



# CSI3105: Software Testing

Module 2d: Strings/ Languages/ Reg Ex

Creative  
thinkers  
made here.

# String, Languages and Regular Expressions



1. Read the material from the slides (from the textbook) on BB as well as pages 98 and 99 of the textbook.

# String – Slides from the textbook

Strings play an important role in testing. A string serves as a test input.

- Examples: 1011; AaBc; “Hello world”.

## Alphabet

- A collection of symbols is known as an **alphabet**. We use an upper case letter such as X and Y to denote alphabets.
- Though alphabets can be infinite, we are concerned only with finite alphabets. For example,  $X=\{0, 1\}$  is an alphabet consisting of two symbols 0 and 1. Another alphabet is  $Y=\{\text{dog, cat, horse, lion}\}$  that consists of four symbols “dog”, “cat”, “horse”, and “lion”.

A string over an alphabet X is any sequence of zero or more symbols that belong to X. For example, 0110 is a string over the alphabet  $\{0, 1\}$ . Also, **dog cat dog dog lion** is a string over the alphabet  $\{\text{dog, cat, horse, lion}\}$ .

We will use lower case letters such as p, q, r to denote strings. The length of a string is the number of symbols in that string.

- Given a string s, we denote its length by  $|s|$ . Thus,  $|1011|=4$  and  $|\text{dog cat dog}|=3$ . A string of length 0, also known as an **empty string**, is denoted by  $\varepsilon$ .

# String – Slides from the textbook

## String Concatenation

- Let  $s_1$  and  $s_2$  be two strings over alphabet  $X$ . We write  $s_1.s_2$  to denote the **concatenation** of strings  $s_1$  and  $s_2$ .
- For example, given the alphabet  $X=\{0, 1\}$ , and two strings 011 and 101 over  $X$ , we obtain  $011.101=011101$ . It is easy to see that  $|s_1.s_2|=|s_1|+|s_2|$ . Also, for any string  $s$ , we have  $s.\epsilon=s$  and  $\epsilon.s=s$ .

# Languages - Slides from the textbook

A set  $L$  of strings over an alphabet  $X$  is known as a **language**. A language can be finite or infinite.

For example, a set of all strings consisting of zeros and ones is the language of binary numbers.

- The following sets are finite languages over the binary alphabet  $\{0, 1\}$ :
  - $\emptyset$ : The empty set
  - $\{\epsilon\}$ : A language consisting only of one string of length zero
  - $\{00, 11, 0101\}$ : A language containing three strings

# Regular Expressions - Slides from the textbook

A pattern of special characters (called metacharacters) which is used to match strings in a search. Any non-meta character matches itself

Examples:

- $/a|x /$  -- matches a, x, aa, aax, aaxx, aaa, xx, etc  
**(Note that | -- metacharacter for logical OR)**
- $abc^*$  --- ab, abc, abcc, abccc .....
- $(abc)^*$  -- “ ”, abc, abcabc, abcabcabc, ....
- $(a|b)(c|d)$  -- ac, ad, bc, bd
- $(0|1)^*$  --- all binary strings

Metacharacter	Matches...
.	<b>Any one character, except new line</b>
[a-z]	<b>Any one of the enclosed characters (e.g. a-z)</b>
*	<b>Zero or more of preceding character</b>
?	<b>Zero or one of the preceding characters</b>
+	<b>One or more of the preceding characters</b>

# Regular Expressions - Metacharacter

RE Metacharacter	Matches...
<code>^</code>	<b>beginning of line</b>
<code>\$</code>	<b>end of line</b>
<code>\char</code>	<b>Escape the meaning of <i>char</i> following it</b>
<code>[^]</code>	<b>One character <u>not</u> in the set</b>
<code>\&lt;</code>	<b>Beginning of word anchor</b>
<code>\&gt;</code>	<b>End of word anchor</b>
<code>( ) or \(\)</code>	<b>Tags matched characters to be used later (max = 9)</b>
<code>  or \ </code>	<b>Or grouping</b>
<code>x\{m\}</code>	<b>Repetition of character x, m times (x,m = integer)</b>
<code>x\{m,\}</code>	<b>Repetition of character x, at least m times</b>
<code>x\{m,n\}</code>	<b>Repetition of character x between m and n times</b>