# CMSC 330: Organization of Programming Languages

#### Ruby Regular Expressions

#### String Processing in Ruby

 Scripting languages provide many useful libraries for manipulating strings

- The Ruby String class provides many useful, such as
  - Concatenate two strings
  - grabbing substrings
  - Searching and Replacing

# String Operations in Ruby

- What if we want to find more complicated patterns?
  - Find Steve, Stephen, Steven, Stefan, Esteve
  - Count the number of words that have even number vowels.

We need Regular Expressions

#### Regular Expressions

- A way of describing patterns or sets of strings
  - Searching and matching
  - Formally describing strings
    - > The symbols (lexemes or tokens) that make up a language
- Common to lots of languages and tools
  - awk, sed, perl, grep, Java, OCaml, C libraries, etc.
    - > Popularized (and made fast) as a language feature in Perl
- Based on some really elegant theory
  - Future lecture

#### Regular Expressions

- Regular expressions consist of
  - Constants
    - > empty set Ø:
    - empty string ε
    - > literal character a in Σ,a finite alphabet
  - Operations over these sets
    - > Concatenation: a b
    - > Union: a | b
    - Kleene star: a\*
  - We can build complicated patterns by recursively applying the 3 operation on those 3 constants

#### Example Regular Expressions in Ruby

- /Ruby/
  - Matches exactly the string "Ruby"
  - Concatenation: /r u b y/
- ▶ /Ruby | OCaml/
  - Matches either "Ruby" or "OCaml"

- /(ab)\*/
  - 0 or more occurrences of "ab", matches "", "ab", "abab", "ababab", ...

# Using Regular Expressions

Regular expressions are instances of Regexp

Basic matching using =~ method of String

#### Repetition in Regular Expressions

- \*: zero or more
- +: one or more
  - so /e+/ is the same as /ee\*/
- ?: zero or one occurrence
- {x} means repeat the search for exactly x occurrences
- {x,} means repeat the search for at least x occurrences
- {x, y} means repeat the search for at least x occurrences and at most y occurrences

#### Watch Out for Precedence

/(Ruby)\*/ means {"", "Ruby", "RubyRuby", ...}

- /Ruby\*/ means {"Rub", "Ruby", "Rubyy", ...}
- Best to use parentheses to disambiguate
  - Note that parentheses have another use, to extract matches, as we'll see later

#### **Character Classes**

- /[abcd]/
  - {"a", "b", "c", "d"} (Can you write this another way?)
- ► /[a-zA-Z0-9]/
  - Any upper or lower case letter or digit
- /[^0-9]/
  - Any character except 0-9 (the ^ is like not and must come first)
- /[\t\n]/
  - Tab, newline or space
- /[a-zA-Z\_\\$][a-zA-Z\_\\$0-9]\*/
  - Java identifiers (\$ escaped...see next slide)

#### **Special Characters**

```
any character
                                            Using /^pattern$/
             beginning of line
Λ
                                            ensures entire
             end of line
                                            string/line must
                                            match pattern
\$
             just a $
             digit, [0-9]
\d
             whitespace, [\t\r\n\f\s]
\s
             word character, [A-Za-z0-9]
\w
\D
             non-digit, [^0-9]
\S
             non-space, [^\t\r\n\f\s]
             non-word, [^A-Za-z0-9_]
\W
```

#### Potential Character Class Confusions

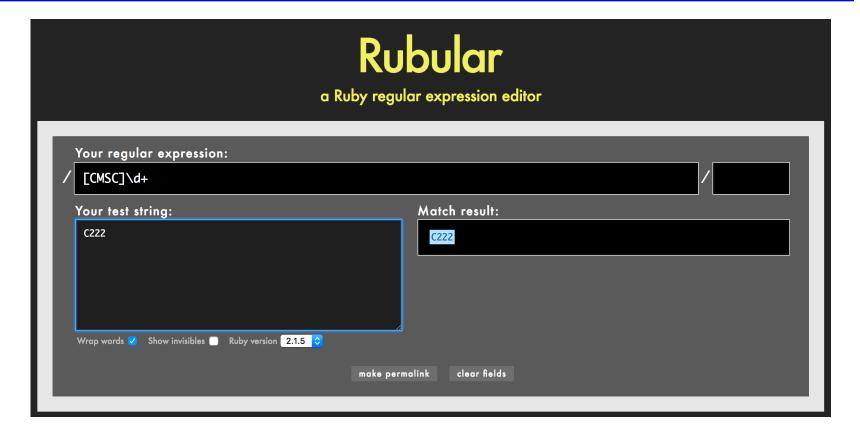
Inside character classes: not Outside character classes: beginning of line Inside regular expressions: character class Outside regular expressions: array Note: [a-z] does not make a valid array Inside character classes: literal characters ( )  $\rightarrow$  Note /(0..2)/ does not mean 012 Outside character classes: used for grouping Inside character classes: range (e.g., a to z given by [a-z])

Outside character classes: subtraction

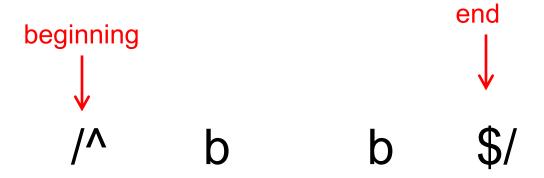
#### Summary

- ▶ Let *re* represents an arbitrary pattern; then:
  - /re/ matches regexp re
  - $/(re_1|re_2)/$  match either  $re_1$  or  $re_2$
  - /(re)\*/ match 0 or more occurrences of re
  - /(re)+/ match 1 or more occurrences of re
  - /(re)?/ match 0 or 1 occurrences of re
  - /(re){2}/ match exactly two occurrences of re
  - /[a-z]/ same as (a|b|c|...|z)
  - / [^0-9]/ match any character that is not 0, 1, etc.
  - ^, \$ match start or end of string

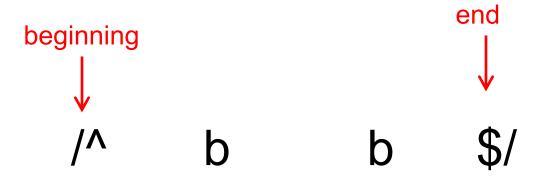
#### Try out regexps at rubular.com



Contains 2 b's, may not be consecutive.



Contains 2 b's, may not be consecutive.



Starts with c, followed by one vowel, and any number of letters

/^C

\$/

Starts with c, followed by one vowel, and any number of letters

/^c

\$/

Starts with c, followed by one vowel, and any number of letters

All letters are in alphabetic order

All letters are in alphabetic order

/^a\*b\*c\*d\*e\*f\*g\*h\*i\*j\*k\*I\*m\*n\*o\*p\*r\*t\*\$/

Contains sss or ccc

Contains sss or ccc

contains 2 ab

$$/(ab){2}/$$

contains 2 b

Starts with a, 0 or 1 letter after that

Starts with a, 0 or 1 letter after that

Contains one or more ab or ba

Contains one or more ab or ba

steve, steven, or stephen

steve, steven, or stephen

/^ste(ve|phen|ven)\$/

Even length string

Even number of vowels

Even number of vowels

```
/^([^aouei]*[aouei][^aouei]*[aouei][^aouei]*)*$/
```

starts with not-b, followed by one or more a's

starts with not-b, followed by one or more a's

#### Quiz 1

How many different strings could this regex match?

```
/^Hello, Anyone awake?$/
```

- A. 1
- в. **2**
- c. 4
- D. More than 4

How many different strings could this regex match?

e or nothing

```
/^Hello, Anyone awake?$/
```

- A. 1
- в. 2
- c. 4
- D. More than 4

Which regex is **not** equivalent to the others?

```
A. ^[crab]$
B. ^c?r?a?b?$
c. ^(c|r|a|b)$
D. ^([cr]|[ab])$
```

Which regex is **not** equivalent to the others?

```
A. ^[crab]$
B. ^c?r?a?b?$
c. ^(c|r|a|b)$
d. ^([cr]|[ab])$
```

Which string does not match the regex?

$$/[a-z]{4}\d{3}/$$

- A. "cmsc\d\d\d"
- B. "cmsc330"
- c. "hellocmsc330"
- D. "cmsc330world"

# Which string does not match the regex?

Recall that without ^ and \$, a regex will match any **sub**string

$$/[a-z]{4}\d{3}/$$

- A. "cmsc\d\d\d"
- B. "cmsc330"
- c. "hellocmsc330"
- D. "cmsc330world"

# Extracting Substrings based on R.E.'s Method 1: Back References

Two options to extract substrings based on R.E.'s:

- Use back references
  - Ruby remembers which strings matched the parenthesized parts of r.e.'s
  - These parts can be referred to using special variables called back references (named \$1, \$2,...)

# Back Reference Example

sets min = \$1 and max = \$2

#### Input

Min: 1 Max: 27

Min: 10 Max: 30

Min: 11 Max: 30

Min: a Max: 24

#### Output

mini=1 maxi=27

mini=10 maxi=30

mini= maxi=

mini= maxi=

Extra space messes up match

Not a digit; messes up match

#### **Back References are Local**

#### Warning

- Despite their names, \$1 etc are local variables
- (Normally, variables starting with \$ are global)

```
def m(s)
    s =~ /(Foo)/
    puts $1  # prints Foo
    end
    m("Foo")
    puts $1  # prints nil
```

#### Back References are Reset

#### Warning 2

If another search is performed, all back references are reset to nil

```
gets =~ /(h)e(ll)o/
puts $1
puts $2
gets =~ /h(e)llo/
puts $1
puts $2
gets =~ /hello/
puts $1
```

```
hello
h
ll
hello
e
nil
hello
nil
```

47

CMSC 330 - Fall 2020

What is the output of the following code?

```
s = "help I'm stuck in a text editor"
s =~ /([A-Z]+)/
puts $1
```

- A. help
- B.
- c. **l'm**
- D. I'm stuck in a text editor

What is the output of the following code?

```
s = "help I'm stuck in a text editor"
s =~ /([A-Z]+)/
puts $1
```

- A. help
- B.
- c. I'm
- D. I'm stuck in a text editor

What is the output of the following code?

```
"Why was 6 afraid of 7?" = \sim /\d\s(\w+).*(\d)/puts $2
```

- A. afraid
- B. Why
- c. 6
- D. 7

What is the output of the following code?

```
"Why was 6 afraid of 7?" = \sim /\d\s(\w+).*(\d)/puts $2
```

- A. afraid
- B. Why
- c. **6**
- D. **7**

# Method 2: String.scan

- Also extracts substrings based on regular expressions
- Can optionally use parentheses in regular expression to affect how the extraction is done
- Has two forms that differ in what Ruby does with the matched substrings
  - The first form returns an array
  - The second form uses a code block
    - > We'll see this later

#### First Form of the Scan Method

- str.scan(regexp)
  - If regexp doesn't contain any parenthesized subparts, returns an array of matches
    - > An array of all the substrings of *str* which matched

```
s = "CMSC 330 Fall 2018"
s.scan(/\S+ \S+/)
# returns array ["CMSC 330", "Fall 2018"]
```

```
s.scan(/\S{2}/)
# => ["CM", "SC", "33", "Fa", "11", "20", "18"]
```

# First Form of the Scan Method (cont.)

- If regexp contains parenthesized subparts, returns an array of arrays
  - Each sub-array contains the parts of the string which matched one occurrence of the search

```
s = "CMSC 330 Fall 2018"
s.scan(/(\S+) (\S+)/) # [["CMSC", "330"],
# ["Fall", "2018"]]
```

- Each sub-array has the same number of entries as the number of parenthesized subparts
- > All strings that matched the first part of the search (or \$1 in back-reference terms) are located in the first position of each sub-array

#### Practice with Scan and Back-references

```
> 1s -1
drwx----
            2 sorelle sorelle
                                      4096 Feb 18 18:05 bin
-rw----
            1 sorelle sorelle
                                      674 Jun 1 15:27 calendar
             3 sorelle sorelle
drwx----
                                      4096 May 11 2006 cmsc311
            2 sorelle sorelle
drwx----
                                      4096 Jun 4 17:31 cmsc330
drwx----
         1 sorelle sorelle
                                      4096 May 30 19:19 cmsc630
drwx----
              1 sorelle sorelle
                                      4096 May 30 19:20 cmsc631
```

#### Extract just the file or directory name from a line using

```
• SCan name = line.scan(/\S+$/) # ["bin"]
```

back-references

What is the output of the following code?

```
s = "Hello World"
t = s.scan(/\w{2}/).length
puts t
```

- A. 3
- в. **4**
- c. 5
- D. 6

What is the output of the following code?

```
s = "Hello World"
t = s.scan(/\w{2}/).length
puts t
```

- A. 3
- в. 4
- c. 5
- d. **6**

What is the output of the following code?

```
s = "To be, or not to be!"
a = s.scan(/(\S+) (\S+)/)
puts a.inspect
```

```
A. ["To","be,","or","not","to","be!"]
B. [["To","be,"],["or","not"],["to","be!"]]
c. ["To","be,"]
D. ["to","be!"]
```

What is the output of the following code?

```
s = "To be, or not to be!"
a = s.scan(/(\S+) (\S+)/)
puts a.inspect
```

```
A. ["To","be,","or","not","to","be!"]
B. [["To","be,"],["or","not"],["to","be!"]]
c. ["To","be,"]
D. ["to","be!"]
```

#### Second Form of the Scan Method

Can take a code block as an optional argument

- str.scan(regexp) { |match| block }
  - Applies the code block to each match
  - Short for str.scan(regexp).each { |match| block }
  - The regular expression can also contain parenthesized subparts

# Example of Second Form of Scan

```
34
12
        23
                 input file:
19 77 87
                 will be read line by line, but
11 98
                 column summation is desired
   45
sum a = sum b = sum c = 0
while (line = gets)
  line.scan(/(d+)s+(d+)/s+(d+)/) { |a,b,c|
    sum a += a.to i
                           converts the string
    sum b += b.to i \leftarrow
                           to an integer
    sum c += c.to i
```

printf("Total: %d %d %d\n", sum a, sum b, sum c)

#### Sums up three columns of numbers

end

# Practice: Amino Acid counting in DNA

Write a function that will take a filename and read through that file counting the number of times each group of three letters appears so these numbers can be accessed from a hash.

(assume: the number of chars per line is a multiple of 3)

gcggcattcagcacccgtatactgttaagcaatccagatttttgtgtataacataccggc catactgaagcattcattgaggctagcgctgataacagtagcgctaacaatgggggaatg tggcaatacggtgcgattactaagagccgggaccacacaccccgtaaggatggagcgtgg taacataataatccgttcaagcagtgggcgaaggtggagatgttccagtaagaatagtgg gggcctactacccatggtacataattaagagatcgtcaatcttgagacggtcaatggtac cgagactatatcactcaactccggacgtatgcgcttactggtcacctcgttactgacgga

# Practice: Amino Acid counting in DNA

```
get the
          def countaa(filename)
file
          file = File.new(filename, "r")
handle
             lines = file.readlines
                                                 initialize
                                                 the hash, or
array
             hash = Hash.new
                                                 you will get
of
             lines.each{ |line|
lines
                                                 an error when
                   acids = line.scan(/.../)
from
                                                 trying to
                   acids.each{ |aa|
                                                 index into an
the
                      if hash[aa] == nil
£9¥e
                                                 array with a
                            hash[aa] = 1
<del>each</del>
                                                 string
line in
                      else
the
                            hash[aa] += 1
                                                 get an array
fale
                                                 of triplets
                      end
each
                                                 in the line
triplet
in the
          end
line
```

CMSC 330 - Fall 2020

63