PROJECT 2:    81646751(XI SUN) / 36088298(HONG CHANG)

Our Filter Rule of URL:

```
if(!mystring.matches("http://.*\\.ics\\.uci\\.edu.*"))
                    return false;
if(mystring.matches(".*\\.java"))
                    return false;
if(mystring.matches(".*ftp\\.ics\\.uci\\.edu.*"))
                    return false;
if(mystring.matches(".*seraja\\.ics\\.uci\\.edu.*"))
                    return false;
if(mystring.matches(".*fano\\.ics\\.uci\\.edu.*"))
                    return false;
if(mystring.contains("?"))
                    return false;
if(mystring.matches(".*edu/.*\\.(html|htm|php|jsp)"))
                    return true;
if(mystring.matches(".*edu/[^\\.]*"))
                    return true;
```

Q1. How much time did it take to crawl the entire domain?

It took 2h43min19sec (9799 seconds ) to crawl the entire domain after applying above filter rules.

Q2. How many unique pages did you find in the entire domain? (Uniqueness is established by the URL)

NOTE: This does not contain some pages that we have used our filter scheme to eliminate.\

There are unique 27912 pages in entire domain.

Q3. How many subdomains did you find?

NOTE: This does not contain some subdomains that we have used our filter scheme to eliminate. For example, we do not consider about ftp.ics.uci.edu. We find www.ics.uci.edu is the largest domain, which contains 23273 different pages.

There are 66 subdomains in entire domain. (URL, # unique pages)

alderis.ics.uci.edu,7
archive.ics.uci.edu,1182
asterix.ics.uci.edu,9
auge.ics.uci.edu,16
awareness.ics.uci.edu,500
calendar.ics.uci.edu,2
cert.ics.uci.edu,30
cgvw.ics.uci.edu,2
chime.ics.uci.edu,1
cleo.ics.uci.edu,5
cml.ics.uci.edu,1
computableplant.ics.uci.edu,19
cradl.ics.uci.edu,4
dblp.ics.uci.edu,2
deepthought.ics.uci.edu,6
drzaius.ics.uci.edu,6
duttgroup.ics.uci.edu,1
emme.ics.uci.edu,15

esl.ics.uci.edu,6
evoke.ics.uci.edu,50
flamingo.ics.uci.edu,11
fr.ics.uci.edu,10
frost.ics.uci.edu,59
galen.ics.uci.edu,100
graphics.ics.uci.edu,63
graphmod.ics.uci.edu,162
hana.ics.uci.edu,15
hcc.ics.uci.edu,3
hci.ics.uci.edu,1
hobbes.ics.uci.edu,7
hombao.ics.uci.edu,4
honors.ics.uci.edu,21
i-sensorium.ics.uci.edu,6
ipubmed.ics.uci.edu,4
isg.ics.uci.edu,11
jujube.ics.uci.edu,7

kdd.ics.uci.edu,101
luci.ics.uci.edu,166
metaviz.ics.uci.edu,8
mlearn.ics.uci.edu,258
mondego.ics.uci.edu,5
motifmap.ics.uci.edu,1
ngs.ics.uci.edu,14
phoenix.ics.uci.edu,19
ppopp2013.ics.uci.edu,14
psearch.ics.uci.edu,4
sami.ics.uci.edu,9
sconce.ics.uci.edu,16
sdcl.ics.uci.edu,118
sherlock.ics.uci.edu,5
sli.ics.uci.edu,91
snekker.ics.uci.edu,1
soc.ics.uci.edu,11
sourcerer.ics.uci.edu,10

| | | |
|---|---|---|
| sprout.ics.uci.edu,37 | vcp.ics.uci.edu,1190 | wics.ics.uci.edu,2 |
| student-council.ics.uci.edu,55 | vip.ics.uci.edu,9 | www-db.ics.uci.edu,3 |
| tastier.ics.uci.edu,1 | vision.ics.uci.edu,127 | www.ics.uci.edu,23273 |
| testlab.ics.uci.edu,9 | w3.ics.uci.edu,1 | xtune.ics.uci.edu,6 |

**Q4. What is the longest page?**

We detect that http://www.ics.uci.edu/~xhx/project/MotifMap/SNP/motif_sites_overlap_db_snp.list.html is longest page.

**Q5. What are the 500 most common words in this domain?**

NOTE:
1) According to the paper *word length, sentence length and frequency- ZIPF revised,* 99% of lexical words has length <15 character. We believe in practical search that input has very long length can be ignored, so in our tokenized file, we omit the word that has length > 15.
2) We have implemented our method that can choose including / not including numerical words for our statistics work. Because we find if we do not filter numbers, most of top words are biological terms or DNA/RNA sequence that containing numbers. Since this question asked us to output words, so we assume we only need to output alphabetical words, however, you can switch the method for further search engine projects in GenerateTokenFile.java
3) We use stopwordlist from http://www.ranks.nl/resources/stopwords.html (long version)
4) According to our observation, top words are related to biology field, student services and classes.

*Top 500 frequent words list:*
NOTE: The frequency rank list from high to low as following sequence:
left column> right column, then top row> bottom row, for example, mrna has frequency> protein, protein has frequency > data, which in second column.

| | | |
|---|---|---|
| mrna | chrx | ru |
| protein | type | will |
| hypothetical | finger | associated |
| domain | und | repeat |
| factor | zinc | polypeptide |
| receptor | transcription | group |
| family | alpha | phosphatase |
| transcript | class | potassium |
| homolog | beta | uci |
| member | channel | subunit |
| variant | box | cell |
| binding | gene | transmembrane |
| drosophila | nuclear | married |
| kinase | carrier | antigen |
| open | data | ics |
| chromosome | solute | oncogene |
| reading | subfamily | homeobox |
| frame | growth | interacting |

| | | |
|---|---|---|
| krakow | leukemia | version |
| computer | web | view |
| membrane | dec | guanine |
| rna | determining | problem |
| inhibitor | system | lymphoma |
| site | eos | immunoglobulin |
| rich | canon | nucleotide |
| image | multiple | org |
| syndrome | equiv | sodium |
| region | serine | early |
| software | specific | rho |
| ii | sry | jul |
| leucine | -box | ring |
| set | suppressor | fibroblast |
| size | university | protease |
| ras | public | gtpase |
| lim | http | string |
| java | mouse | interleukin |
| technical | jun | activator |
| cadherin | file | synthase |
| dehydrogenase | code | protocadherin |
| tyrosine | nov | student |
| calcium | oct | adhesion |
| method | yeast | ligand |
| original | number | activated |
| details | mar | david |
| gamma | protein-coupled | actin |
| time | transporter | collagen |
| mitochondrial | homeo | sulfate |
| enhancer | project | response |
| ribosomal | cerevisiae | algorithms |
| sequence | feb | oxidase |
| iso | aug | matrix |
| acid | complex | computing |
| dna | apr | delta |
| systems | domains | atpase |
| viral | elegans | basic |
| glutamate | school | listing |
| superfamily | regulator | methods |
| voltage-gated | cancer | polymerase |
| activating | bren | sema |
| forkhead | element | transporting |
| integration | design | tree |
| sex | large | object |
| tumor | molecule | deleted |
| course | students | t-cell |
| program | hormone | crw |
| science | server | learning |
| work | sep | classes |
| regulatory | user | small |
| neuronal | irvine | help |
| jan | ankyrin | source |
| b-cell | list | lab |
| motif | precursor | enzyme |

| | | |
|---|---|---|
| pou | programming | ataxin |
| cll | people | morphogenetic |
| f-box | heparan | package |
| paper | static | short |
| activin | candidate | point |
| cytochrome | -like | retinoic |
| ionotropic | pdf | transforming |
| homology | isr | test |
| frames | resource | high |
| example | application | order |
| leucine-rich | processing | week |
| myeloid | california | neurexin |
| search | paired | de |
| eppstein | cycle | slit |
| similarity | subject | well |
| contact | single | cs |
| int | find | machine |
| split | cytoplasmic | interactive |
| field | necrosis | dual |
| component | three | glypican |
| lang | ubiquitin | analysis |
| sara | bone | pdz |
| iii | histone | collection |
| expressed | problems | brain |
| pleckstrin | phd | state |
| virus | saturday | t-box |
| engineering | access | synaptotagmin |
| synthetase | including | cysteine |
| support | reductase | translation |
| neural | product | department |
| phospholipase | fibronectin | technology |
| graph | phosphoprotein | prev |
| signal | email | g-protein |
| induced | algorithm | context |
| yancees | division | html |
| myosin | international | orphan |
| btb | differentiation | muscle |
| values | slide | previous |
| kruppel-like | network | message |
| ecotropic | cation | human |
| conference | general | acidic |
| lecture | case | distributed |
| wingless-type | homeodomain | informatics |
| mmtv | based | july |
| process | contactin | transducin-like |
| development | usm | banns |
| june | returns | function |
| assignment | interface | double |
| acm | final | wd |
| protein-like | endothelial | translocation |
| light | applications | management |
| model | points | john |
| sorting | latrophilin | personal |
| metalloprotease | write | text |

october
required
read
networks
questions
requirements
exchanger
carcinoma
link
create
sub-family
eukaryotic
files
thyroid
vision
threonine
rescue
disease
description
iroquois
note
tz
graduate
low
long
start
ets
ieee
breast
plasma
form
click
bruno-like
projects
specificity
abraham
resources
itr
copy
uc
anion

langsam
ps
dataguard
br
directory
repeats
integrin
glucosamine
zipper
document
theory
password
exchange
special
lead
change
cassette
cyclin
initiation
current
signalling
dominant
chain
models
return
chemokine
putative
collaboration
fragile
kh
control
programs
atp-binding
reference
apoptosis
embryonic
cyclase
papers
death
office
inducible

tm
good
space
request
proteasome
son
abstract
riken
sets
working
smad
april
protocol
cdna
dead
internet
march
max
security
variable
nexin
co-repressor
avian
second
policy
heavy
gef
glycoprotein
include
it's
database
rar-related
elongation
olfactory
interferon
pr
ny
insulin-like
notes

## Q6. What are the 20 most common 2-grams?

*Top 20 2-gram list*
NOTE: The frequency rank list from high to low as following sequence:
left column> right column, then top row> bottom row.

protein mrna
hypothetical protein
transcript variant
mrna hypothetical
reading frame

open reading
chromosome open
variant mrna
binding protein
frame mrna

mrna protein
mrna chrx
homolog drosophila
zinc finger
transcription factor

| domain mrna | finger protein | mrna | chromosome |
| family member | drosophila mrna | | |

Q7. (extra problem) How many unique pages excluding similar contents?

There are 9387 unique pages(~ 59.7% of original) excluding similar contents in the whole domain.

NOTE: We used shingle algorithm to compute the similarity between different pages. Because of scale of the executing time, we chose 0.5 as our resemblance threshold and containment threshold. We did the computing and compressing at the same time, so we can do accomplish this task in several hours. (details could be discussed during face-to-face talk)