

# Lexical Analysis

Compiler Design

# Lexical Analysis

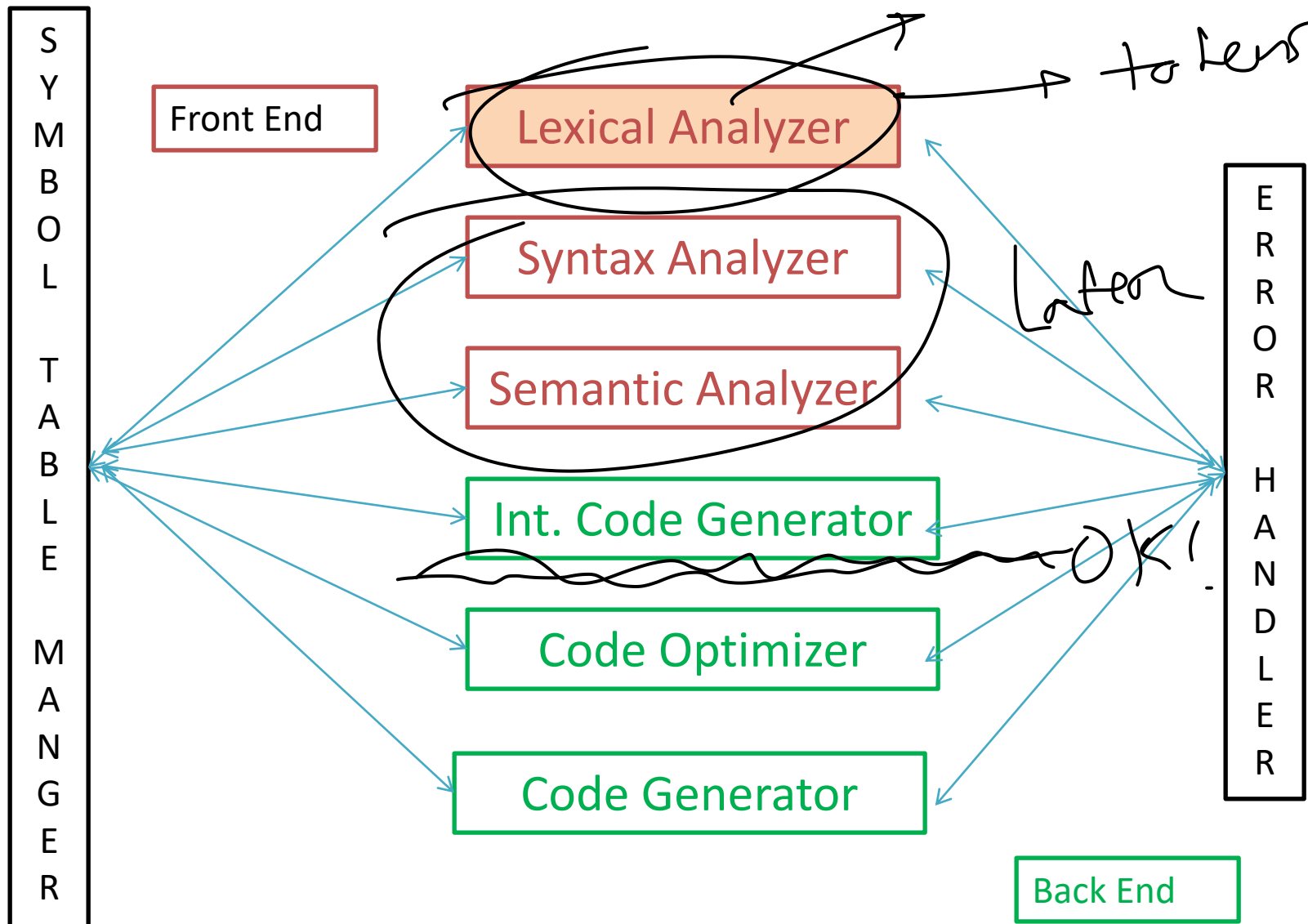
**Lexical analysis** is a method to convert high level input program into a sequence of Tokens.

The process or program which performs Lexical analysis is known as **Lexical analyzer** or **Scanner**.

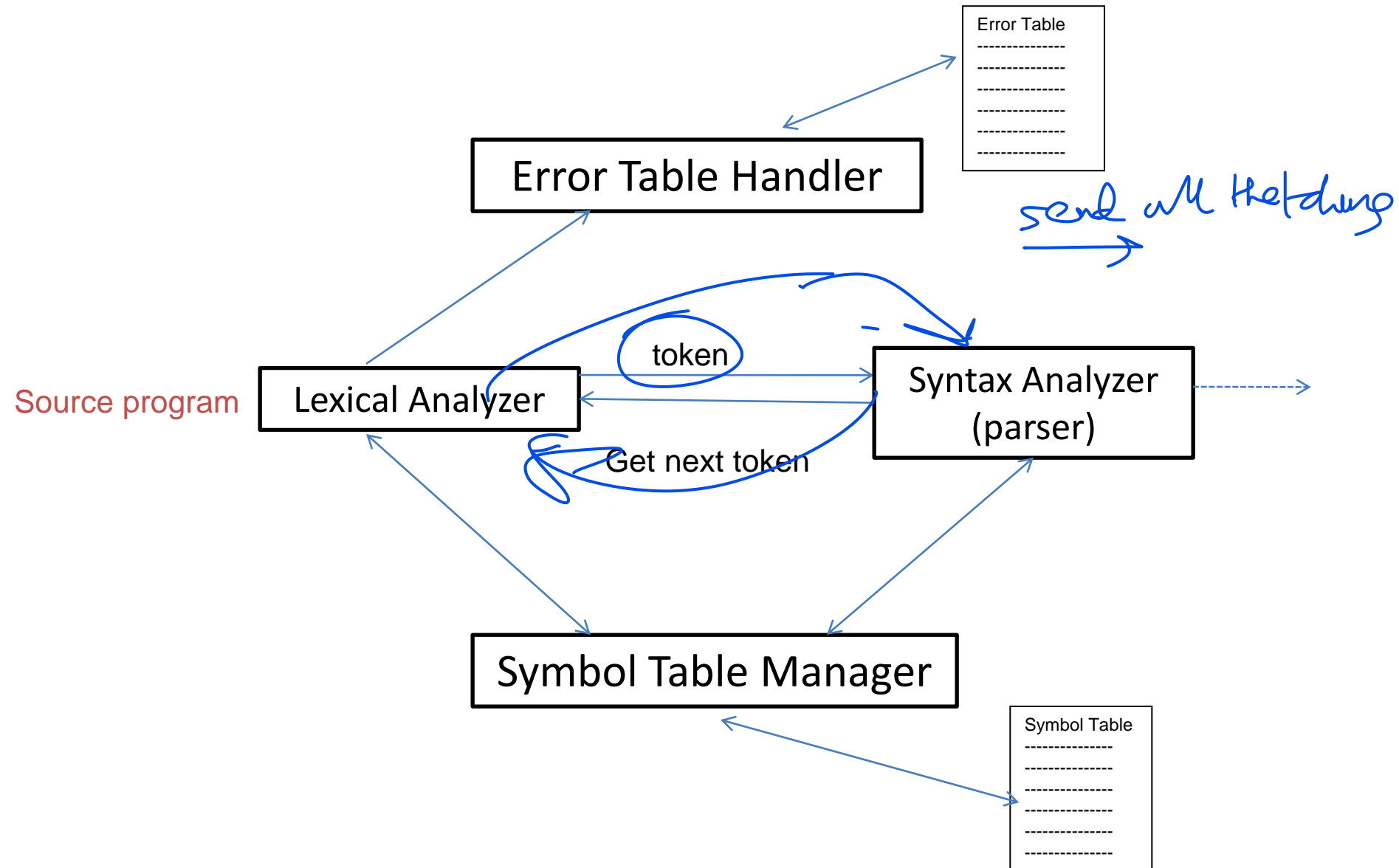
Example:

- *Lex* in Linux platform
- *Flex* in Windows platform
- You can write your own Lexical analyzer

# Structure of a Compiler



# Role of a Lexical Analyzer



# Objective

1. Simplicity
2. Efficiency
3. Portability

# Lexical Errors

Types of Errors-

- Lexical: misspelling of an identifier, keyword or operator.
- Syntactic: arithmetic expression with unbalanced parenthesis
- Semantic: operator applied to an incompatible operand
- Logical: infinite recursive call



# Error Recovery Strategies

- Panic mode error recovery
- Deleting an extraneous character
- Inserting a missing character
- Replacing an incorrect character by another
- Transposing two adjacent character

# Input Buffering

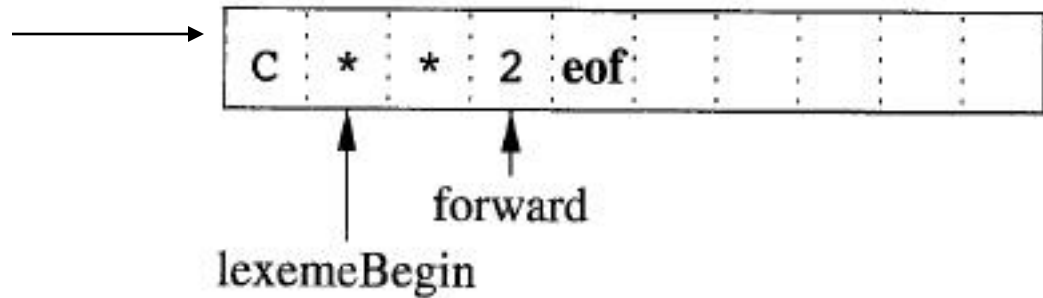
Instead of reading single character from a file a block of data is retrieved into a variable by system call.

Usually Buffer size is taken same as disk block size.

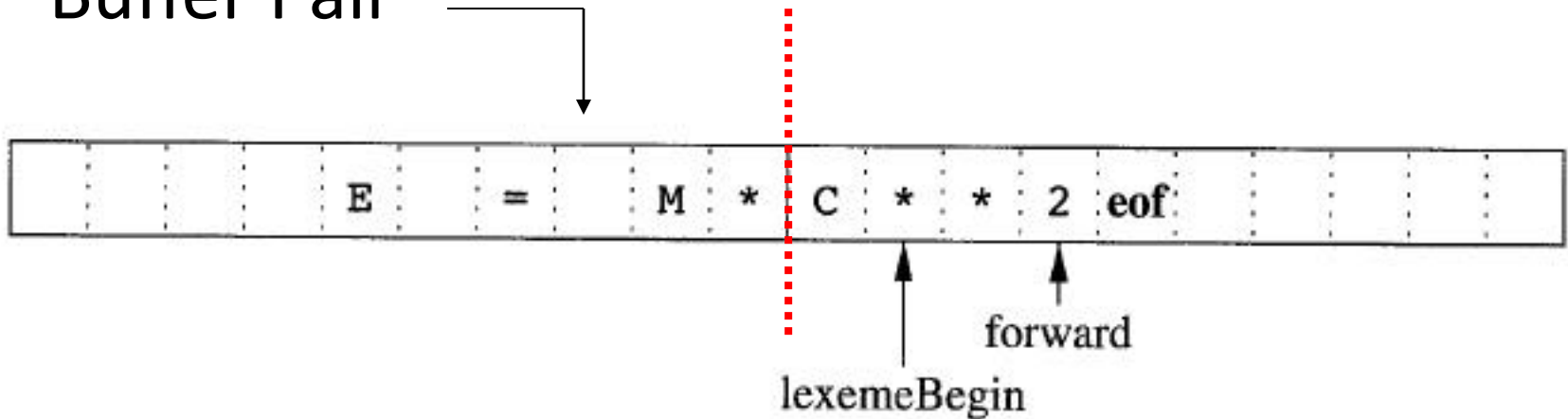


# Input Buffering types

- Single Buffer



- Buffer Pair



# Tokens, Patterns, Lexemes

## Tokens

**const**  
**if**  
**relation**  
**id**  
**num**  
**literal**

## Sample Lexemes

**const**  
**if**  
**<, <=, =, !=, >=**  
**pi, student\_name**  
**3.14156, 084499**  
**“Abdur Rahman”**

# Parts of string

- Prefix
- Suffix
- Substring
- Proper prefix, suffix, and substring
- Subsequence

# Operation on Languages

OPERATION	DEFINITION AND NOTATION
<i>Union of <math>L</math> and <math>M</math></i>	$L \cup M = \{s \mid s \text{ is in } L \text{ or } s \text{ is in } M\}$
<i>Concatenation of <math>L</math> and <math>M</math></i>	$LM = \{st \mid s \text{ is in } L \text{ and } t \text{ is in } M\}$
<i>Kleene closure of <math>L</math></i>	$L^* = \bigcup_{i=0}^{\infty} L^i$
<i>Positive closure of <math>L</math></i>	$L^+ = \bigcup_{i=1}^{\infty} L^i$

# Language

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

# Natural Language

- The syntax of a **natural language**, that is, a spoken language, such as English, French, German, or Spanish, is extremely complicated.
- In fact, it does not seem possible to specify all the rules of syntax for a natural language.

# Formal Language

- Formal language, which, unlike a natural language, is specified by a *well-defined set of rules of syntax*. Rules of syntax are important not only in linguistics, the study of natural languages, but also in the study of programming languages.
- This set of **rules of syntax** is called a grammar.

# Phrase-structure grammar

A phrase-structure grammar  $\mathbf{G} = \{\mathbf{V}, \mathbf{T}, \mathbf{S}, \mathbf{P}\}$ , where-

**V**: vocabulary

**T**: a subset of  $V$  consisting of terminal elements

**S**: start symbol

**P**: a set of productions or rules

**N**:  $V - T$ , is known as non-terminal symbols, every production must contain at least one nonterminal on its left side.



# Types of Phrase-structure grammar

It can be classified according to the types of productions that are allowed.

For productions:  $w1 \rightarrow w2$

- Type 0: no restrictions — phrase ~~structure~~ <sup>structure</sup>
- Type 1:  $len(w1) \leq len(w2)$  — context sensitive
- Type 2:  $w1 = A$  — context free
- Type 3:  $w1 = A$  and  $w2 =$  either  $aB$  or  $a$  — regular

$A, B$  are nonterminal and  $a$  is a terminal symbol.

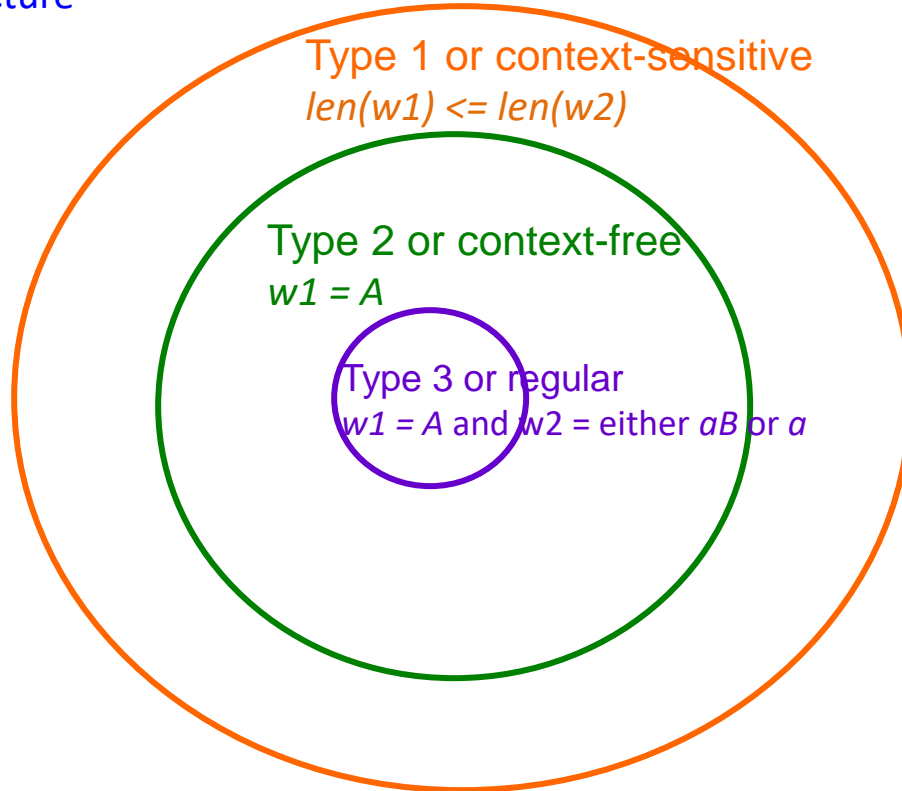
# Types of Phrase-structure grammar

Type 0 or phrase-structure

Type 1 or context-sensitive  
 $\text{len}(w1) \leq \text{len}(w2)$

Type 2 or context-free  
 $w1 = A$

Type 3 or regular  
 $w1 = A$  and  $w2 = \text{either } aB \text{ or } a$



# A Sample Grammar (context free)

1. *sentence* => *noun\_phrase* *verb\_phrase*
2. *noun\_phrase* => *article* *adjective* *noun* |
3. *article* *noun*
4. *verb\_phrase* => *verb* *adverb* |
5. *verb*

*the large rabbit hops quickly*

*the large rabbit hops quickly*

article

adjective

noun

verb

adverb



noun phrase



verb phrase



sentence

