

Bioinformatics one-liners

Useful bash one-liners for useful for bioinformatics.

<https://github.com/stephenturner/oneliners>

Download the [PDF](#) here.

Sources:

- <http://sed.sourceforge.net/sed1line.txt>
- <https://github.com/lh3/seqtk>
- <http://lh3lh3.users.sourceforge.net/biounix.shtml>
- <http://genomespot.blogspot.com/2013/08/a-selection-of-useful-bash-one-liners.html>

awk, sed

Sum column 1 of file.txt:

```
awk '{sum+=$1} END {print sum}' file.txt
```

Number each line in file.txt:

```
sed = file.txt | sed 'N;s/\n/ /'
```

Get unique entries in file.txt based on column 1 (takes only the first instance):

```
awk '!arr[$2]++' file.txt
```

Replace all occurrences of **foo** with **bar** in file.txt:

```
sed 's/foo/bar/g' file.txt
```

Convert a FASTQ file to FASTA:

```
sed -n '1~4s/^@/>/p;2~4p' file.fq > file.fa
```

Extract every 4th line starting at the second line (extract the sequence from FASTQ file):

```
sed -n '2~4p' file.fq
```

Basic sequence statistics. Print total number of reads, total number unique reads, percentage of unique reads, most abundant sequence, its frequency, and percentage of total in file.fq:

```
cat myfile.fq | awk '((NR-2)%4==0){read=$1;total++;count[read]++}END{for(read in count){if(!max||count[read]>max){max=count[read];maxRead=read};if(count[read]==1){unique++}};print total,unique,unique*100/total,maxRead,count[maxRead],count[maxRead]*100/total}'
```

Convert .bam back to .fastq:

```
samtools view file.bam | awk 'BEGIN {FS="\t"} {print "@" $1 "\n" $10 "\n+\n" $11}' > file.fq
```

sort, uniq, cut, etc.

Count the number of unique lines in file.txt

```
cat file.txt | sort | uniq | wc -l
```

Find number of lines shared by 2 files:

```
sort file1 file2 | uniq -d
```

Find the most common strings in column 2:

```
cut -f2 file.txt | sort | uniq -c | sort -k1nr | head
```

Pick 10 random lines from a file:

```
shuf file.txt | head -n 10
```

Print rows where column 3 is larger than column 5 in file.txt:

```
awk '$3>$5' file.txt
```

Compute the mean of column 2:

```
awk '{x+=$2}END{print x/NR}' file.txt
```

Extract fields 2, 4, and 5 from file.txt:

```
awk '{print $2,$4,$5}' input.txt
```

Print all possible 3mer DNA sequence combinations:

```
echo {A,C,T,G}{A,C,T,G}{A,C,T,G}
```

find, xargs, and GNU parallel

Download GNU parallel at <https://www.gnu.org/software/parallel/>.

Search for .bam files anywhere in the current directory recursively:

```
find . -name "*.bam"
```

Delete all .bam files:

```
find . -name "*.bam" | xargs rm
```

Rename all .txt files to .bak (backup *.txt before doing something else to them, for example):

```
find . -name "*.txt" | sed "s/\.txt$//" | xargs -i echo mv {}.txt {}.bak | sh
```

Chastity filter raw Illumina data (grep reads containing **:N:**, append (-A) the three lines after the match containing the sequence and quality info, and write a new filtered fastq file):

```
find *fq | parallel "cat {} | grep -A 3 '^@.*[^:]*:N:[^:]*:' | grep -v '^\\-\\-$' > {}.filt.fq"
```

Run FASTQC in parallel 12 jobs at a time:

```
find *.fq | parallel -j 12 "fastqc {} --outdir ."
```

Index your bam files in parallel, but only echo the commands (`--dry-run`) rather than actually running them:

```
find *.bam | parallel --dry-run 'samtools index {}'
```

seqtk

Download seqtk at <https://github.com/lh3/seqtk>. Seqtk is a fast and lightweight tool for processing sequences in the FASTA or FASTQ format. It seamlessly parses both FASTA and FASTQ files which can also be optionally compressed by gzip.

Convert FASTQ to FASTA:

```
seqtk seq -a in.fq.gz > out.fa
```

Convert ILLUMINA 1.3+ FASTQ to FASTA and mask bases with quality lower than 20 to lowercases (the 1st command line) or to **N** (the 2nd):

```
seqtk seq -aQ64 -q20 in.fq > out.fa  
seqtk seq -aQ64 -q20 -n N in.fq > out.fa
```

Fold long FASTA/Q lines and remove FASTA/Q comments:

```
seqtk seq -c160 in.fa > out.fa
```

Convert multi-line FASTQ to 4-line FASTQ:

```
seqtk seq -l0 in.fq > out.fq
```

Reverse complement FASTA/Q:

```
seqtk seq -r in.fq > out.fq
```

Extract sequences with names in file `name.lst`, one sequence name per line:

```
seqtk subseq in.fq name.lst > out.fq
```

Extract sequences in regions contained in file `reg.bed`:

```
seqtk subseq in.fa reg.bed > out.fa
```

Mask regions in `reg.bed` to lowercases:

```
seqtk seq -M reg.bed in.fa > out.fa
```

Subsample 10000 read pairs from two large paired FASTQ files (remember to use the same random seed to keep pairing):

```
seqtk sample -s100 read1.fq 10000 > sub1.fq  
seqtk sample -s100 read2.fq 10000 > sub2.fq
```

Trim low-quality bases from both ends using the Phred algorithm:

```
seqtk trimfq in.fq > out.fq
```

Trim 5bp from the left end of each read and 10bp from the right end:

```
seqtk trimfq -b 5 -e 10 in.fa > out.fa
```

Untangle a FASTQ file. If a FASTQ file has paired-end reads intermingled, and you want to separate them into separate /1 and /2 files, and assuming the /1 reads precede the /2 reads:

```
seqtk seq -l0 tangled.fq | gawk '{if ((NR-1) % 8 < 4) print >> "separate_1.fq";  
else print >> "separate_2.fq"}'
```

Other generally useful aliases for your .bashrc

Never type `cd ../../..` again:

```
alias ..='cd ..'  
alias ...='cd ../../'  
alias ....='cd ../../..'  
alias .....='cd ../../../../'  
alias .....='cd ../../../../..'
```

Ask before removing or overwriting files:

```
alias mv="mv -i"  
alias mf="mv -i"  
alias cp="cp -i"  
alias rm="rm -i"
```

My favorite `ls` aliases:

```
alias ls="ls -lp --color=auto"  
alias l="ls -lhGgo"  
alias ll="ls -lh"  
alias la="ls -lhGgoA"  
alias lt="ls -lhGgotr"  
alias lS="ls -lhGgoSr"  
alias l.="ls -lhGgod .*"   
alias lhead="ls -lhGgo | head"  
alias ltail="ls -lhGgo | tail"  
alias lmore='ls -lhGgo | more'
```

Use `cut` on space- or comma- delimited files:

```
alias cuts="cut -d \" \""  
alias cutc="cut -d \",\""
```

Pack and unpack tar.gz files:

```
alias tarup="tar -zcf"  
alias tardown="tar -zxf"
```