# Regular Expressions Exercises

What will the following code print?

```
import re
line = "IT 102 students: Alan Jason Farihah Liam William"
results = re.findall('W[a-z]+', line, re.I|re.M)
if results:
    print ("results: ", results)
else:
    print ("Nothing found")
```

Download the first book of the Bible:
http://www.vatican.va/archive/bible/genesis/documents/bible_genesis_e
n.html and save it as Genesis.txt


**Exercises:**

1. Search for all words in Genesis that start with capital "G". Note:
if you want only unique results, you can convert results to a set:
results = set(results)

2. Search for word(s) in Genesis that consist of 15 or more
characters. What have you found?

3. Search for word(s) that contain "app" or "comp" substrings. What
words have you found?

3. Complete the first three exercises at http://regex.alf.nu

# regularexpressions

## Anchors

| | |
|---|---|
| ^ | Start of string |
| \A | Start of string |
| $ | End of string |
| \Z | End of string |
| \b | Word boundary |
| \B | Not word boundary |
| \< | Start of word |
| \> | End of word |

## Character Classes

| | |
|---|---|
| \c | Control character |
| \s | White space |
| \S | Not white space |
| \d | Digit |
| \D | Not digit |
| \w | Word |
| \W | Not word |
| \x | Hexadecimal digit |
| \O | Octal digit |

## POSIX

| | |
|---|---|
| [:upper:] | Upper case letters |
| [:lower:] | Lower case letters |
| [:alpha:] | All letters |
| [:alnum:] | Digits and letters |
| [:digit:] | Digits |
| [:xdigit:] | Hexadecimal digits |
| [:punct:] | Punctuation |
| [:blank:] | Space and tab |
| [:space:] | Blank characters |
| [:cntrl:] | Control characters |
| [:graph:] | Printed characters |
| [:print:] | Printed characters and spaces |
| [:word:] | Digits, letters and underscore |

## Assertions

| | |
|---|---|
| ?= | Lookahead assertion |
| ?! | Negative lookahead |
| ?<= | Lookbehind assertion |
| ?!= or ?<! | Negative lookbehind |
| ?> | Once-only Subexpression |
| ?() | Condition [if then] |
| ?()| | Condition [if then else] |
| ?# | Comment |

## Quantifiers

| | |
|---|---|
| * | 0 or more |
| + | 1 or more |
| ? | 0 or 1 |
| {3} | Exactly 3 |
| {3,} | 3 or more |
| {3,5} | 3, 4 or 5 |

## Quantifier Modifiers

"x" below represents a quantifier

| | |
|---|---|
| x? | Ungreedy version of "x" |

## Escape Character

| | |
|---|---|
| \ | Escape Character |

## Metacharacters (must be escaped)

| | | |
|---|---|---|
| ^ | [ | . |
| $ | { | * |
| ( | \ | + |
| ) | | | ? |
| < | > | |

## Special Characters

| | |
|---|---|
| \n | New line |
| \r | Carriage return |
| \t | Tab |
| \v | Vertical tab |
| \f | Form feed |
| \xxx | Octal character xxx |
| \xhh | Hex character hh |

## Sample Patterns

| Pattern | Will Match |
|---|---|
| ([A-Za-z0-9-]+) | Letters, numbers and hyphens |
| (\d{1,2}\/\d{1,2}\/\d{4}) | Date (e.g. 21/3/2006) |
| ([^\s]+(?=\.(jpg\|gif\|png))\.\2) | jpg, gif or png image |
| (^[1-9]{1}$\|^[1-4]{1}[0-9]{1}$\|^50$) | Any number from 1 to 50 inclusive |
| (#?([A-Fa-f0-9]){3}(([A-Fa-f0-9]){3})?) | Valid hexadecimal colour code |
| ((?=.*\d)(?=.*[a-z])(?=.*[A-Z]).{8,15}) | String with at least one upper case letter, one lower case letter, and one digit (useful for passwords). |
| (\w+@[a-zA-Z_]+?\.[a-zA-Z]{2,6}) | Email addresses |
| (\<(/?[^\>]+)\>) | HTML Tags |

Note: These patterns are intended for reference purposes and have not been extensively tested. Please use with caution and test thoroughly before use.

## Groups and Ranges

| | |
|---|---|
| . | Any character except new line (\n) |
| (a\|b) | a or b |
| (...) | Group |
| (?:...) | Passive Group |
| [abc] | Range (a or b or c) |
| [^abc] | Not a or b or c |
| [a-q] | Letter between a and q |
| [A-Q] | Upper case letter between A and Q |
| [0-7] | Digit between 0 and 7 |
| \n | nth group/subpattern |

Note: Ranges are inclusive.

## Pattern Modifiers

| | |
|---|---|
| g | Global match |
| i | Case-insensitive |
| m | Multiple lines |
| s | Treat string as single line |
| x | Allow comments and white space in pattern |
| e | Evaluate replacement |
| U | Ungreedy pattern |

## String Replacement (Backreferences)

| | |
|---|---|
| $n | nth non-passive group |
| $2 | "xyz" in /^(abc(xyz))$/ |
| $1 | "xyz" in /^(?:abc)(xyz)$/ |
| $` | Before matched string |
| $' | After matched string |
| $+ | Last matched string |
| $& | Entire matched string |

**Homework**

Complete as many levels at regex.alf.nu as you can.