*Advanced Programming Techniques*
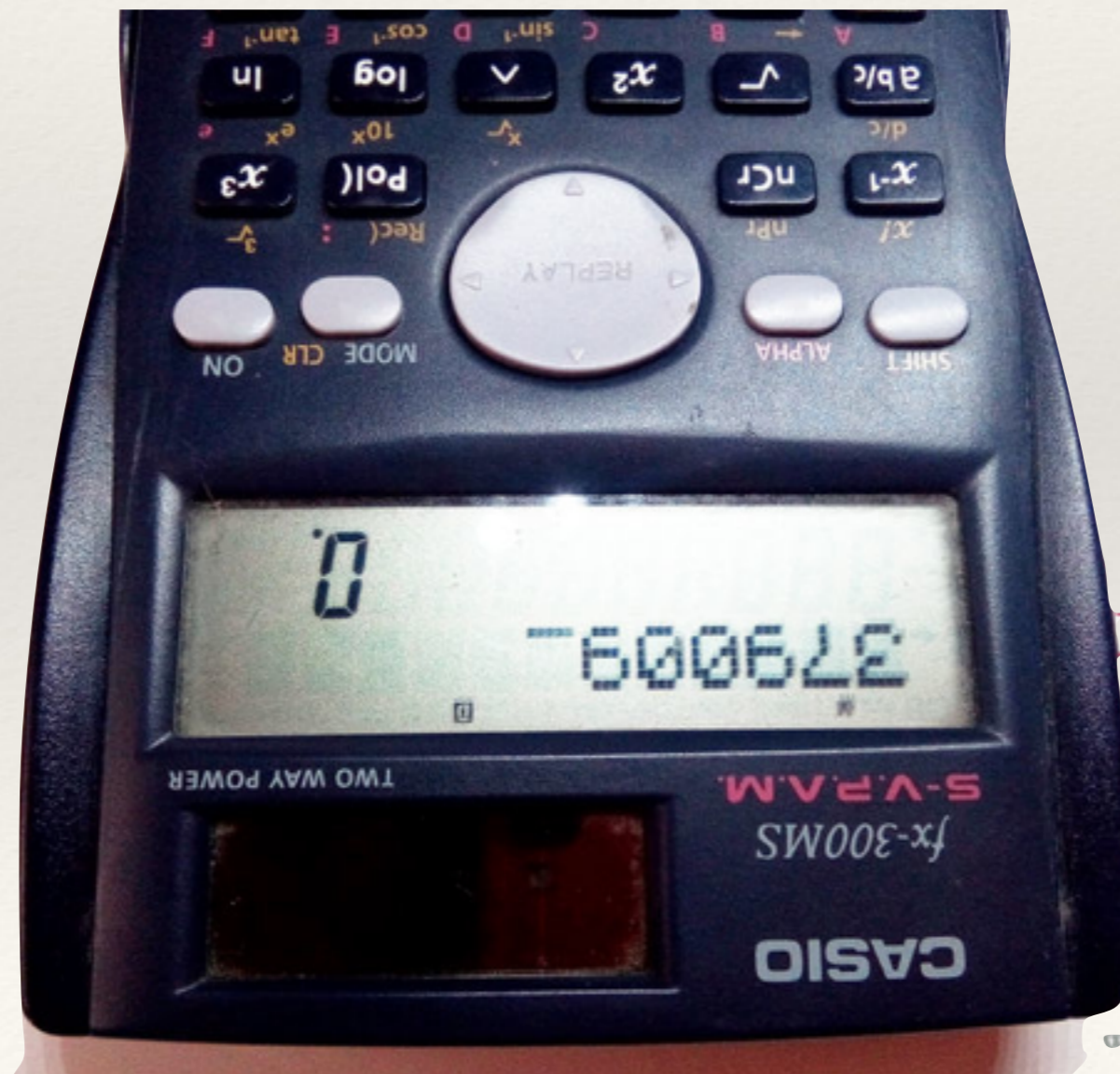
# Regular Expressions

Christopher Moretti

"It's not a silly question if you can't answer it."

–Jostein Gaarder, *Sophie's World*

# Sounds like 1134? Not in this font.

❖ Based on ideas from automata theory in regular languages pioneered by Stephen Kleene *34

❖ Practically applied in late-60's in various settings (compilers, editors) and started to enter into languages in the 70's (e.g., AWK)

❖ Mechanized text filtering and pattern match

  ❖ Pervasive in *nix tools (e.g. sed, grep)

  ❖ Key built-in for scripting syntax (pattern-action languages)

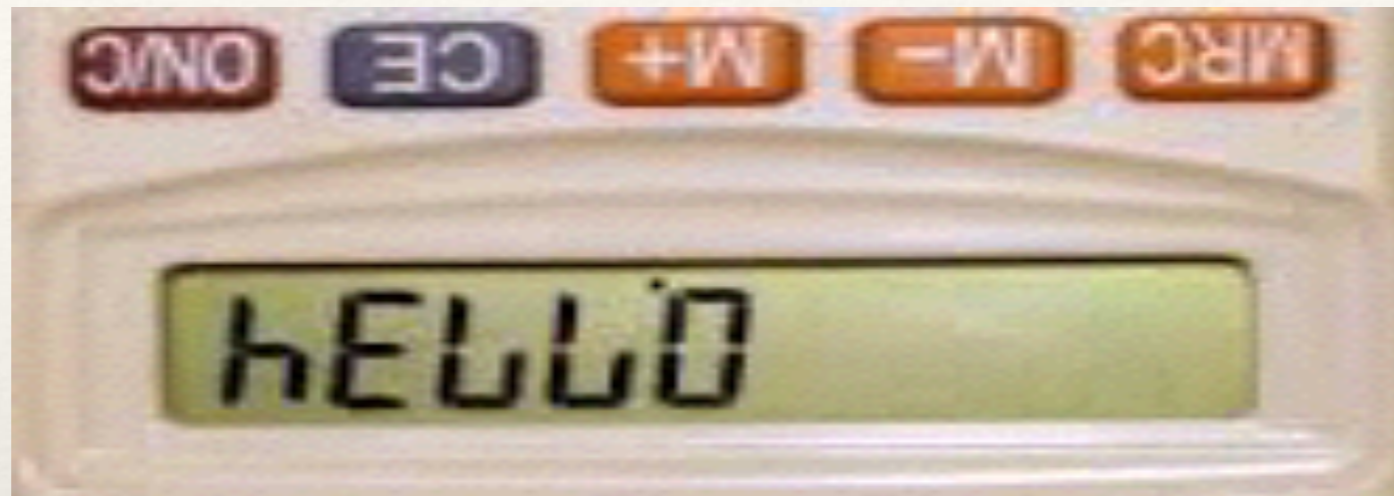  ❖ Available in just about any other language as a library

# grep Regular Expressions

`c`      Any character matches itself, except for meta-characters: `. [ ] ^ $ * \`

`rs`      Matches r followed by s

`.`      Matches any single character

`[rst]`      Matches one of r, s, or t    — range shorthands are allowed: [0-9a-z]

`[^rst]`      Matches one character other than r, s, or t

`^`      Matches start of the line when at start of pattern (not special otherwise)

`$`      Matches end of the line when at end of pattern (not special otherwise)

`x*`      Kleene closure: 0 or more repetitions of x

`\c`      Matches c unless c is ( ) or a digit — \ is the escape character

`\(x\)`      Matches the pattern x and saves the substring that matches (x') as \#

`x{m,n}`      Matches x repeated between m and n times

# Reasonable Toy Examples

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

# RE-Building our Vocabulary



1. Contains only the letters beghilosz

2. Can use the letters in any order

3. Can use any letter multiple times

4. Assume no length restriction

```
^[beghilosz]*$
```

# grep in Action

```
tars: wc -l /usr/share/dict/words
  235886 /usr/share/dict/words

tars: grep '^[beghilosz]*$' words | wc -l
    451



opus: wc -l /usr/share/dict/words
479828 /usr/share/dict/words

opus: grep '^[beghilosz]*$' words | wc -l
975
```

Note the single-quotes around the REs —
we don't want the shell to  take matters into its own hands

opus: grep '^[beghilosz]*$' words | column -c 220

```
b             biggies     bogglebo   egises      ghis        gobies      hello        hoo       leo       logoes       oilish     shi        slobs
bb            biggish     boggles    ego         ghole       gobo        helloes      hool      leos      logoi        oilless    shiel      sloe
bbl           bilbi       bogglish   egoize      gi          goboes      hellos       hoolee    les       logos        oils       shiels     sloes
bbls          bilbie      boghole    egol        gib         gobos       hells        hoolie    lese      logs         ol         shies      slog
bbs           bilbies     bogie      egoless     gibbles     gobs        heloe        hoose     less      loli         ole        shih       slogs
be            bilbo       bogies     egos        gibbol      goebbels    helosis      hoosh     lessee    loll         oleo       shill      sloo
bebless       bilboes     bogle      eh          gibbose     goel        helzel       hose      lessees   lollies      oleos      shilloo    sloosh
bebog         bilbos      bogles     eigh        gibe        goes        heo          hosel     lesses    lolls        oleose     shills     slosh
beboss        bile        bogo       eisegeses   gibel       gog         hes          hoseless  lez       loo          oles       shish      sloshes
bee           biles       bogs       eisegesis   gibes       goggle      hg           hosels    lezes     loob         olio       sho        so
beebee        bilge       boh        eisell      gibleh      goggles     hi           hoses     lezzie    looie        ologies    shoe       sob
beebees       bilges      boho       el          gibli       gogo        hie          hoss      lezzies   looies       ooh        shoebill   sobole
beeish        bilio       boil       elb         gibs        gogos       hies         hs        lg        loos         oohs       shoebills  soboles
beele         bill        boils      elegies     gie         goi         higgle       i         lh        looses       oollies    shoeless   sobs
beelol        billie      boise      elegise     gies        goil        higgles      ib        lhb       loosish      oologies   shoes      soe
bees          billies     bol        elegises    giesel      gol         high         ibis      li        lose         oologize   shog       sog
beg           bills       bole       elegize     gig         golee       highhole     ibisbill  lib       losel        oooo       shoggie    soh
begiggle      bilo        boles      elegizes    gigge       goles       highish      ibises    libel     losels       oos        shoggle    soho
begloze       bilobe      bolis      elhi        gighe       goli        highs        ie        libelees  loses        oose       shogi      soil
bego          bilos       boll       eli         giggish     goll        hile         ieee      libellee  losh         ooze       shogs      soils
begob         bilsh       bolls      eligible    giggle      goloe       hili         igloo     libellees loss         oozes      shole      soilless
begobs        bio         bollies    eligibles   giggles     golosh      hill         igloos    libels    losses       os         shoo       sol
begs          biog        bolo       ell         gigglish    goloshe     hillbillies  ihi       lie       lossless     ose        shoogle    sole
beige         biol        bolos      elle        giglio      goloshes    hillos       ihs       liebig    o            osi        shooi      solei
beigel        biologese   bolshie    ells        gigolo      goo         hills        ii        liege     ob           oslo       shool      soleil
beiges        biologies   bolshies   eloge       gigolos     goog        his          iii       liegeless obb          osse       shools     soleless
bel           biologize   boo        els         gigs        googlies    hish         il        lieges    obe          oz         shoos      solgel
belee         bios        boob       else        gile        googol      hisis        ile       lies      obeish       ozs        si         soli
belibel       biose       boobie     elses       gilgie      googols     hiss         ill       liesh     obeli        s          sib        solio
belie         biosis      boobies    eo          gill        gool        hissel       illess    lig       obelise      sb         sibb       solo
belies        bis         boobish    eole        gillie      gools       hisses       illish    lige      obelises     se         sibbs      solos
bell          bise        booboisie  eos         gills       goos        hizz         ills      ligge     obelize      see        sibs       sols
belle         bises       booboo     es          gilo        goose       hizzie       io        lile      obelizes     seel       sie        soogee
belles        bish        booboos    ese         gils        goosebill   hl           ios       liles     obes         seels      siege      sool
belli         bisie       boobs      eses        gilse       goosegog    ho           is        lilies    obese        sees       sieges     sooloos
bellies       bisso       boogie     ess         gio         gooses      hob          ise       lill      obi          seesee     sig        soosoo
bello         biz         boogies    esse        gis         goosish     hobbies      ish       lills     obis         seg        sigh       sos
bells         bize        boohoo     essee       gise        gos         hobbil       isis      lis       obl          seggio     sighless   sosh
bels          bizel       boohoos    esses       gisel       gosh        hobble       isize     lish      obli         seghol     sighs      sosie
bes           bizes       bool       g           gish        goss        hobbles      isl       lisle     oblige       sego       sigil      soso
besee         bizz        boolies    ge          gizz        gozell      hoblob       isle      lisles    obligee      segol      sigill     sosoish
beseige       bl          boos       geb         gl          gozill      hobo         isleless  liss      obligees     segos      sigils     soss
beshell       bleb        boose      gebbie      glb         h           hoboe        isles     lisses    obliges      segs       sigloi     sossiego
besiege       blebs       boosies    gee         glebe       hb          hoboes       isls      lizzie    obol         sei        siglos     sossle
besieges      blee        booze      gees        glebeless   he          hobos        iso       ll        oboe         seige      sil        sozzle
besigh        bleeze      boozes     geez        glebes      hebe        hobs         isogloss  llb       oboes        seis       sile       z
besoil        bleo        bos        gegg        glee        hebes       hoe          isoglosses lo       oboles       seise      sill       zee
bess          bless       bose       geggee      gleg        hee         hoes         isohel    lob       oboli        seises     sillibib   zees
bessel        blesse      bosh       gel         gliosis     heh         hog          isohels   lobe      obolos       seize      sillibibs  zeiss
besses        blesses     boshes     gelee       gliss       hehs        hogg         isolog    lobeless  obols        seizes     sillies    zel
bezel         blibe       boss       gelees      glob        heel        hoggee       isologs   lobes     obs          sel        silo       zho
bezels        bliss       bosses     gell        globe       heelless    hoggie       isz       lobi      obsess       sele       silos      zig
bezil         blisses     bossies    gelose      globes      heels       hoggish      izle      loblollies obsesses    sell       sis        zigs
bezils        blissless   bozo       gelosie     globose     heeze       hoggs        l         lobo      oe           selle      sise       zill
bezzi         blizz       bozos      gels        globs       heezes      hogo         lb        lobolo    oes          sellie     sisel      zills
bezzle        blo         bozze      geo         glogg       heezie      hogs         lbs       lobolos   oesogi       sells      sises      zizel
bezzo         blob        bs         geobios     gloggs      hei         hoh          le        lobos     og           sels       sish       zizz
bg            blobs       bsh        geog        glos        heigh       hoho         lebes     lobose    ogee         sess       sisi       zizzle
bi            blooie      bz         geol        gloss       heii        hohs         lee       lobs      ogees        sg         siss       zizzles
bib           bls         e          geologies   glosses     heil        hoi          lees      loe       ogle         sh         sissies    zo
bibb          bo          ebb        geologise   glossies    heils       hoise        leese     loeil     ogles        she        sissoo     zobo
bibble        bob         ebbs       geologize   glossless   heishi      hoises       leg       loess     oh           shee       size       zogo
bibbs         bobbie      eblis      gess        glossologies heize      hol          lege      loesses   ohelo        sheel      sizes      zoll
bibi          bobbies     eboe       gesso       gloze       hel         hole         leges     loge      ohio         sheol      sizz       zolle
bibl          bobbish     ee         gessoes     glozes      helbeh      holeless     legge     loges     oho          sheols     sl         zoo
bible         bobble      eel        gez         go          hele        holes        legible   loggie    ohs          shes       slee       zoogeog
bibles        bobbles     eelbob     ghee        gob         helio       holi         legis     loggish   oie          shh        sleigh     zool
bibless       bobo        eels       ghees       gobble      helios      holies       legless   loglog    oii                     sleighs    zoologies
bibliog       bobol       eg         ghi         gobbles     heliosis    holl         legs      logie     oil                     slish      zoologize
bibliologies  bobs        egg        ghibli      gobi        hell        holloes      lei       logis     oilhole                 slob       zoosis
biblos        bobsleigh   eggless    ghiblis                 hellhole    holloo       leis      logo      oilholes                slobbish   zoozoo
bibs          boe         eggs       ghillie                 hellholes   holloos      leiss                                                  zs
big           bog         eggshell   ghillies                hellish     hollos
biggie        boggle      egis                                           hols
```

# Back-references

`\(rst\)` Matches rst and saves that matching string as \1 to be used later in pattern

`\(abc\)\(rst\)\1\(xyz\)\3\2\1`

```
echo abcrstabcxyzxyzrstabc | grep '\(abc\)\(rst\)\1\(xyz\)\3\2\1'
abcrstabcxyzxyzrstabc
```

`\(ious\)\1`

```
grep '\(ious\)\1' words
homoiousious
```

`\(ious\).*\1`

```
grep '\(ious\).*\1' words
homoiousious
```

`\(.\).*\1.*\1.*\1.*\1.*\1.*\1.*\1`

```
grep '\(.\).*\1.*\1.*\1.*\1.*\1.*\1.*\1' words
bras-dessus-bras-dessous
humuhumunukunukuapuaa
pneumonoultramicroscopicsilicovolcanoconiosis
possessionlessness
```

# Nested Back-references

`\(r\(.\)\2\).*\1`

So we have an r, then the same character twice,
then any string,
then the r and **same** same character twice

```
grep '\(r\(.\)\2\).*\1' words
broomroot
free-reed
greegree
greegrees
Greentree
groo-groo
proof-proof
proofroom
reel-to-reel
six-three-three
three-reel
tree-creeper
```

# POSIX Extended Standard

`(x)`            Grouping

`x+`             Matches x repeated 1 or more times

`x?`             Matches x 0 or 1 times

`a|b`            Matches a OR b

`[:alnum:]`      Alphanumeric characters: [A-Za-z0-9]

`[:alpha:]`      Alphabetic characters: [A-Za-z]

`[:blank:]`      Space and Tab

`[:digit:]` `[:lower:]` `[:upper:]` `[:xdigit:]`

# egrep

egrep (equivalent to 'grep -E') is "extended" grep
Some changes/non-standards/caveats:

```
echo foo | egrep '[:digit:]'
grep: character class syntax is [[:space:]], not [:space:]
echo foo | egrep '[[:digit:]]'
echo f00 | egrep '[[:digit:]]'
f00

echo tartar | grep '(a..a)*'
echo 't(arta)r' | grep '(a..a)*'
t(arta)r
echo tartar | egrep '(a..a)*'
tartar

echo tartar | grep '\(t..\)\1'
tartar
echo tartar | egrep '\(t..\)\1'
grep: Invalid back reference
echo tartar | egrep '(t..)\1'
tartar
```

# Party Tricks

## Words with all the vowels in order?

```
grep 'a.*e.*i.*o.*u.*y' words
```

**abstemiously**
**adventitiously**
**anticensoriously**
**anticeremoniously**
**antireligiously**
**auteciously**
**autoeciously**
f**acetiously**
h**alf-ingeniously**
h**alf-rebelliously**
h**alf-seriously**
non**abstemiously**
non**adventitiously**
non**facetiously**
non**sacrilegiously**
over**abstemiously**
p**ancreaticoduodenostomy**
p**areciously**
p**aroeciously**
pseudos**acrilegiously**
qu**asi-conscientiously**
qu**asi-enviously**
qu**asi-mysteriously**
qu**asi-rebelliously**
qu**asi-religiously**
qu**asi-seriously**
s**acrilegiously**
un**abstemiously**
un**facetiously**
uns**acrilegiously**

## No, what I really meant was …

```
grep 'a[^aeiouy]*e[^aeiouy]*i[^aeiouy]*o[^aeiouy]*u[^aeiouy]*y' words
```

**abstemiously**
f**acetiously**
h**alf-seriously**
non**abstemiously**
non**facetiously**
over**abstemiously**
p**areciously**
un**abstemiously**
un**facetiously**

## No, what I really, *really* meant was …

```
grep '^[^aeiouy]*a[^aeiouy]*e[^aeiouy]*i[^aeiouy]*o[^aeiouy]*u[^aeiouy]*y[^aeiouy]*$' words
```

**abstemiously**
**facetiously**
**half-seriously**
**pareciously**

# Words with all their letters in order?

```
grep '^a*b*c*d*e*f*g*h*i*j*k*l*m*n*o*p*q*r*s*t*u*v*w*x*y*z*$' words | wc -l
1370
grep $pattern_above words | awk '{print length, $0;}' | sort -n | tail | column
6 knoppy  7 alloquy 7 begorry 7 billowy   7 egilops
6 knotty  7 beefily 7 belloot 7 deglory   8 aegilops
```



## I before E? Except after C? And in words like "neighbor" and "weigh"?

```
grep 'ei' words | wc -l
5677
```
2.94:1

```
grep 'ie' words | wc -l
16690
grep 'cei' words | wc -l
300
```
2.87:1 … that difference begets a named exception?

```
grep 'cie' words | wc -l
863
grep 'eigh' words | wc -l
429
```
143:1 … okay, this one's pretty legit

```
grep 'iegh' words | wc -l
3
grep 'iegh' words
abiegh
driegh
skiegh
```

## Floating point number?

```
[-+]?([0-9]+\.?[0-9]*|\.[0-9]+)([Ee][-+]?[0-9]+)?
```

# Need Practice?

([a-z]*), (\S*), and (\w*)

Flags: ☑ g - Global ☐ i - Insensitive ☐ m - Mult

egular expressions can be a pain. This tool is designed t
evelopers **learn**, **practice**, **and compose** re
xpressions.

he HiFi RegExp tool is 100% JavaScript using jQuery. T
reated by New Media Campaigns, the team behind HiFi,
eneration CMS for web **designers**, **developers**
gencies. Enjoy!

| Matched Text | $1 | $2 |
|---|---|---|
| learn, practice, and... | learn | practice |
| designers, developer... | designers | developers |

```javascript
var regex = /([a-z]*), (\S*), and (\w*)/g;
var input = "your input string";
if(regex.test(input)) {
  var matches = input.match(regex);
  for(var match in matches) {
        alert(matches[match]);
  }
} else {
  alert("No matches found!");
}
```

*https://regexcrossword.com/challenges/*

# Just a Masochist?

/bul[rn]tl[coy]el[mtg]aljlisoln[hl]l[ae]dllevlshl[lnd]il[po]olls/
matches the last names of elected US presidents but not their opponents.

http://xkcd.com/1313/

# Day-to-day grep Workflows

- ❖ It is undeniably fun to find the perfect regexp.

  - ❖ It is undeniably more common to pipe them together until you exhaustively get exactly what you want

  - ❖ "Uh, name starts with J, upperclassman, BSE …"

```
. . . | grep J      | egrep '16|7' | egrep '(COS)|(ORF)|(ELE)|(CBE)|(MAE)|(CEE)'
        Janet 17 COS   Janet 17 COS   Janet 17 COS
        James 16 ORF   James 16 ORF   James 16 ORF
        Jenna 17 WWS   Jenna 17 WWS
        Julia 18 COS   Joaoa 16 HIS
        Jesse 19 BSE   …
        Joaoa 16 HIS
        …
```

# Day-to-day grep "Workflows"

```
grep $p p.sh                           grep '$p' p.sh
^C (after 30 seconds of hanging)       $python —c '#do nothing'


grep university cslist.txt  #immediately returns no matches
grep —i university cslist.txt
…
Johnson, Elizabeth Chief Reader   XAVIER  UNIVERSITY
Leyzberg, Daniel Reader PRINCETON  UNIVERSITY
Liu, David  Reader Indiana University Purdue University Fort Wayne
…


grep —i hu cslist.txt
Baker, Jeffrey  Reader Huntsville High School
Cunningham, Susie  Reader Indiana Academy of Science Math & Humanities
Hu, Chenglie  Reader CARROLL  UNIVERSITY
Huggins, James   Question Leader  KETTERING  UNIVERSITY
Hughes, Mary Ann Reader Albert Gallatin Area School District
Wang, Huanjing   Reader WESTERN KENTUCKY UNIVERSITY


grep —w Hu cslist.txt
Hu, Chenglie  Reader CARROLL  UNIVERSITY
```

# It's 2016, why do people still say this?

❖ Do you cringe every time you hear someone read a URL?

  ❖ H-T-T-P … Colon …

  ❖ Backslash Backslash

  ❖ W-W-W

❖ No doubt these people are prone to type it that way, too. And you are going to end up responsible for cleaning up their mess



SO MUCH PUN.COM

http://xkcd.com/208/

# Everybody Stand Back. Uh, …

```
http://good.url
this is a tent! /\
\/ for \/endetta
http:\\bad.url
```

```
grep \ slashandburn
^C
```

```
grep '\' slashandburn
grep: trailing backslash (\)
```

```
grep '\\' slashandburn
this is a tent! /\
\/ for \/endetta
http:\\bad.url
```

```
grep '\\\' slashandburn
grep: trailing backslash (\)
```

NB: you can use other delimiters to avoid having to escape every slash in your RE and replacement:
  `sed s#'\\\\'#'//'#g`

```
grep '\\\\' slashandburn
http:\\bad.url
```

```
sed s/'\\\\'/'//'/g slashandburn
sed: 1: "s/\\\\////g": bad flag in substitute command: '/'
```

```
sed s/'\\\\'/'\/\/'/g slashandburn
http://good.url
this is a tent! /\
\/ for \/endetta
http://bad.url
```

Wait, forgot to escape a space. Wheeeeee [taptaptap] eeeeee.

http://xkcd.com/208/

# There's No Escaping Escaping



```
\                    — BACKSLASH
\\                   — REAL BACKSLASH
\\\                  — REAL REAL BACKSLASH
\\\\                 — ACTUAL BACKSLASH, FOR REAL THIS TIME
\\\\\                — ELDER BACKSLASH
\\\\\\               — BACKSLASH WHICH ESCAPES THE SCREEN AND ENTERS YOUR BRAIN
\\\\\\\\             — BACKSLASH SO REAL IT TRANSCENDS TIME AND SPACE
\\\\\\\\\\           — BACKSLASH TO END ALL OTHER TEXT
\\\\\\\\\\\\\...     — THE TRUE NAME OF BA'AL, THE SOUL-EATER
```

I searched my .bash_history for the line with the highest ratio of special characters to regular alphanumeric characters, and the winner was:

**cat out.txt | grep -o "\\\[[(].*\\\[\\])][^)\]]*$"**

... I have no memory of this and no idea what I was trying to do, but I sure hope it worked.

# The grep Family

- ❖ fgrep
  - ❖ parallel search, but not actually full regexp (just fixed strings)
- ❖ agrep
  - ❖ "approximate" grep: search with errors permitted
- ❖ relatives that use similar regular expressions

  | | |
  |---|---|
  | ❖ ed | original Unix editor |
  | ❖ sed | stream editor |
  | ❖ vi, emacs, ... | editors |
  | ❖ lex | lexical analyzer generator |

- ❖ simpler variants
  - ❖ filename "wild cards" in Unix and other shells
  - ❖ "LIKE" operator in SQL, Visual Basic, etc.