



Universiteti i Prishtinës

# Programimi në Internet

## *Strings and Regular Expressions in PHP*

Dr. Ing. Lule Ahmedi

# ***Strings and their Workaround***

Questions that often arise:

- ⑥ How to format and manipulate text?
- ⑥ How to search words, phrases, or other patterns within a string?
- ⑥ How to store data into and retrieve them from a database
- ⑥ How to build search engines

In this lecture:

- ⑥ Formatting strings
- ⑥ Joining and splitting strings
- ⑥ Comparing strings
- ⑥ Matching and replacing substrings
- ⑥ Using regular expressions

# Formatting Strings

`trim()`, `ltrim()`, `chop()`:

- ⑥ Strip whitespaces (`\n`, `\r`, `\t`, `\v`, `\0`, spaces) from the start and/or end of a string
- ⑥ Characters other than whitespaces can be stripped out if specified as a second argument

HTML formatting with `nl2br()`:

- ⑥ Replaces all the newlines with `<br />` tag
- ⑥ Otherwise a single line displayed to the browser

Change the case of a string:

- ⑥ `strtoupper()`, `strtolower()`, `ucfirst()`, `ucwords()`

Escape special characters like quotes, backslashes, etc.:

- ⑥ `AddSlashes(string input)`, `StripSlashes(input)`

# Format a String for Printing

- ⑥ `echo()`, `print()` (the latter returns a boolean value)
- ⑥ `printf()` - format a string before printing it to the browser
- ⑥ `sprintf()` - format a string and return it
- ⑥ `printf("Total amount of order is %s (with shipping %.2f)", $total, $total_shipping);`
  - △ `%s` and `%.2f` are *type specifiers*

# Joining and Splitting Strings

- ⑥ Want to look at words in a sentence (say for spellchecking)
- ⑥ Want to split a domain name, a complete URL, or email address into its components
- ⑥ `array explode(string separator, string input [, int limit]);`
  - △ Splits a string *input* into (*number* of) pieces on a specified *separator* string

```
$email_array = explode('@', $email);  
  
if ($email_array[1]=='bigcustomer.com')  
$toaddress =  
    'bob@example.com'; else $toaddress = 'feedback@example.com';
```

## ***Joining and Splitting Strings (cont.)***

- ⑥ implode() or join() - the reverse effects of explode()
- ⑥ `$new_email = implode('@', $email_array);`
- ⑥ `strtok()` similar to `explode()`, but gets pieces (tokens) from a string one at a time
- ⑥ `string substr(string string, int start [, int length] );`
- ⑥ `substr('Your customer service is excellent', 0, 4)` would result into 'Your'

# Comparing Strings

- ⑥ `int strcmp(string str1, string str2);`
- ⑥ `strcasecmp( )` - is identical, but not case sensitive
- ⑥ `strnatcmp( )` - compares strings according to a "natural order", not according to a "lexicographical order"
- ⑥ `strlen(string length)` - gets the length of a string

# Matching and Replacing Substrings

- ⑥ `string strstr(string haystack, string needle);`
- ⑥ `int strpos(string haystack, string needle, int [offset] );`
- ⑥ `mixed str_replace(mixed needle, mixed new_needle, mixed haystack);`
  - △ Will replace all the instances of *needle* on *haystack* with *new\_needle* and return the new version of *haystack*
  - △ `$feedback = str_replace($offcolor, '%!@', $feedback);`
- ⑥ `string substr_replace(string string, string replacement, int start, int [length]);`
  - △ Will replace part (defined with the last two parameters) of *string* with *replacement*
  - △ `$test = substr_replace($test, 'X', -1);`



# Regular Expressions

Two styles of syntax supported in PHP:

- ⑥ POSIX style: compiled by default
  - △ Easier to learn and execute faster, but binary-unsafe
- ⑥ Perl-compatible regular expressions style : compiles in the PCRE library

## Pattern matching

- ⑥ Using string functions: *limited* to exact (substring) match
- ⑥ Regular expressions necessary for more complex ones

# *What are Regular Expressions?*

A way of describing a pattern in a piece of text

- ⑥ **Exact match**, e.g., searching for terms like "shop", "delivery"

Matching regexps is similar to `strstr()`, i.g., match a string somewhere within another string

- ⑥ **Substring match**, e.g., "shop" matches the regexp "shop", but also "h", "ho", ..

More advanced regexps:

- ⑥ Use **meta-characters** to indicate, say,
  - △ For a pattern to occur at the start or end of a string
  - △ For a part of the pattern to be repeated
  - △ For characters in a pattern to be of particular type

# Character Sets and Classes

**Character sets** - to match any character of a particular type (a kind of wildcard)

- ⑥ `'.'` - a wildcard for any single character (except `\n`)
  - △ `.at` matches `'cat'`, `'#at'`, `'sat'`, ..
  - △ Often used for filename matching in operating systems
- ⑥ `[a-z]`, `[a-zA-Z]`, or `[aeiou]` - a character class, i.e., a set of characters (can be a range, like is `[a-z]`, or a set of ranges, like is `[a-zA-Z]`) to which a (*single*) matched character must belong to; the latter is
- ⑥ `[^a-z]` - matches any character that is not between a and z
- ⑥ A number of **predefined character classes** exists:  
`[:alnum:]`, `[:alpha:]`, `[:space:]`, etc.

# *Repetition and Subexpressions*

A repetition:

- ⑥ There might be multiple occurrences of a particular string or class of characters
  - △ ' \* ' zero or more occurrences of a pattern
  - △ ' + ' one or more occurrences of a pattern
- ⑥ Example: `[ [ :alnum: ] ] +` means "at least one alphanumeric character"

Subexpressions:

- ⑥ Expressions within paranthesis to address patterns on a more fine granular basis
  - △ For each subexpression one pattern

# Counted Subexpressions

- ⑥  $\{n\}$  - exact number (n) of repetitions
- ⑥  $\{m, n\}$  - a range from m to n of repetitions
- ⑥  $\{m, \}$  - opened ended range of repetitions, at least m of them
  - △ Example:  $(\text{very})\{1, 3\}$  - matches 'very ', 'very very ', and 'very very very '

# ***Anchoring to the Beginning or End of a String***

- ⑥ The caret symbol '^' - a particular subexpression should appear at the start of a searched string
- ⑥ '\$' - analogously, at the end

## Examples:

- ⑥ ^bob - matches bob at the start of a string
- ⑥ com\$ - matches com at the end of a string
- ⑥ ^[a-z]\$ - matches any single character from a to z, in the string of its own

# Branching

- ⑥ A vertical pipe '|' - represents a choice in a regexp
  - △ (com) | (edu) | (net) - matches either of the strings within parentheses

# ***Matching Special Characters***

- ⑥ Special characters, such as '.', '{', or '\$'
  - △ To match them, a slash '\' should be put in front of it, say \., or \{, or \\$



# RegExpr Example #1

Detect particular terms in the customer feedback:

- ⑥ To match on 'shop', 'customer feedback', or 'retail'
  - △ Using string functions: three different searches, one for each term
  - △ Using regular expressions: in one pass
    - `shop|customer service|retail`

## ***RegExpr Example #2: E-mail Validator***

- ⑥ Encode the standardized format of an email address

`^[a-zA-Z0-9_\-\.]+@[a-zA-Z0-9\-\ ]+\.[a-zA-Z0-9_\-\.]+`

- ⑥ `^[a-zA-Z0-9_\-\.]+` - start the string with at least one letter, number, underscore, hyphen, or dot, or some combination of those (username)
- ⑥ `@` - matches a literal '@'
- ⑥ `[a-zA-Z0-9\-\ ]+` - matches the first part of the host name including alphanumeric characters and hyphens
- ⑥ `\.` - matches a literal '.' (slashed out since '.' is a special character, like `_` and `-` as well)
- ⑥ `[a-zA-Z0-9_\-\.]+$` - matches the rest of a domain name, including letters, numbers, hyphens, and more dots if required, up until the end of the string

# Finding Substrings

- ⑥ `int ereg(string pattern, string search, array [matches]);`
- ⑥ It searches for matches of the *pattern* in the *search* string
- ⑥ If searches are found for subexpressions of *pattern*, they will be stored in the array *matches*, one subexpression per array element
- ⑥ `eregi()` is identical, but case insensitive

# Finding Substrings Example

```
<?php
if (!ereg(
    ('^[a-zA-Z0-9_\-\.]+@[a-zA-Z0-9_\-\.]+$',
    $email))
{
    echo 'That's not a valid email address. Please return'
        .' to the previous page and try again.';
    exit;
}

$toaddress = 'feedback@example.com'; // default
if (ereg('shop|customer service|retail', $feedback))
    $toaddress = 'retail@example.com';
else if (ereg('deliver.*|fulfil.*', $feedback))
    $toaddress = 'fulfilment@example.com';
else if (ereg('bill|account', $feedback))
    $toaddress = 'account@example.com';

if (ereg('bigcustomer\.com', $email))
    $toaddress = 'bob@example.com';
?>
```

# Replacing Substrings

Similar to the `str_replace()` string function

- ⑥ `string ereg_replace(string pattern, string replacement, string search);`
- ⑥ It searches for the *pattern* in the *search* string and replaces it with the string *replacement*
- ⑥ `ereg_replace()` is identical, but not case sensitive

# Splitting Strings

- ⑥ `array split(string pattern, string search, int [max]) ;`
- ⑥ It splits the string *search* into substrings on the *pattern* and returns the (*max* number of) substrings in an array

To split the host name into its components:

```
$domain = 'cst.see-university.edu.mk';  
$arr = split ('\\.', $domain);  
  
while (list($key, $value) = each ($arr))  
    echo '<br/>'.$value;
```

# ***String Functions vs. RegExprs***

- ⑥ Regular expressions run less efficiently than the string functions with similar functionality
  - △  $\Rightarrow$  For simple applications, use string expressions where possible