

# Welcome!

- **today**
  - administrative stuff
  - course overview
  - regular expressions and grep
- **check out the course web page (CS, not Blackboard) and Piazza**
  - notes, readings and assignments are posted only on the web page; monitor the web page and Piazza every day
  - Assignment 1 is posted; due 10:00 pm Thursday Feb 14
  - initial project information is posted (more on Thursday)
- **please do the survey if you haven't already**

# People



**Christopher  
Moretti**



**Allison  
Chang '18**



**Jace  
Lu**



**Lance  
Goodridge '17**

## Very Tentative Outline

week 1	regular expressions, grep; project info
week 2	scripting: shell, AWK, Python
week 3	web technology: HTTP, CSS, Javascript
week 4	client libraries and server frameworks
week 5	user interfaces; phone apps
week 6	<b>databases</b> ; software engineering
(spring break)	
week 7	networks
week 8	advanced C++, Java; Go
week 9	APIs, design patterns
week 10	XML, JSON, REST, DSLs
week 11	??
week 12	??
(start of reading period)	
May 8-9?	project presentations
May 12	project submission (2 days before Dean's Date)

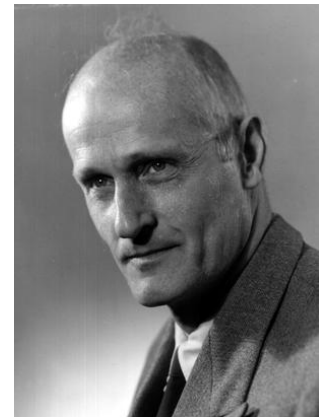
# House rules

- please turn cell phones off
- please don't use your laptop, tablet, phone, ...
  - it distracts you
  - it distracts your neighbors
  - it distracts me
- please don't snore (sleeping is ok)
- please sit towards the front, not in the back
- please stay away if you're sick !!!
- please ask questions about anything at any time

# Regular expressions and grep

- **regular expressions**

- based on ideas from automata theory pioneered by Stephen Kleene \*34
- notation
- mechanization
- pervasive in Unix tools
- in all scripting languages, often as part of the syntax
- in general-purpose languages, as libraries
- basic implementation is remarkably simple
- efficient implementation requires good theory and good practice



- **grep is the prototypical tool**

- written by Ken Thompson @ Bell Labs ~1972



# Grep regular expressions

- c** any character matches itself, except for *metacharacters* `. [ ] ^ $ * \`
- $r_1 r_2$**  matches  $r_1$  followed by  $r_2$
- .** matches any single character
- [ . . . ]** matches one of the characters in set ...  
shorthand like a-z or 0-9 includes any character in the range
- [ ^ . . . ]** matches one of the characters not in set  
[<sup>^</sup>0-9] matches non-digit
- ^** matches beginning of line when ^ begins pattern  
no special meaning elsewhere in pattern
- \$** matches end of line when \$ ends pattern  
no special meaning elsewhere in pattern
- \*** any regular expression followed by \* matches 0 or more
- \c** matches c unless c is ( ) or digit
- \ ( . . . ) \** tagged regular expression that matches ...  
the matched strings are available as \1, \2, etc.

## Examples of matching

<code>xy</code>	<code>xy</code> anywhere in string
<code>^xy</code>	<code>xy</code> at beginning of string
<code>xy\$</code>	<code>xy</code> at end of string
<code>^xy\$</code>	string that contains only <code>xy</code>
<code>^</code>	matches any string, even empty
<code>^\$</code>	empty string
<code>.</code>	non-empty string, i.e., at least 1 char
<code>xy.\$</code>	<code>xy</code> plus any char at end of string
<code>xy\.\$</code>	<code>xy.</code> at end of string
<code>\\xy\\</code>	<code>xy</code> anywhere in string
<code>[xX]y</code>	<code>xy</code> or <code>Xy</code> anywhere in string
<code>xy[0-9]</code>	<code>xy</code> followed by one digit
<code>xy[^0-9]</code>	<code>xy</code> followed by a non-digit
<code>xy[0-9][^0-9]</code>	<code>xy</code> followed by digit, then non-digit
<code>xy1.*xy2</code>	<code>xy1</code> then any text then <code>xy2</code>
<code>^xy1.*xy2\$</code>	<code>xy1</code> at beginning and <code>xy2</code> at end

# "Regular expressions" are not always regular

- there's a precise definition but lots of casual usage
- **R:**  $c$   $R_1R_2$   $R_1 | R_2$   $(R)$   $R^*$ 
  - equivalent to a finite automaton
  - this is what egrep provides
- shorthands like `[A-Z]`, `\d`, `[:alnum:]`, etc., don't change properties
  - can't count, can't recognize repeated strings, ...
- can have subsets that do less (coming up)
- can do much more than pure REs:
  - supersets (back-referencing in grep)
  - libraries that are Turing-complete (Java, Python, etc.)
  - extra-lingual processing (commandline arguments like `grep -i -v --color`)



# egrep: fancier regular expressions

$r^+$  one or more occurrences of  $r$   
 $r^?$  zero or one occurrences of  $r$   
 $r_1 | r_2$   $r_1$  or  $r_2$   
 $(r)$   $r$  (grouping)

grammar:

$r$ :  $c$   $.$   $^$   $\$$   $[ccc]$   $[\^ccc]$   
 $r^*$   $r^+$   $r^?$   
 $r_1 r_2$   
 $r_1 | r_2$   
 $(r)$

precedence:

$*$   $+$   $?$  higher than concatenation, which is higher than  $|$

$([0-9]^+\backslash.?[0-9]^*|\backslash.[0-9]^+)([Ee][^-+]?[0-9]^+)?$



Al Aho \*67

# Python RE's

<code>^</code>	Matches beginning of line.
<code>\$</code>	Matches end of line.
<code>.</code>	Matches any single character except newline. Using <code>m</code> option allows it to match newline as well.
<code>[...]</code>	Matches any single character in brackets.
<code>[^...]</code>	Matches any single character not in brackets
<code>re*</code>	Matches 0 or more occurrences of preceding expression.
<code>re+</code>	Matches 1 or more occurrence of preceding expression.
<code>re?</code>	Matches 0 or 1 occurrence of preceding expression.
<code>re{n}</code>	Matches exactly <code>n</code> number of occurrences of preceding expression.
<code>re{n,}</code>	Matches <code>n</code> or more occurrences of preceding expression.
<code>re{n, m}</code>	Matches at least <code>n</code> and at most <code>m</code> occurrences of preceding expression.
<code>alb</code>	Matches either <code>a</code> or <code>b</code> .
<code>(re)</code>	Groups regular expressions and remembers matched text.
<code>(?imx)</code>	Temporarily toggles on <code>i</code> , <code>m</code> , or <code>x</code> options within a regular expression. If in parentheses, only that area is affected.
<code>(?-imx)</code>	Temporarily toggles off <code>i</code> , <code>m</code> , or <code>x</code> options within a regular expression. If in parentheses, only that area is affected.
<code>(?: re)</code>	Groups regular expressions without remembering matched text.
<code>(?imx: re)</code>	Temporarily toggles on <code>i</code> , <code>m</code> , or <code>x</code> options within parentheses.
<code>(?-imx: re)</code>	Temporarily toggles off <code>i</code> , <code>m</code> , or <code>x</code> options within parentheses.
<code>(?#...)</code>	Comment.
<code>(?= re)</code>	Specifies position using a pattern. Doesn't have a range.
<code>(?! re)</code>	Specifies position using pattern negation. Doesn't have a range.
<code>(?&gt; re)</code>	Matches independent pattern without backtracking.
<code>\w</code>	Matches word characters.
<code>\W</code>	Matches nonword characters.
<code>\s</code>	Matches whitespace. Equivalent to <code>[\t\n\r\f]</code> .
<code>\S</code>	Matches nonwhitespace.
<code>\d</code>	Matches digits. Equivalent to <code>[0-9]</code> .
<code>\D</code>	Matches nondigits.
<code>\A</code>	Matches beginning of string.
<code>\Z</code>	Matches end of string. If a newline exists, it matches just before newline.
<code>\z</code>	Matches end of string.
<code>\G</code>	Matches point where last match finished.
<code>\b</code>	Matches word boundaries when outside brackets. Matches backspace (0x08) when inside brackets.
<code>\B</code>	Matches nonword boundaries.
<code>\n, \t, etc.</code>	Matches newlines, carriage returns, tabs, etc.
<code>\1...\9</code>	Matches <code>n</code> th grouped subexpression.
<code>\10</code>	Matches <code>n</code> th grouped subexpression if it matched already. Otherwise refers to the octal representation of a character code.

# The grep family

- **grep**
- **egrep**
  - fancier regular expressions, trades compile time and space for run time
- **fgrep**
  - parallel search for many fixed strings
- **agrep**
  - "approximate" grep: search with errors permitted
- **relatives that use similar regular expressions**
  - ed                                    original Unix editor
  - sed                                    stream editor
  - vi, emacs, sam, ...            editors
  - lex, flex                            lexical analyzer generator
  - awk, perl, python, ...        all scripting languages
  - Java, C# ...                        libraries in mainstream languages
- **simpler variants**
  - filename "wild cards" in Unix and other shells (assignment 1 this year)
  - "LIKE" operator in SQL, Visual Basic, etc.

# Important ideas from regexps & grep

- **tools: let the machine do the work**
  - good packaging matters
- **notation: makes it easy to say what to do**
  - may organize or define implementation
- **hacking can make a program faster, sometimes, usually at the price of more complexity**
- **a better algorithm can make a program go a lot faster**
- **don't worry about performance if it doesn't matter (and it often doesn't)**
- **when it does,**
  - use the right algorithm
  - use the compiler's optimization
  - code tune, as a last resort