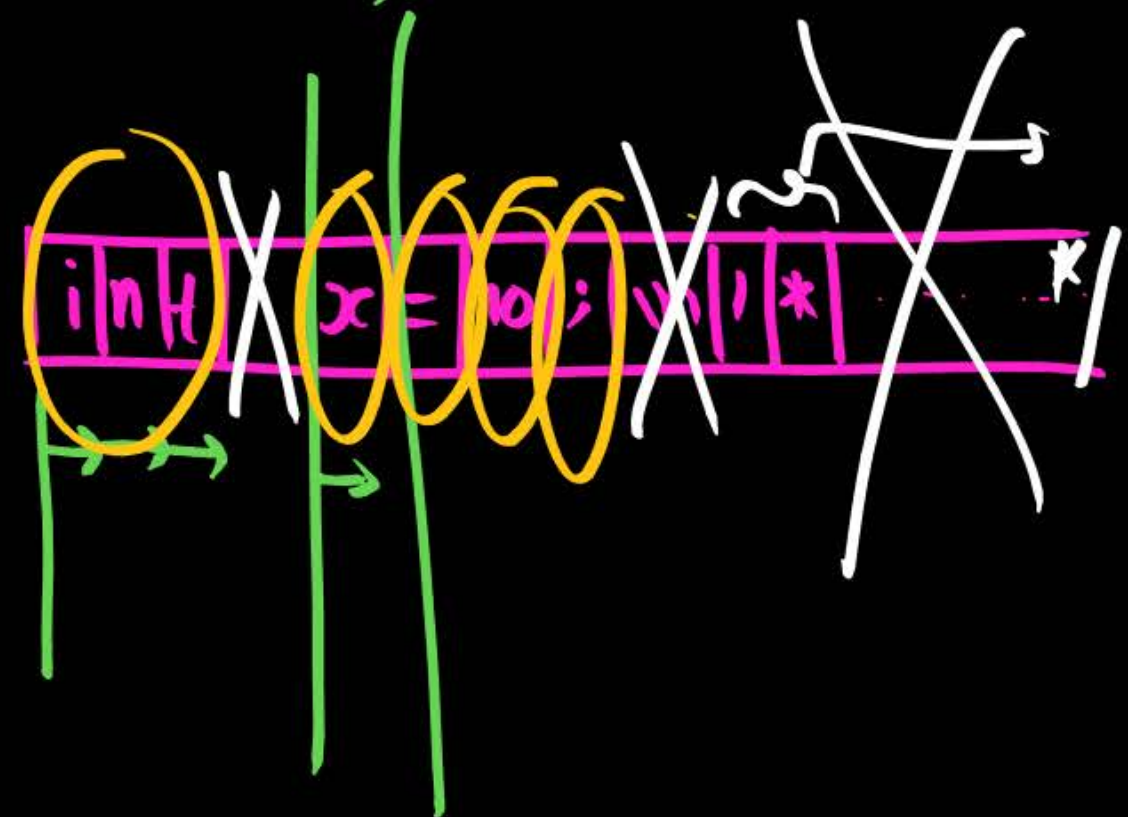
Lexical Analysis



- It takes sequence of characters as i/p and produces sequence of tokens
(character stream) (token stream)
- It may produce Lexical Error if any.
- It groups one or more characters into a token
- *** {
 - It uses longest prefix Rule (Maximal munch Rule) to produce token
 - It knows tokens, comments, white spaces
 - It produces tokens but ignores white spaces & comments (deletes)}



```
int x=10;  
/*comment*/  
main()  
{  
  int y=20;  
}
```



T₁ T₂ T₃ T₄ T₅ T₆ T₇ T₈ T₉ T₁₀ T₁₁ T₁₂ T₁₃ T₁₄ T₁₅

Lexeme: Smallest unit of program
TOKEN: type of Lexeme

Lexeme \cong TOKEN



Types of Tokens:

- i) can't start with digit
- ii) includes only letters, digits, underscore

1) Identifiers

- iii) should not include any other character
- iv) Keyword never be identifier

32 Keywords

2) Keywords

goto	default	case	struct	const
if	while		union	volatile
else	for		enum	sizeof
switch	break		register, auto,	
do	continue		extern, static	
	return			

3) Operators

+ , - , * , / , % , ++ , -- , < , <= , > , >= , == , != ,
&& , || , ! , & , | , ^ , ~ , << , >> , = , += , -= , ...

4) Literals (constants)

Integer → Decimal 235
 → Octal 0126
 → Hexa 0x1a2f

Floating
→ 20.34
→ 20.34f

character → 'A'
 → 'a'

string → "abc"

5) Special symbols

; , (,) , { , } , :

Comments



Diagram showing a box containing a diagonal line (/) and an asterisk (*), representing the beginning of a comment.

Begin
of comment

Diagram showing a box containing an asterisk (*) and a diagonal line (/), representing the end of a comment.

end
of comment

Deleted by
lexical analyzer

Lexical Error:

①

1abc → Invalid Sequence → not a token ⇒ Lexical Error
 a1bc → Identifier → Token
 1234 → Integer literal

Note: space always separates a token

②

1 abc ⇒ 2 tokens
 Integer Identifier

③

0[✓]2[✓]9^x3 → lexical error
 ↓
 octal only

0x2[✓]f[✓] → token
 ↓
 hexadecimal

0x2[✓]h^x → lexical error

013[✓]5[✓] → token

Lexical Error

④

/*

comment

Begin
of comment

is started

/*

There is no end rule
for comment

↓
Lexical error

⑤

/*

comment

Start

/* is started

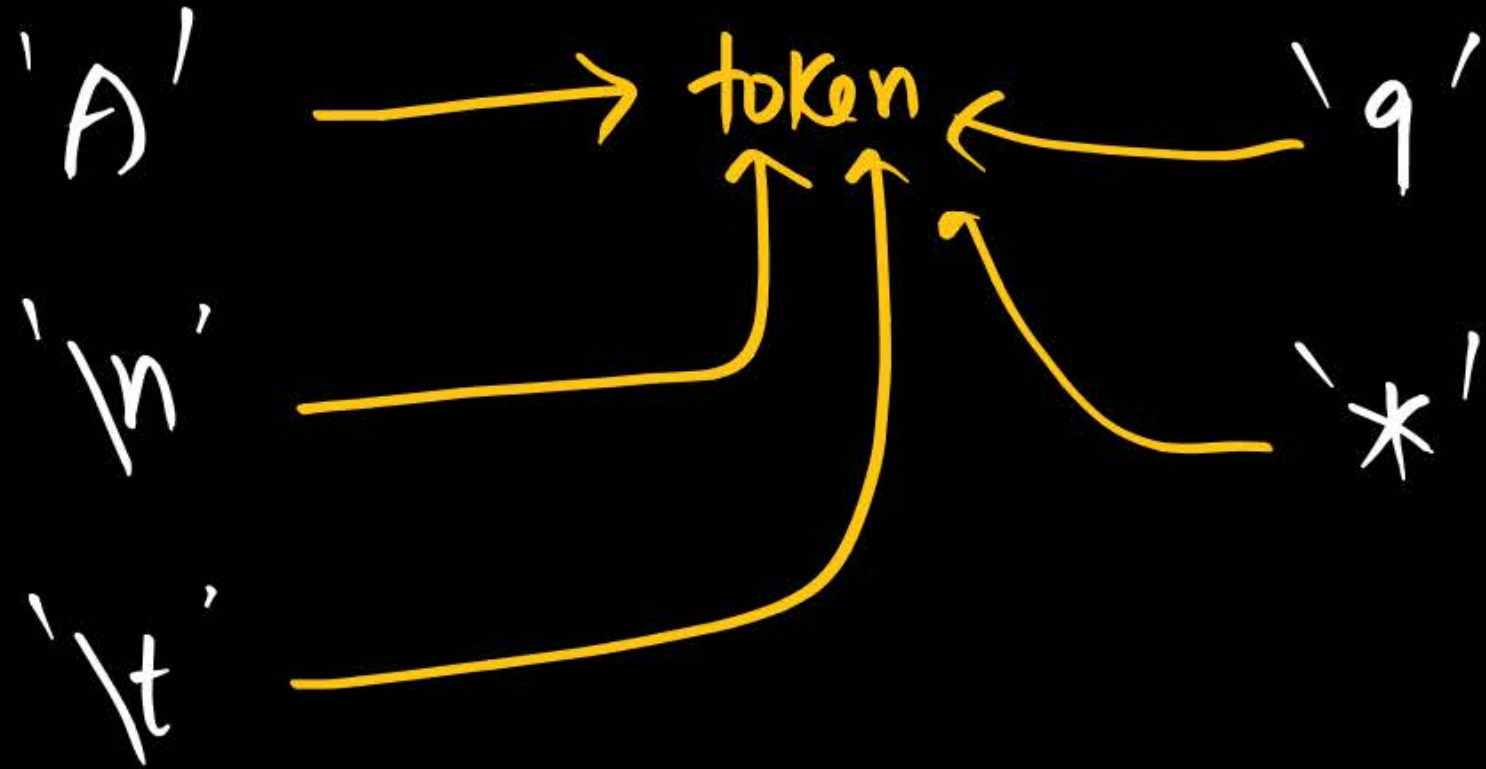
*/

end

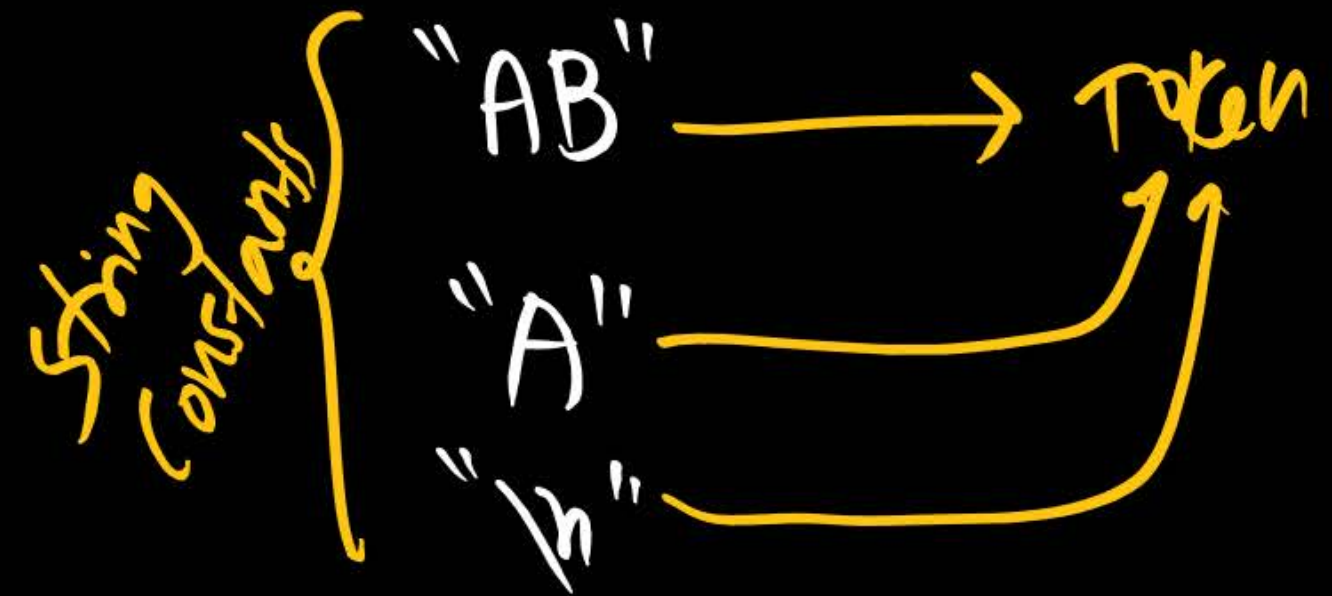
There is no lexical error

6

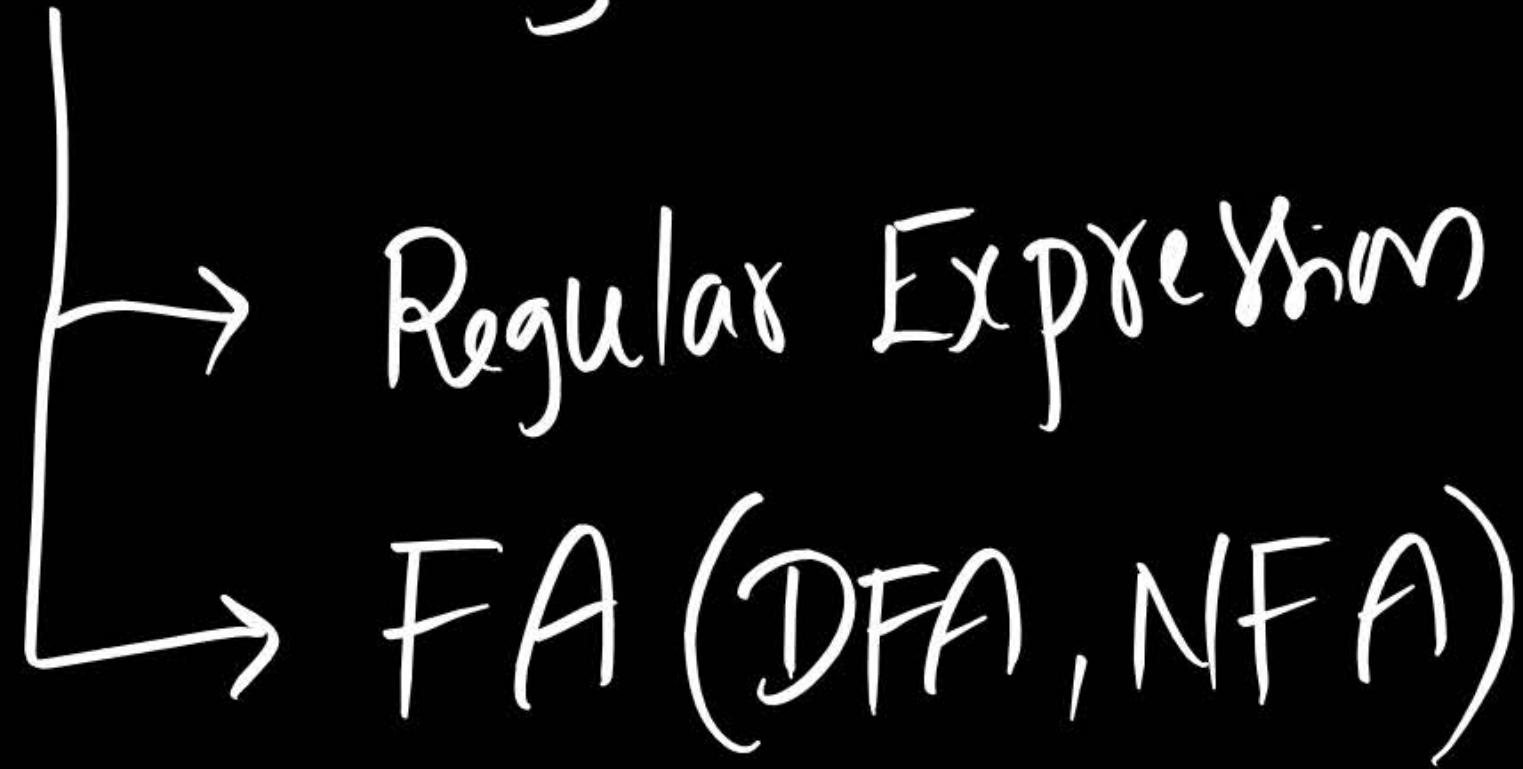
in
one character

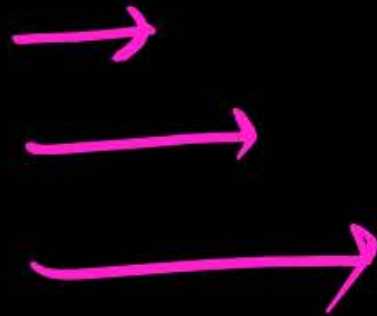
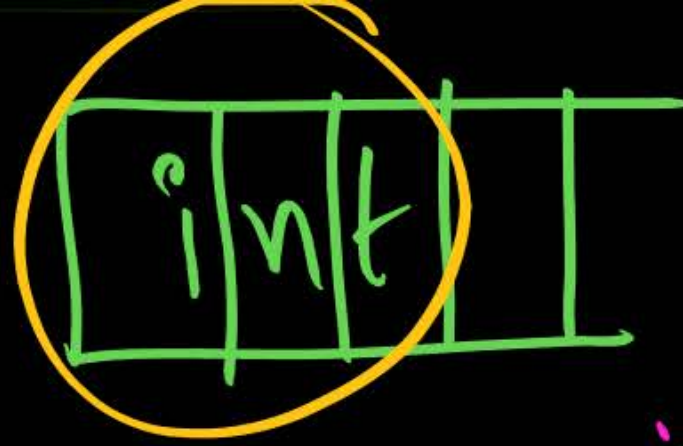


'AB' → Lexical error



Lexical Analysis





i ✓
in ✓
int ✓

longest prefix

Finding no. of tokens in the code.



① $\begin{matrix} T_1 & T_2 & T_3 \\ \boxed{\text{int}} & \boxed{x} & \boxed{;} \end{matrix} \Rightarrow 3 \text{ tokens}$ 

② $\boxed{\text{int}} \boxed{\text{temp}} \boxed{;} \Rightarrow 3$

③ $\underline{\text{int}} \underline{x}, \underline{y}; \Rightarrow 5$

Finding no. of tokens in the code.



④ main()
{ main() ;
} $\Rightarrow 9$

No compilation Error
Runtime error
(stack overflow)

⑤ int = float ; $\Rightarrow 4$

No lexical Error
Syntax Error
Compilation Error

⑥ int x = y ; $\Rightarrow 5$

No lexical error
No Syntax error
Semantic error (y is not declared)
Compilation error

int x = y;

y is undeclared
Declaration
error

lexeme	token	line	type	address
x	id ₁	100	int	
y	id ₂	100		

function names
variable names

Symbol Table:

→ It is data structure

→ It stores attribute information of each identifier

Depends on Scope, Symbol tables will be maintained.



Finding no. of tokens in the code.



⑦

x++;

⇒ 3 tokens

No lexical error

No syntax error

Semantic error (x is not declared)

⑧

int *x;

⇒ 4 tokens

no compilation error

⑨

int **x;

⇒ 5 tokens

no error

$(* (* x))$



** is not operator in C

* \rightarrow dereference

Finding no. of tokens in the code.



⑩

printf ("x=%d", x); \Rightarrow 7 tokens

Semantic Error

⑪

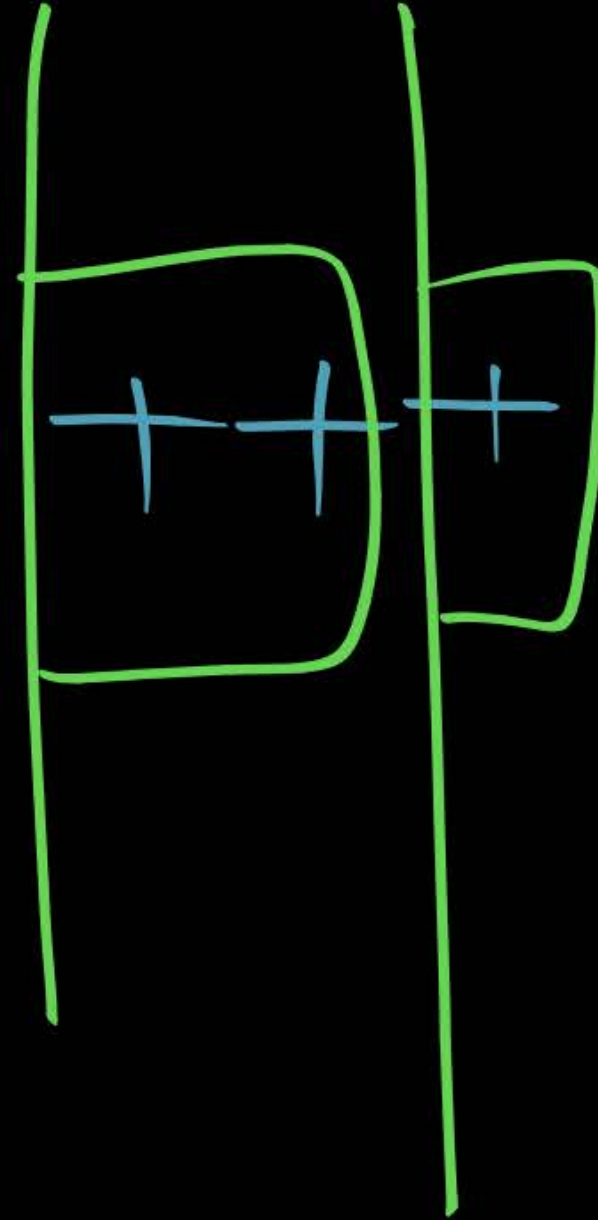
int x;
2 = x; \Rightarrow 7 tokens

Semantic Error
(lvalue required)
Compilation Error

⑫

int x;
x = x + + + 2; \Rightarrow 10 tokens

no Compilation error



2 plans



Finding no. of tokens in the code.



⑬ int x, y, z = 10;

\Rightarrow 9

no error

⑭ int x = y = z = 10;

\Rightarrow 9

Semantic error

⑮ int x = 1, y = 2;
x = (1, 2, 3);

\Rightarrow 19

no error

Finding no. of tokens in the code.



①⑥ /*comment
int x=10 /*comment*/ \Rightarrow 0 tokens
No error

①⑦ in /*comment */ t = 10; \Rightarrow in t = 10; \Rightarrow S
Semantic error

①⑧
int x;
char y;
y = x;
 \Rightarrow 10 tokens

int size 2 or 4
char size 1
Semantic error

in t = 10 ;

~~Syntax Error~~

Semantic Error

Note :

```
typedef int in ;
in t = 10 ;
```

no compilation error

Finding no. of tokens in the code.



(19)

int x=1;

char y;

y = (char) x;

⇒ 15 tokens

No error

(20)

char *p = "gate";

⇒ 6 tokens

Errors

Compile time Errors

Preprocessor Errors

`#include <stdio.h>`

not preprocessor command

Compilation Errors

- Lexical
- Syntax
- Semantic

Run time Errors

- File not found
- Divide by zero
- out of bound access in array
- Segmentation fault
- pointer related errors
- stack overflow
- infinite loop

Question:



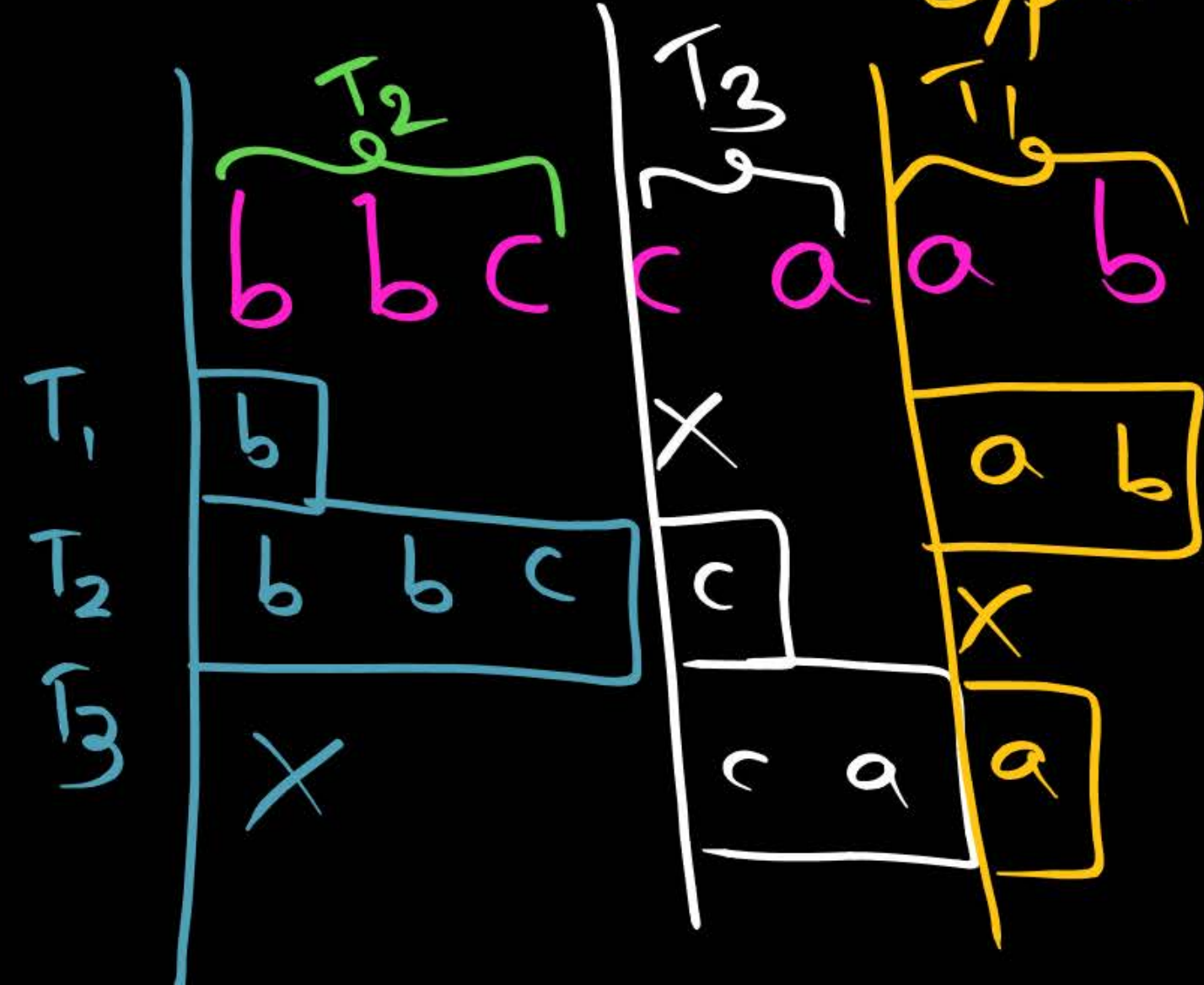
$T_1: a^*b$

$T_2: b^*c$

$T_3: c^*a$

I/P: bbccaaab

O/P: $T_2 T_3 T_1$



→ Lexical Analysis

Next: Syntax Analysis

