

## Theory of Computational

### Languages

- In English, we distinguish 3 different entities: letters, words, and sentences.
  - Groups of letters make up words and groups of words make up sentences.
  - However, not all collections of letters form valid words, and not all collections of words form valid sentences.
- This situation also exists with computer languages.
  - Certain (but not all) strings of characters are recognizable words (e.g., IF, ELSE, FOR, WHILE ...); and certain (but not all) strings of words are recognizable commands.
- To construct a general theory of **formal languages**, we need to have a definition of a **language structure**, in which the decision of whether a given string of units constitutes a valid larger unit is not a matter of guesswork, but is based on explicitly stated rules.
- In this model, language will be considered as symbols with formal rules, and not as expressions of ideas in the minds of humans.
- The term “**formal**” emphasizes that it is the **form** of the string of symbols that we are interested in, not the meaning.

### Basic Definitions

**Alphabet** : A finite non-empty set of symbols (letters), is called an alphabet. It is denoted by  $\Sigma$  ( Greek letter sigma).

Example:  $\Sigma = \{a, b\}$

$\Sigma = \{0, 1\}$  //important as this is the language //which the computer understands.

$\Sigma = \{i, j, k\}$

**Strings:** Concatenation of finite symbols from the alphabet is called a string.

Example: If  $\Sigma = \{a, b\}$  then a, abab, aaabb, abababababababab

**Words:** Words are strings belonging to some language.

Example: If  $\Sigma = \{x\}$  then a language L can be defined as

$L = \{x^n : n=1, 2, 3, \dots\}$  or  $L = \{x, xx, xxx, \dots\}$

Here x, xx, ... are the words of L

(All words are strings, but not all strings are words).

### EMPTY STRING or NULL STRING

- We shall allow a string to have no letters. We call this **empty string** or **null string**, and denote it by the symbol  $\Lambda$ .

- For all languages, the **null word**, if it is a word in the language, is the word that has no letters. We also denote the null word by  $\Lambda$ .
- Two words are considered the same if all their letters are the same and in the same order.
- For clarity, we usually do not allow the symbol  $\Lambda$  to be part of the alphabet of any language.

### Discussion of null

- The language that has no words is denoted by the standard symbol for null set,  $\emptyset$ .
- It is not true that  $\Lambda$  is a word in the language  $\emptyset$  since this language has no words at all.
- If a certain language  $L$  does not contain the word  $\Lambda$  and we wish to add it to  $L$ , we use the operation “+” to form  $L + \{\Lambda\}$ . This language is **not** the same as  $L$ .
- However, the language  $L + \emptyset$  is the same as  $L$  since no new words have been added.

### Introduction to Defining Languages

- The rules for defining a language can be of two kinds:
  - They can tell us how to test if a string of alphabet letters is a valid word, or
  - They can tell us how to construct all the words in the language by some clear procedures.
- 

### Defining Languages

**Example:** Consider this alphabet with only one letter  $\Sigma = \{x\}$

- We can define a language by saying that any nonempty string of alphabet letters is a word  $L_1 = \{x, xx, xxx, xxxx, \dots\}$  or  $L_1 = \{x^n \text{ for } n = 1, 2, 3, \dots\}$

Note that because of the way we have defined it, the language  $L_1$  does not include the null word  $\Lambda$ .

**Example:** The language  $L$  of strings of odd length, defined over  $\Sigma = \{a\}$ , can be written as  $L = \{a, aaa, aaaaa, \dots\}$

**Example:** The language  $L$  of strings that does not start with  $a$ , defined over  $\Sigma = \{a, b, c\}$ , can be written as  $L = \{b, c, ba, bb, bc, ca, cb, cc, \dots\}$

**Example:** The language  $L$  of strings of length 2, defined over  $\Sigma = \{0, 1, 2\}$ , can be written as  $L = \{00, 01, 02, 10, 11, 12, 20, 21, 22\}$

**Example:** The language  $L$  of strings ending in 0, defined over  $\Sigma = \{0, 1\}$ , can be written as  $L = \{0, 00, 10, 000, 010, 100, 110, \dots\}$

**Example:** The language **EQUAL**, of strings with number of a's equal to number of b's, defined over  $\Sigma = \{a, b\}$ , can be written as:  $\{\Lambda, ab, aabb, abab, baba, abba, \dots\}$

**Example:** The language **EVEN-EVEN**, of strings with even number of a's and even number of b's, defined over  $\Sigma = \{a, b\}$ , can be written as:  $\{\Lambda, aa, bb, aaaa, aabb, abab, abba, baab, baba, bbaa, bbbb, \dots\}$

### Concatenation:

- Let us define an operation, **concatenation**, in which two strings are written down side by side to form a new longer string.  
xxx concatenated with xx is the word xxxxx  
 $x^n$  concatenated with  $x^m$  is the word  $x^{n+m}$
- For convenience, we may label a word in a given language by a new symbol. For example, xxx is called a, and xx is called b
- Then to denote the word formed by concatenating a and b, we can write  $ab = xxxxx$
- It is not true that when two words are concatenated, they produce another word. For example, if the language is  $L_2 = \{x, xxx, xxxxx, \dots\} = \{x^{2n+1} \text{ for } n = 0, 1, 2, \dots\}$   
then  $a = xxx$  and  $b = xxxxx$  are both words in  $L_2$ , but their concatenation  $ab = xxxxxxxx$  is not in  $L_2$

### Concatenation makes new Words?

- Note that in this simple example, we have:  $ab = ba$   
But in general, this relationship does NOT hold for all languages (e.g., **houseboat** and **boathouse** are two different words in English).

**Example:** Consider another language by beginning with the alphabet

$$\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$

### Define the language:

- $L_3 = \{\text{any finite string of alphabet letters that does not start with the letter zero}\}$
- This language  $L_3$  looks like the set of positive integers:  
 $L_3 = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, \dots\}$
- If we want to define  $L_3$  so that it includes the string (word) 0, we could say  
 $L_3 = \{\text{any finite string of alphabet letters that, if it starts with a 0, has no more letters after the first}\}$

### Definition: Length

- We define the function **length** of a string to be the number of letters in the string.

**Example:**

- If  $a = xxxx$  in the language  $L_1$ , then  $\text{length}(a) = 4$
- If  $c = 428$  in the language  $L_3$ , then  $\text{length}(c) = 3$
- If  $d = 0$  in the language  $L_3$ , then  $\text{length}(d) = 1$
- In any language that includes the null word  $\Lambda$ , then  $\text{length}(\Lambda) = 0$
- For any word  $w$  in any language, if  $\text{length}(w) = 0$  then  $w = \Lambda$ .
- Recall that the language  $L_1$  does not contain the null string  $\Lambda$ . Let us define a language like  $L_1$  but that does contain  $\Lambda$ :  
$$L_4 = \{ \Lambda, x, xx, xxx, xxxx, \dots \} = \{ x^n \text{ for } n = 0, 1, 2, 3, \dots \}$$
- Here we have defined that:  $x^0 = \Lambda$  (NOT  $x^0 = 1$  as in algebra)
- In this way,  $x^n$  always means the string of  $n$  alphabet letters  $x$ 's.
- Remember that even  $\Lambda$  is a word in the language, it is not a letter in the alphabet.

**Definition: Reverse:**

- If  $a$  is a word in some language  $L$ , then **reverse**( $a$ ) is the same string of letters spelled backward, even if this backward string is not a word in  $L$ .

**Example:**

- $\text{reverse}(xxx) = xxx$
- $\text{reverse}(145) = 541$
- Note that 140 is a word in  $L_3$ , but  $\text{reverse}(140) = 041$  is NOT a word in  $L_3$

**Definition: Palindrome:**

- Let us define a new language called Palindrome over the alphabet  $\Sigma = \{ a, b \}$   
 $\text{PALINDROME} = \{ \Lambda, \text{ and all strings } x \text{ such that } \text{reverse}(x) = x \}$
- If we want to list the elements in PALINDROME, we find  
 $\text{PALINDROME} = \{ \Lambda, a, b, aa, bb, aaa, aba, bab, bbb, aaaa, abba, \dots \}$

**Palindrome**

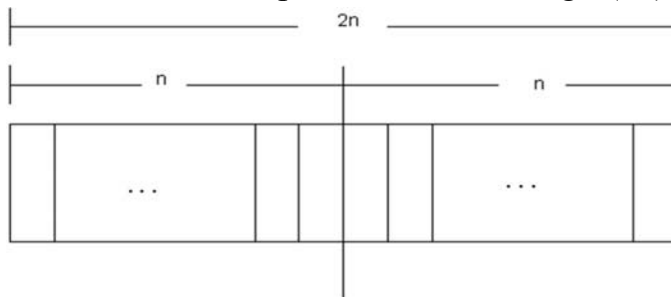
- Sometimes two words in PALINDROME when concatenated will produce a word in PALINDROME
  - **abba** concatenated with **abbaabba** gives **abbaabbaabba** (in PALINDROME)
- But more often, the concatenation is not a word in PALINDROME
  - **aa** concatenated with **aba** gives **aaaba** (NOT in PALINDROME)
- The language PALINDROME has interesting properties that we shall examine later.

### ❖ Task

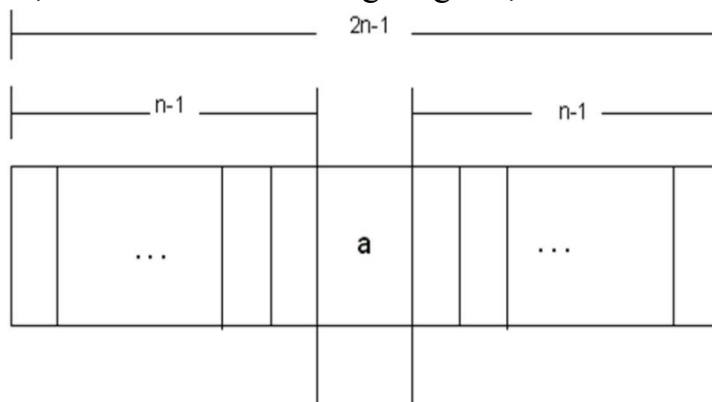
Q) Prove that there are as many palindromes of length  $2n$ , defined over  $\Sigma = \{a,b,c\}$ , as there are of length  $2n-1$ ,  $n = 1,2,3\dots$ . Determine the number of palindromes of length  $2n$  defined over the same alphabet as well.

### Solution

To calculate the number of palindromes of length  $(2n)$ , consider the following diagram,



- which shows that there are as many palindromes of length  $2n$  as there are the strings of length  $n$  *i.e.* the required number of palindromes are  $3^n$  (as there are three letters in the given alphabet, so the number of strings of length  $n$  will be  $3^n$ ).
- To calculate the number of palindromes of length  $(2n-1)$  with  $a$  as the middle letter, consider the following diagram,



- Which shows that there are as many palindromes of length  $2n-1$ , with  $a$  as middle letter, as there are the strings of length  $n-1$ , *i.e.* the required number of palindromes are  $3^{n-1}$ .

Similarly the number of palindromes of length  $2n-1$ , with  $b$  or  $c$  as middle letter, will be  $3^{n-1}$  as well. Hence the total number of palindromes of length  $2n-1$  will be:  $3^{n-1} + 3^{n-1} + 3^{n-1} = 3(3^{n-1}) = 3^n$ .

### ❖ Kleene Closure

**Definition:** Given an alphabet  $\Sigma$ , we define a language in which any string of letters from  $\Sigma$  is a word, even the null string  $\Lambda$ . We call this language the **closure** of the alphabet  $\Sigma$ , and denote this language by  $\Sigma^*$ .

**Examples:** If  $\Sigma = \{ x \}$  then  $\Sigma^* = \{ \Lambda, x, xx, xxx, \dots \}$

If  $\Sigma = \{ 0, 1 \}$  then  $\Sigma^* = \{ \Lambda, 0, 1, 00, 01, 10, 11, 000, 001, \dots \}$

If  $\Sigma = \{ a, b, c \}$  then  $\Sigma^* = \{ \Lambda, a, b, c, aa, ab, ac, ba, bb, bc, ca, cb, cc, aaa, \dots \}$

### Lexicographic order

- Notice that we listed the words in a language in size order (i.e., words of shortest length first), and then listed all the words of the same length alphabetically.
- This ordering is called **lexicographic** order, which we will usually follow.
- The star in the closure notation is known as the **Kleene star**.
- We can think of the Kleene star as an **operation** that makes, out of an alphabet, an *infinite* language (i.e., *infinitely many* words, each of *finite* length).

### Kleene Closure

- Let us now generalize the use of the Kleene star operator to sets of words, not just sets of alphabet letters.

**Definition:** If  $S$  is a set of words, then  $S^*$  is the set of all finite strings formed by concatenating words from  $S$ , where any word may be used as often as we like, and where the null string  $\Lambda$  is also included.

**Example:** If  $S = \{ aa, b \}$  then

$S^* = \{ \Lambda \text{ plus any word composed of factors of } aa \text{ and } b \}, \text{ or}$

$S^* = \{ \Lambda \text{ plus any strings of } a\text{'s and } b\text{'s in which the } a\text{'s occur in even clumps} \},$   
or

$S^* = \{ \Lambda, b, aa, bb, aab, baa, bbb, aaaa, aabb, baab, bbaa, bbbb, aaaab, aabaa, aabbb, baaaa, baabb, bbaab, bbbba, bbbbb, \dots \}$

Note that the string  $aabaaab$  is not in  $S^*$  because it has a clump of  $a$ 's of length 3.

**Example:** Let  $S = \{ a, ab \}$ . Then

$S^* = \{ \Lambda \text{ plus any word composed of factors of } a \text{ and } ab \}, \text{ or}$

$S^* = \{ \Lambda \text{ plus all strings of } a\text{'s and } b\text{'s except those that start with } b \text{ and those that contain a double } b \}, \text{ or}$

$S^* = \{ \Lambda, a, aa, ab, aaa, aab, aba, aaaa, aaab, abaa, abab, aaaaa, aaaab, aaaba, aabaa, aabab, abaaa, abaab, ababa, \dots \}$

- Note that for each word in  $S^*$ , every  $b$  must have an  $a$  immediately to its left, so the double  $b$ , that is  $bb$ , is not possible; neither any string starting with  $b$ .

### How to prove a certain word is in the closure language $S^*$

- We must show how it can be written as a concatenation of words from the base set  $S$ .
- In the previous example, to show that  $abaab$  is in  $S^*$ , we can factor it as follows:  
 $abaab = (ab)(a)(ab)$
- These three factors are all in the set  $S$ , therefore their concatenation is in  $S^*$ .

Note that the parentheses, ( ), are used for the sole purpose of demarcating the ends of factors.

- Observe that if the alphabet has no letters, then its closure is the language with the null string as its only word; that is  
if  $\Sigma = \emptyset$  (the empty set), then  $\Sigma^* = \{ \Lambda \}$
- Also, observe that if the set  $S$  has the null string as its only word, then the closure language  $S^*$  also has the null string as its only word; that is  
if  $S = \{ \Lambda \}$ , then  $S^* = \{ \Lambda \}$  because  $\Lambda\Lambda = \Lambda$ .
- Hence, the Kleene closure always produces an infinite language unless the underlying set is one of the two cases above.

### Kleene Closure of different sets

- The Kleene closure of two different sets can end up being the same language.

Example: Consider two sets of words  $S = \{ a, b, ab \}$  and  $T = \{ a, b, bb \}$

Then, both  $S^*$  and  $T^*$  are languages of all strings of  $a$ 's and  $b$ 's since any string of  $a$ 's and  $b$ 's can be factored into syllables of  $(a)$  or  $(b)$ , both of which are in  $S$  and  $T$ .

### Positive Closure

- If we wish to modify the concept of closure to refer only the concatenation of **some (not zero)** strings from a set  $S$ , we use the notation  $+$  instead of  $*$ .
- This “plus operation” is called **positive closure**.

Example: if  $\Sigma = \{ x \}$  then  $\Sigma^+ = \{ x, xx, xxx, \dots \}$

Observe that:

1. If  $S$  is a language that **does not** contain  $\Lambda$ , then  $S^+$  is the language  $S^*$  without the null word  $\Lambda$ .
2. If  $S$  is a language that **does** contain  $\Lambda$ , then  $S^+ = S^*$
3. Likewise, if  $\Sigma$  is an alphabet, then  $\Sigma^+$  is  $\Sigma^*$  without the word  $\Lambda$ .

### $S^{**}$ ?

- What happens if we apply the closure operator twice?  
– We start with a set of words  $S$  and form its closure  $S^*$

- We then start with the set  $S^*$  and try to form its closure, which we denote as  $(S^*)^*$  or  $S^{**}$

### **Theorem 1:**

For any set  $S$  of strings, we have  $S^* = S^{**}$

- Before we prove the theorem, recall from Set Theory that
  - $A = B$  if  $A$  is a subset of  $B$  **and**  $B$  is a subset of  $A$
  - $A$  is a subset of  $B$  if for all  $x$  in  $A$ ,  $x$  is also in  $B$

### **Proof of Theorem 1:**

- Let us first prove that  $S^{**}$  is a subset of  $S^*$ :  
Every word in  $S^{**}$  is made up of factors from  $S^*$ . Every factor from  $S^*$  is made up of factors from  $S$ . Hence, every word from  $S^{**}$  is made up of factors from  $S$ . Therefore, every word in  $S^{**}$  is also a word in  $S^*$ . This implies that  $S^{**}$  is a subset of  $S^*$ .
- Let us now prove that  $S^*$  is a subset of  $S^{**}$ :  
In general, it is true that for any set  $A$ , we have  $A$  is a subset of  $A^*$ , because in  $A^*$  we can choose as a word any factor from  $A$ . So if we consider  $A$  to be our set  $S^*$  then  $S^*$  is a subset of  $S^{**}$
- Together, these two inclusions prove that  $S^* = S^{**}$ .

### **Example:**

#### **Defining language of EVEN**

Step 1: 2 is in **EVEN**.

Step 2: If  $x$  is in **EVEN** then  $x+2$  and  $x-2$  are also in **EVEN**.

Step 3: No strings except those constructed in above, are allowed to be in **EVEN**.

### **Example:**

#### **Defining the language factorial**

Step 1: As  $0!=1$ , so 1 is in **factorial**.

Step 2:  $n!=n*(n-1)!$  is in **factorial**.

Step 3: No strings except those constructed in above, are allowed to be in **factorial**.

#### **Defining the language PALINDROME, defined over $\Sigma = \{a,b\}$**

Step 1:  $a$  and  $b$  are in **PALINDROME**



Step 2: if  $x$  is palindrome, then  $s(x)\text{Rev}(s)$  and  $xx$  will also be palindrome, where  $s$  belongs to  $\Sigma^*$

Step 3: No strings except those constructed in above, are allowed to be in palindrome

**Defining the language  $\{a^n b^n\}$ ,  $n=1,2,3,\dots$ , of strings defined over  $\Sigma=\{a,b\}$**

Step 1:  $ab$  is in  $\{a^n b^n\}$

Step 2: if  $x$  is in  $\{a^n b^n\}$ , then  $axb$  is in  $\{a^n b^n\}$

Step 3: No strings except those constructed in above, are allowed to be in  $\{a^n b^n\}$

**Defining the language  $L$ , of strings ending in  $a$ , defined over  $\Sigma=\{a,b\}$**

Step 1:  $a$  is in  $L$

Step 2: if  $x$  is in  $L$  then  $s(x)$  is also in  $L$ , where  $s$  belongs to  $\Sigma^*$

Step 3: No strings except those constructed in above, are allowed to be in  $L$

**Defining the language  $L$ , of strings beginning and ending in same letters, defined over  $\Sigma=\{a, b\}$**

Step 1:  $a$  and  $b$  are in  $L$

Step 2:  $(a)s(a)$  and  $(b)s(b)$  are also in  $L$ , where  $s$  belongs to  $\Sigma^*$

Step 3: No strings except those constructed in above, are allowed to be in  $L$

**Defining the language  $L$ , of strings containing  $aa$  or  $bb$ , defined over  $\Sigma=\{a, b\}$**

Step 1:  $aa$  and  $bb$  are in  $L$

Step 2:  $s(aa)s$  and  $s(bb)s$  are also in  $L$ , where  $s$  belongs to  $\Sigma^*$

Step 3: No strings except those constructed in above, are allowed to be in  $L$

**Defining the language  $L$ , of strings containing exactly  $aa$ , defined over  $\Sigma=\{a, b\}$**

Step 1:  $aa$  is in  $L$

Step 2:  $s(aa)s$  is also in  $L$ , where  $s$  belongs to  $b^*$

Step 3: No strings except those constructed in above, are allowed to be in  $L$