

## Input Buffering :-

→ Deals about speedup the reading process.

### Reading Tokens :-

look one or more characters  
beyond the next Lexeme  
before we be sure that the  
right lexeme are seen.

Ex:- look atleast one characters before  
 == ~~we~~ before deciding an Identifier. additional

### Other issues in Reading input :-

\* Single character operator like  
 $-$ ,  $=$ ,  $<$  could be the begining of a  
 two character operators like,  
 $\neq$ ,  $==$ ,  $\geq$ ,  $\leq$ .

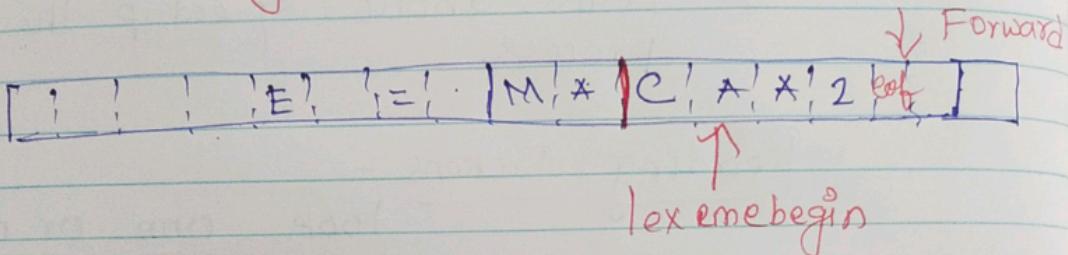
A Solution:- Two buffer Scheme  
 to Handle large  
 lookahead safely.

### Buffer Pair:-

- \* Reduce the amount of overhead to process a single input character
- \* a two buffer scheme that are loaded with input alternatively.

(31)

## Using a pair of Input Buffers



- \* Let buffer size is  $N$ , usually size of a disk block e.g 4096 bytes.
- \* ~~Instead of one read system command to N characters, into a buffer, instead of one system call per character.~~
- \* Eof: \* if  $N$  characters remain in the input file, then a special character represented by eof.  
↓  
marks the end of the source program
- \* Eof is a character, different from character present in source program.

### Pointers used in Buffer:-

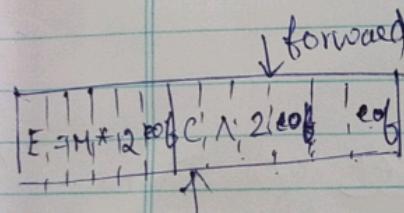
① Lexeme Begin:- \* marks the beginning of the current lexeme

\* Who extent we are attempting to determine.

② forward :- Scans until a pattern match is found.

- \* Once the next lexeme is determined, forward is set to the character at its right end.
- \* Then record the lexeme.
- \* Then lexeme begin is set to the next character immediately after the lexeme is found.

### Sentinels:-



check, each time when advance the forward pointer, that we have not moved off one of the buffers.

Specifications of Tokens:- Regular expression are an important notation for specifying lexeme patterns  
Strings:- An alphabet is a any finite set & Languages of symbols

### Examples of Symbols:-

- \* Letters
- \* digits
- \* punctuations.
- \*  $0, 1, 3 \rightarrow$  binary alphabet.
- \* ASCII
- \* UNICODE

String :- String over alphabet is a finite sequence of symbols drawn from that alphabet.

Sentence & word are same, used instead of string.

Let  $s$  is a string

\*  $|s| \rightarrow$  length of the String.

if  $s = \text{"banana"}$  then

$$|s| = 6,$$

\*  $\underline{\underline{\epsilon}} \rightarrow$  empty string (or)  $\emptyset$ .

$$|\epsilon| = 0$$

Language :- any Countable set of strings over some fixed alphabet.

Terms for Parts of Strings:

prefix :- string obtained by removing zero or more symbols from the end of  $s$ .

Example :-  $s = \text{"banana"}$

Prefix of  $s$  :-

ban, banana,  $\epsilon$ .

(34)

Suffix:- a string obtained by removing zero or more symbols from the beginning of  $s$ .

Suffix of  $s$  :- 'ana, banana, e, .

Suffix :-  $s$  is obtained by deleting any prefix and any suffix from  $s$

Suffix of  $s$  :- banana, na, e

Proper prefix, Suffix and Substring of  $s$  :-  
 $s$  (string itself),  $e$  are not part of prefix, Suffix, Substring.

Subsequence :-  $s$  is any string formed by deleting zero or more not necessarily consecutive positions of  $s$

Subsequences of  $s$  :- baan

Concatenation :-

Let  $x, y$  are two strings  
 Then  $xy$  is the concatenation.

Example :-

$x = \text{"dog"}, y = \text{"house"}$

$xy = \text{"doghouse"}$

Exponentiation :-

✓  $s^0 = \epsilon$

✓ If  $i > 0$ ,  $s^i$  to be  $s^{i-1}s$ .

✓  $\epsilon s = s$

✓  $s^1 = s$   
 ✓  $s^2 = ss$ ,  $s^3 = sss$ .

## Operations on Languages :-

Most important operations on languages:-

- (i) Union  $\cup$
- (ii) Concatenation
- (iii) closure.  $*$ ,  $^+$   
 Kleen closure      positive closure

### Operation.

Union of L and M

$L \cup M = \{s \mid s \in L \text{ or } s \in M\}$

Concatenation of L and M.

$L \cdot M = \{st \mid s \in L \text{ and } t \in M\}$

Kleen closure of L

$$L^* = \bigcup_{i=0}^{\infty} L^i$$

Positive closure of L.

$$L^+ = \bigcup_{i=1}^{\infty} L^i \quad \{ L^0 \text{ not included} \}$$

## Regular Expressions:-

Regular expressions describing all the languages that can be build from the operators ( $\cup, *, +, \text{any}$ ) applied to symbols of some alphabets.

letter  $\rightarrow$  any letter or underscore  
 digits  $\rightarrow$  any digit

Identifiers:- (Description):

letter - (letter - 1digit)\*

(Starting with letter then followed by letter (zero or more times)).

Rules that define Regular expressions:-

- \* Regular expression  $\sigma$
- \* Regular expression denoted by a language  $L(\sigma)$ .
- \* Regular expression  $\sigma$  is recursively built from smaller regular expression
- \*  $L(\sigma)$  is also recursively defined from the  $\sigma$ .

Two Rules:-

- ①  $\epsilon$  is a regular expression.  
 $L(\epsilon)$  is  $\{\epsilon\}$ .  $\rightarrow$  language whose sole member is the empty string.

⑥ If  $a$  is a symbol in  $\Sigma$ .

$$L(a) = \{a\}$$

(set of alphabet)

The language conti one string.

Use italic for symbols  
boldface for their corresponding regular expression.

Induction:-

- \* Four parts of induction
- \* Larger regular expressions are build from smaller ones.

Let  $r, s$  be regular expressions denoting languages  $L(r), L(s)$ ,

(i)  $(r)/(s)$  is a regular expression denoting the language  $L(r) \cup L(s)$

(ii)  $(r)s$  is a regular expression denoting the language  $L(r)L(s)$

(iii)  $(r)^*$  is a regular expression denoting  $(L(r))^*$

(iv)  $(r)$  is a regular expression denoting  $L(r)$

Parantheses  $\rightarrow$  add additional pair of around expressions

Without changing the language they denote.

NOTE:-

(a) The Unary operator  $*$  has highest precedence and is left associative.

(b) Concatenation has the second highest precedence and left associative.

(c)  $\cup$  has the lowest precedence and is left associative.  
or

Unnecessary parentheses must be removed.

RE (a)  $| (c b)^* c$ ) is replaced by  $a | b^* c$ .

Denotes set of string, single a or zero or more b followed by one c.

### Examples:

(i)  $a | b$  denotes the language  $\{a, b\}$

(ii)  $(a | b) (a | b)$  denotes  $\{aa, ab, ba, bb\}$

String of length two over  $\Sigma$ .

$(aa | bb | ba | bb)$  denotes the same language.

(iii)  $a^*$  → all strings of zero or more a's.

$$(\text{i.e.) } \{ \epsilon, a, aa, aaa \dots \}$$

(iv)  $(a|b)^*$  set of string consisting of zero or more instances of a and b.

$$\{ \epsilon, a, b, aa, ab, bb, ba, aaa \dots \}$$

(v)

$$(a^* b^*)^*$$

(v)  $a|a^*b$  denotes

$$\{ a, b, ab, aab, aaab \dots \}$$

Zero or more a and ending in b.

### NOTE:-

\* A language that can be defined by

\*  $\sigma = S$  RE  $\sigma$  and  $S$  are equivalent.

### Algebraic Laws for RE:-

#### LAW

#### Descriptions.

$$\sigma | \delta = \delta | \sigma$$

| is commutative

$$\sigma | (\delta | t) = (\sigma | \delta) | t$$

| is associative

(40)

$\gamma^{rest}$

$$\gamma(sst) = (\gamma s)t \quad \text{Concatenation is associative}$$

$$\begin{aligned} \gamma(slt) &= \gamma s|st; \\ (S|t)\gamma &= S\gamma|t\gamma \end{aligned} \quad \text{Concatenation is distributive over } |.$$

$$\gamma\epsilon = \gamma\epsilon = \gamma$$

$\epsilon$  is identity for concatenation

$$\gamma^* = (\gamma|\epsilon)^*$$

$\epsilon$  is guaranteed in closure.

$$\gamma^{**} = \gamma^*$$

\* is idempotent.

Example:-

i) Regular definition for the language of C identifiers.

$$\text{letter} \rightarrow [A|B|C|D|\dots|Z]-$$

$$\text{digit} \rightarrow [0|1|2|\dots|9]-$$

$$\text{id} \rightarrow \text{letter} - (\text{letter}|\text{digit})^*$$