GENONE 569

Class 4: Workflow automation tools

Essential UNIX

Function List the contents of a directory Sort a file line by line Search inside files for pattern Write programs that process files line by line Transform a file line by line Cut out selected columns of a file Join two files together based on columns Multi-purpose tools to manipulate tab-delimited files Multi-purpose tool to manipulate BED files Multi-purpose tool to manipulate SAM/BAM files

Command S sort grep awk sed cut join datamash bedtools samtools

Essential commands



Problem: search for strings in other strings

Solution:



grep <regexp> file

- Most characters match themselves The regular expression "test" matches the string 'test', and only that string
- [x] matches any one of a list of characters "[abc]" matches 'a', 'b', or 'c'
- [^x] matches any one character that is not included in x "[^abc]" matches any single character except 'a','b',or 'c'

- "(abc)+" matches 'abc', 'abcabc', 'abcabcabc', etc.
- "." matches any single character Parentheses can be used for grouping
- x y matches x or y "this that" matches 'this' and 'that', but not `thisthat'.

- x* matches zero or more x's "a*" matches '', 'a', 'aa', etc.
- x+ matches one or more x's "a+" matches 'a','aa','aaa', etc.
- x? matches zero or one x's "a?" matches ' ' or ' a'
- $x\{m, n\}$ matches *i* x's, where $m \le i \le n$ "a{2,3}" matches 'aa' or 'aaa'

Example: email addresses

How can we easily tell these two apart?

coletrap@uw.edu

Here's a pattern to match simple email addresses: \w+@(\w+\.)+(com|org|net|edu)

spam@go.away

- "\d" matches any digit; "\D" any non-digit
 "\s" matches any whitespace character; "\S" any non-
- "\s" matches any whitespace whitespace character
- "\w" matches any alphanumeric character; "\W" any nonalphanumeric character
- "^" matches the beginning of the string; "\$" the end of the string
- "\b" matches a word boundary; "\B" matches a character that is not a word boundary

Sample problem 1

- Download "grep_sed_example1.txt"
- Find the lines in the notes that mention 'cell'
- Find the lines in the notes that talk about 'A549'
- Find the lines that talk about either A549 or MCF7 (use regex)
- Find the lines that end with a cell line id (i.e. A549 or MCF7, use regex!)

Change instances of A549 to MCF7

Change only the second instance of A549 to A549 LUNG

Change instances of A549 to MCF7, but without stating "A549" (i.e. use regex)

How do we make all these changes at the same time?

Sample problem 2

Download "grep_sed_example1.txt"



Single-cell RNA-seq with combinatorial cellular indexing



Cao & Packer et al, Science 2017

Pool amplicons & deep sequence to generate single cell 3' digital gene expression profiles



The sci-RNA-seq read layout



sci-RNA-seq: single-cell RNA-seq on whole animals





Anatomy of a pipeline

The sci-RNA-seq pipeline



Repo structure for Cao pipeline

🗸 🛅 sci-RNA-seq-pipeline-scripts		Today at 2:05 PM		Folder
assign-reads-to-genes.jun-pipeline-compatible.py	\bigcirc	Today at 1:10 PM	4 KB	Python Source
📝 assign-reads-to-genes.jun-pipeline-compatible.sh		Today at 1:10 PM	772 bytes	Shell Script
assign-reads-to-genes.py	\bigcirc	Today at 1:10 PM	4 KB	Python Source
assign-reads-to-genes.sh		Today at 1:10 PM	764 bytes	Shell Script
count-rRNA-reads.sh	\bigcirc	Today at 1:58 PM	385 bytes	Shell Script
count-UMI-per-sample.sh		Today at 1:10 PM	818 bytes	Shell Script
debug-assign-reads-to-genes.sh	\bigcirc	Today at 1:10 PM	708 bytes	Shell Script
📝 kallisto.sh		Today at 1:10 PM	339 bytes	Shell Script
👔 knee-plot-merge-samples.R	\bigcirc	Today at 1:10 PM	1 KB	Rez Source
👔 knee-plot.R		Today at 1:10 PM	1 KB	Rez Source
make-fasta-for-kallisto.sh	\bigcirc	Today at 1:10 PM	618 bytes	Shell Script
Nextera.v2.0.P7.oligo.rev.comp		Today at 1:10 PM	1 KB	Document
prelim-umi-count-stats.sh	\bigcirc	Today at 1:10 PM	555 bytes	Shell Script
put-read1-info-in-read2.awk		Today at 1:10 PM	2 KB	Document
put-read1-info-in-read2.fancy-sample-mapping.awk		Today at 1:10 PM	4 KB	Document
put-read1-info-in-read2.fancy-sample-mapping.sh		Today at 1:10 PM	678 bytes	Shell Script
put-read1-info-in-read2.sh		Today at 1:10 PM	657 bytes	Shell Script
📝 rmdup-and-make-split-bed.sh		Today at 1:10 PM	532 bytes	Shell Script
rmdup.awk		Today at 1:10 PM	1 KB	Document
🔹 run-bcl2fastq.sh		Today at 1:10 PM	423 bytes	Shell Script
🔹 run-trim-galore-for-c-elegans.sh		Today at 1:10 PM	370 bytes	Shell Script
🔹 run-trim-galore.sh		Today at 1:10 PM	358 bytes	Shell Script
samtools-filter-sort.sh		Today at 1:10 PM	363 bytes	Shell Script
sci-RNA-seq.P5.oligo.rev.comp		Today at 1:10 PM	1 KB	Document
sci-RNA-seq.RT.oligos		Today at 1:10 PM	1 KB	Document
STAR-alignReads-8-cores.sh		Today at 1:10 PM	610 bytes	Shell Script
STAR-alignReads-moreMem.sh		Today at 1:10 PM	611 bytes	Shell Script
STAR-alignReads.sh		Today at 1:13 PM	611 bytes	Shell Script
UMI-count-rollup.sh		Today at 2:05 PM	289 bytes	Shell Script
sci-RNA-seq.pipeline.elegans.sh		Today at 2:02 PM	13 KB	Shell Script

- Start here

sci-RNA-seq.pipeline.elegans.sh

Put read 1 info (RT well, UMI) into read 2 read name #_____ echo "Moving read 1 info into read 2 name"

cd \$WORKING_DIR

mkdir combined-fastq mkdir file-lists-for-r1-info-munging mkdir put-r1-info-in-r2-logs

ls \$READS_DIR/fastq | grep _R1_ | grep -v Undetermined | split -l 25 -d - file-lists-for-r1-info-munging/





sci-RNA-seq.pipeline.elegans.sh

Put read 1 info (RT well, UMI) into read 2 read name echo "Moving read 1 info into read 2 name"

cd \$WORKING_DIR

mkdir combined-fastq mkdir file-lists-for-r1-info-munging mkdir put-r1-info-in-r2-logs

"pipe": sends the output of one command into the next command as input

grep: prints lines from the input that match the regular expression provided

ls \$READS_DIR/fastq | grep _R1_ | grep -v Undetermined | split -l 25 -d - file-lists-for-r1-info-munging/

split: split up a file into chunks. In this case, chunks of 25 lines





sci-RNA-seq.pipeline.elegans.sh





```
Put-read1-into-read2.awk
```

```
BEGIN {
    read_num = 0;
    hits = 0;
    bases[1] = "A";
    bases[2] = "C";
    bases[3] = "G";
    bases[4] = "T";
    p5 row = substr(PCR COMBO, 1, 1);
    p5_col = substr(PCR_COMB0, 2, 2);
    p7\_row = substr(PCR\_COMB0, 4, 1);
    p7_col = substr(PCR_COMB0, 5, 2);
    single_sample = "";
} {
} END {
}
```

Extract RT barcode & UMI from R2

printf("%d\t%d\t%.3f\t(RT barcode matches, total reads, proportion)\n", hits, read_num, hits / read_num) > "/dev/stderr";





Snakemake

Project Phase II



Fornow