



# Web programming

[www.singidunum.ac.rs](http://www.singidunum.ac.rs)



# Regular Expressions

- A pattern of special characters used to match strings in a search
- Typically made up from special characters called metacharacters
- Widely used in programming:
  - JavaScript form validation
  - Regular expressions are used throughout UNIX:
    - Editors: ed, ex, vi
    - Utilities: grep, egrep, sed, and awk

# Metacharacters

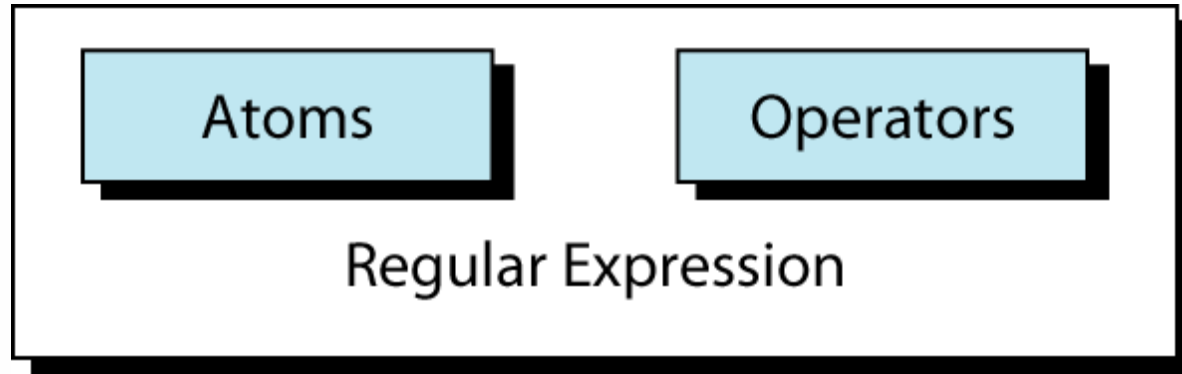
RE Metacharacter	Matches...
.	Any one character, except new line
[a-z]	Any one of the enclosed characters (e.g. a-z)
*	Zero or more of preceding character
? or \?	Zero or one of the preceding characters
+ or \+	One or more of the preceding characters

- any non-metacharacter matches itself

# more Metacharacters

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

# Regular Expression

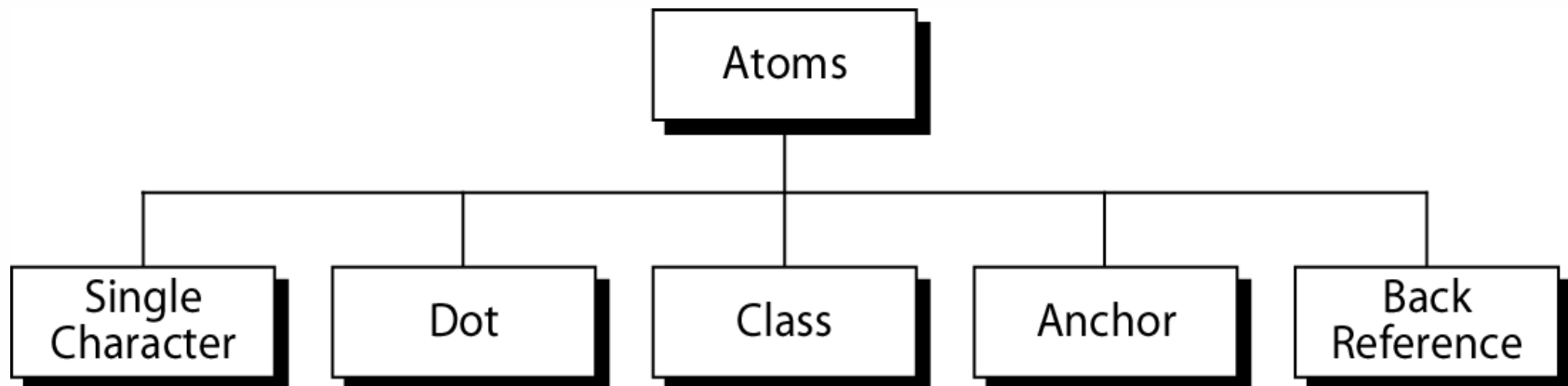


An atom specifies what text is to be matched and where it is to be found.

An operator combines regular expression atoms.

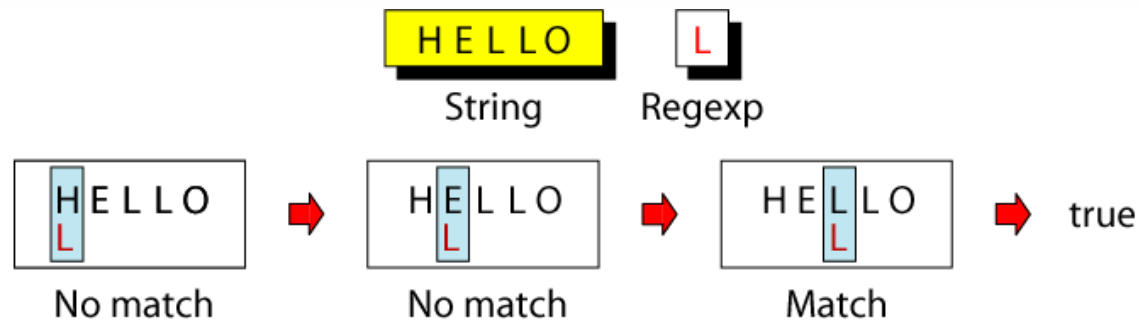
# Atoms

An atom specifies what text is to be matched and where it is to be found.

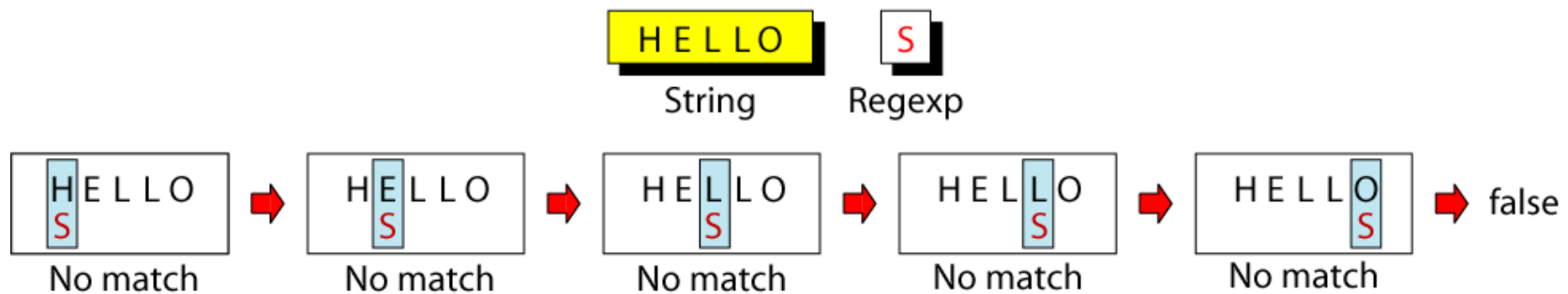


# Single-Character Atom

A single character matches itself



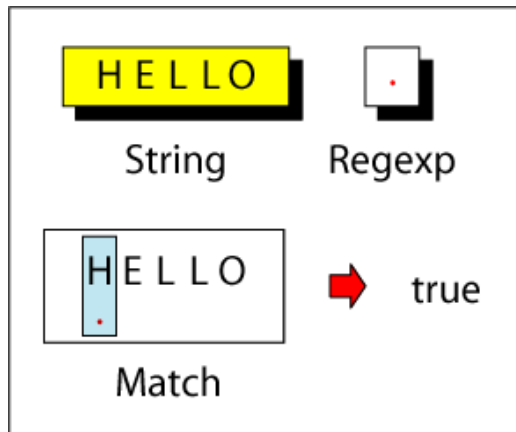
(a) Successful Pattern Match



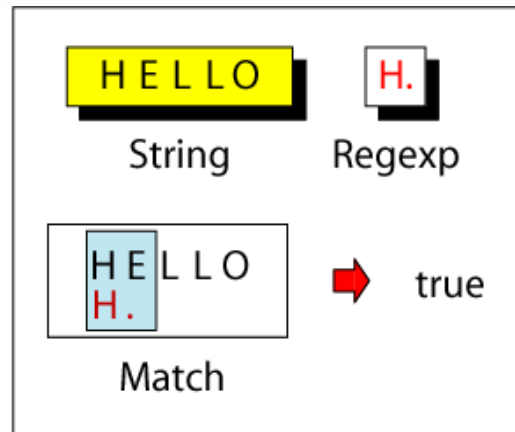
(b) Unsuccessful Pattern Match

# Dot Atom

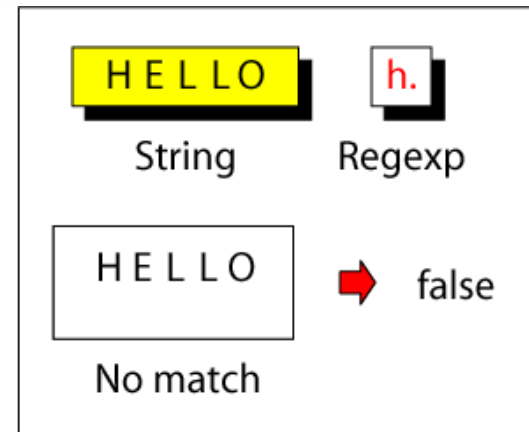
matches **any single character** except for a new line character (`\n`)



(a) Single-Character



(b) Combination-True



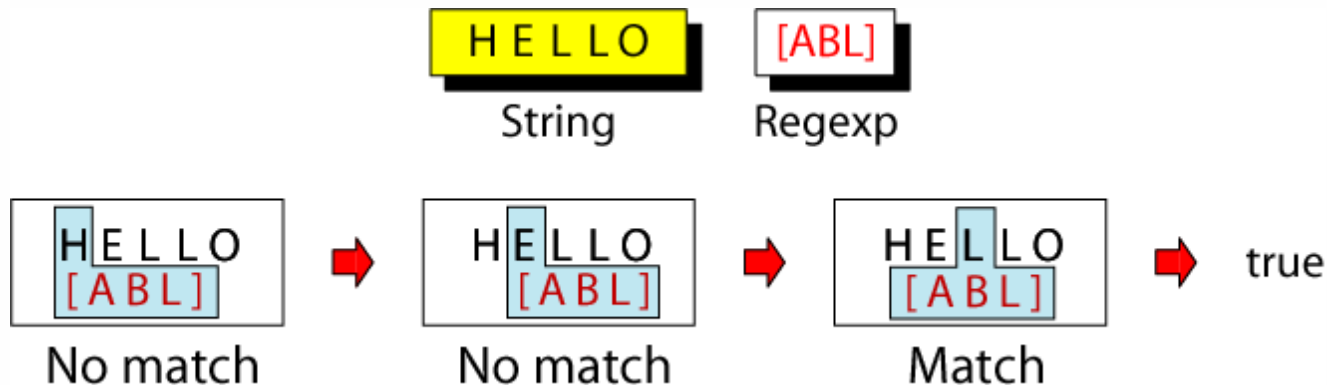
(c) Combination-False



# Class Atom

matches only single character that can be any of the characters defined in a set:

Example: `[ABC]` matches either A, B, or C.



Notes:

- 1) A range of characters is indicated by a dash, e.g. `[A-Q]`
- 2) Can specify characters to be excluded from the set, e.g. `[^0-9]` matches any character other than a number.

# Example: Classes

RegExpr		Means	RegExpr		Means
[A-H]	➔	[ABCDEFGH]	[^AB]	➔	Any character except A or B
[A-Z]	➔	Any uppercase alphabetic	[A-Za-z]	➔	Any alphabetic
[0-9]	➔	Any digit	[^0-9]	➔	Any character except a digit
[a]	➔	[ or a	]a]	➔	] or a
[0-9\ -]	➔	digit or hyphen	[^\^]	➔	Anything except^

# Shortcuts

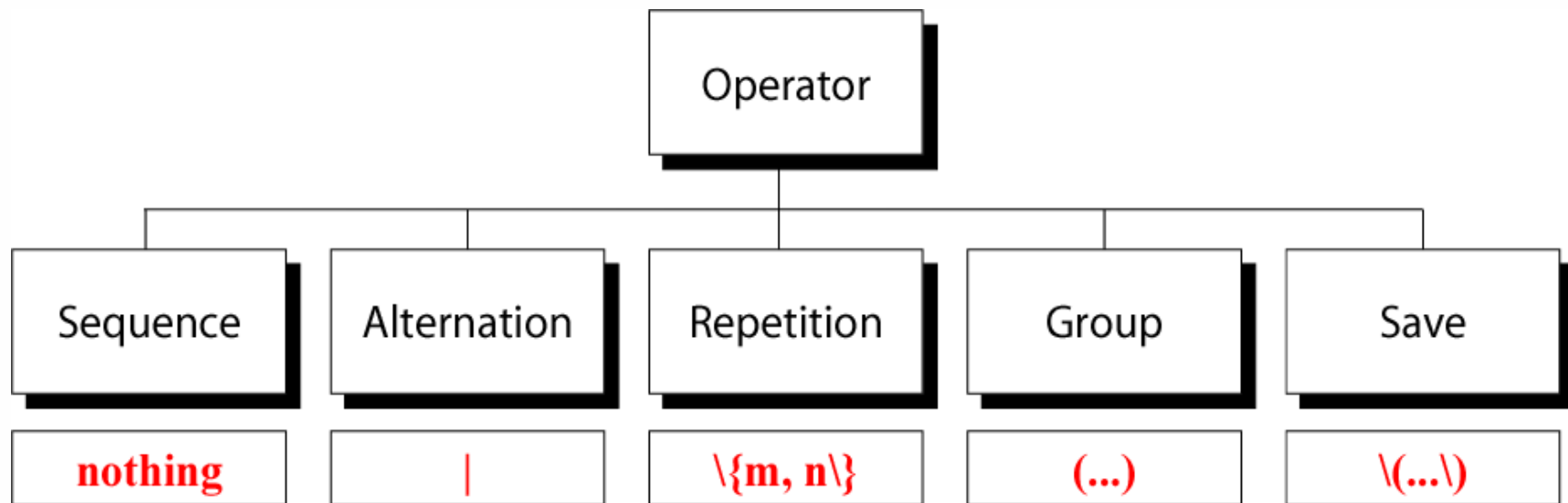
- `\d` => any digit `[0-9]`
- `\w` => any “word” character `[A-Za-z0-9_]`
- `\s` => any white space `[\t\n\r\f ]`
  
- `\D` => any character except a digit `[^\d]`
- `\W` => any character except a “word” character `[^\w]`
- `\S` => any character except a white space `[^\s]`
  
- Can use any of these in conjunction with quantifiers,
- `/\s*/` => any amount of white space

# Anchors

Anchors tell where the next character in the pattern must  
 be located in the text data.

Anchor		Means	Example
<code>^</code>	→	Beginning of line	One line of text. ↑
<code>\$</code>	→	End of line	One line of text. ↑
<code>\&lt;</code>	→	Beginning of word	One line of text. ↑ ↑ ↑ ↑
<code>\&gt;</code>	→	End of word	One line of text. ↑ ↑ ↑ ↑

# Operators



# Sequence Operator

In a sequence operator, if a series of atoms are shown in a regular expression, there is no operator between them.

<code>dog</code>	➔	matches the pattern "dog"
<code>a..b</code>	➔	matches "a" , any two characters, and "b"
<code>[2-4][0-9]</code>	➔	matches a number between 20 and 49
<code>[0-9][0-9]</code>	➔	matches any two digits
<code>^\$</code>	➔	matches a blank line
<code>^.\$</code>	➔	matches a one-character line
<code>[0-9]-[0-9]</code>	➔	matches two digits separated by a "-"

# Alternation Operator: | or \

operator (| or \ ) is used to define one  
**or** more alternatives

**UNIX|unix**



matches "UNIX" or "unix"

**Ms|Miss|Mrs**



matches "Ms" or "Miss" or "Mrs"

# Repetition Operator: $\{...\}$

The repetition operator specifies that the atom or expression immediately before the repetition may be repeated.

$\{m, n\}$

matches previous character m to n times.

$A\{3, 5\}$



matches "AAA", "AAAA", or "AAAAA"

$BA\{3, 5\}$



matches "BAAA", "BAAAA", or "BAAAAA"



# Basic Repetition Forms

## Formats

`\{m\}`



matches previous atom exactly m times

`\{m, \}`



matches previous atom m times or more

`\{, n\}`



matches previous atom n times or less

## Examples

`CA\{5\}`



CAAAAA

`CA\{3, \}`



CAAA, CAAAA, CAAAAA, ...

`CA\{, 2\}`



C, CA, CAA

# Short Form Repetition Operators: \*

## Formats

<b>*</b>	→	special case: matches previous atom zero or more times
<b>+</b>	→	special case: matches previous atom one or more times
<b>?</b>	→	special case: matches previous atom 0 or one time only

## Examples

<b>BA*</b>	→	B, BA, BAA, BAAA, BAAAA, ...
<b>B.*</b>	→	B, BA ... BZ, BAA ... BZZ, BAAA ... BZZZ, ...
<b>.*</b>	→	zero or more characters
<b>.*</b>	→	one or more characters
<b>[0-9]?</b>	→	zero or one digit

# Group Operator

In the group operator, when a group of characters is enclosed in parentheses, the next operator applies to the whole group, not only the previous characters.

Regex		Matches
<code>A(BC)\{3\}</code>	→	<code>ABCBCBC</code>
<code>(F(BC)\{2\}G)\{2\}</code>	→	<code>FBCBCGFBCBCG</code>