

BCGES short courses, session 1: Introduction, file formats, shell scripting, Galaxy

Vincent Plagnol

The aim of this first session is to introduce the tools that we will use to manipulate sequence data. Generally speaking, we will deal with either **R**, bash scripting or **Galaxy**.

R is not usually the preferred way to deal with HTS data. One reason for this is that the general approach of **R** is to load data into the RAM as a first instance, which is not practical when the datasets are very large. Shell scripts (aka command line tools) that read the data on a line per line basis are usually preferred. Nevertheless, **R** has developed tools to overcome these limitations and there is, in fact, surprisingly much that one can do with it. It is also a practical tool for teaching purposes, which is useful in the context of these short courses. So we will use **R** mostly to play with the data and show the sort of things one may want to do with it.

Contents

1 Basic fastq reading and quality control (1h)	2
1.1 Basic shell scripting to read NGS files (20 minutes)	2
1.2 Using the shortRead package in R (15 minutes)	3
1.3 Using the Galaxy server (30 minutes)	4
2 Reading and manipulating BAM files (about 1.5h)	5
2.1 Using shell scripts and samtools (40 minutes)	5
2.2 Understanding the SAM/BAM headers (10 minutes)	5
2.3 Using R and Rsamtools (20 minutes)	6
3 Merging overlapping paired-end reads (20 minutes)	7

1 Basic fastq reading and quality control (1h)

1.1 Basic shell scripting to read NGS files (20 minutes)

Many of the standard HTS formats are simple text files, with tightly defined specifications to allow effective parsing. Some of these files can be compressed and/or indexed to enable quicker access. The first essential step is to be comfortable with reading/compressing/uncompressing various text files.

```
head ../data/fastq_files/fastq1_1.txt

## @A81CH8ABXX:4:1101:1524:1813#NGACCAAT/1
## NACCACTCAGCTCTGGCCAATTATTGCCGTGCAGGAGTGTGGGCTCCTAGTGGCAGGGGGTCTGGAAGTGTGGAAGAAGCAGGCAAACGC
## +
## BQXXQY[V[Yccc_c__V\__ccc__\X~~~~BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
## @A81CH8ABXX:4:1101:1583:1836#NGACCAAT/1
## NTCCCCAAAGCACAGGGCTCAGCTCCAGAGGGAGACGGGCTGGGCTGTCAGCGGGCCCAGGGGCACGCCACTGTTGAGAACAACCTGGTTG
## +
## BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
## @A81CH8ABXX:4:1101:1729:1852#TGACCAAT/1
## NCTATGGACTGTGGTAAAGCTAGGATTAGTAACCAGACATTACTTACCTTGGCTCCGATCTGGTTGCCACACTGGCCAATCTGAATATGC
```

```
tail ../data/fastq_files/fastq1_1.txt

## +
## ggggggggggggggggggggggggggggggggggggggggggggggggggggggfggggggaeede[~d[_Y`~b]a[bZXZ[Y`U^BBBBBBB
## @A81CH8ABXX:4:1101:18871:2349#TGACCAAT/1
## TCTGTGCCGATCCCAGAGTCCCTGGGTGAAGATGATTCTCAATAAGCTGAGCCAATGAAGAGCCTACTCTGATGACCGTGGCCTTGGCT
## +
## eedeee]decdbdbacca`bdbbdac_adWdb__dcee\ceeeeedeeee_ede~ddbdaUWbbd__aa[dee`e\bb_cad`dZcZc\b
## @A81CH8ABXX:4:1101:18909:2349#TGACCAAT/1
## TGGTCACTAGCTCTGGATGGTGTCTCCAACAATGGTCACTTCTCGCGGAGAACACAAACACCAGCATCACAGCGCTGGGTTC CATGGATG
## +
## gggggfgggggggggggdgfdgfgggggggffggggdggggggfgfgfbgggggggggggfgggggcg^gggeffedfgggfeggeeeg
```

If you need more information, shell scripts have useful manual page that can be accessed using the `man` command.

```
man head
man tail
```

One can parse a text file using a variety of commands, and it is useful to become familiar with the following:

```
cat ../data/fastq_files/fastq1_1.txt
less ../data/fastq_files/fastq1_1.txt
more ../data/fastq_files/fastq1_1.txt
less -S ../data/fastq_files/fastq1_1.txt
```

Exercise: Can you see the differences between these different ways of reading a file?

A routine that is often useful is `wc` that counts the words/character/lines of a file. Try:

```
wc ../data/fastq_files/fastq1_1.txt
wc -l ../data/fastq_files/fastq1_1.txt
```

Exercise: Using the man page, find a way to print the first 20 lines of a fastq file (and what about the last 20 lines)?

Exercise: Based on the `wc` output, how many reads do these fastq files contain?

1.2 Using the shortRead package in R (15 minutes)

We start by loading one of the most relevant library, called “ShortReads”. This package may not be installed but it can easily done so by running:

```
source("http://bioconductor.org/biocLite.R")
biocLite("ShortRead")
```

```
library('ShortRead')
```

As a starting point it is possible to read some of the examples fastq files and create relevant R objects.

```
fastq1.1 <- readFastq('../data/fastq_files/fastq1_1.txt')
fastq1.2 <- readFastq('../data/fastq_files/fastq1_2.txt')
```

We can now display the sequences and the qualities. Note that specific classes have been defined to store each of these objects. A lot of work has gone into figuring out how to do this.

```
reads <- sread(fastq1.1)
class(reads)

## [1] "DNAStringSet"
## attr(,"package")
## [1] "Biostrings"

head(as.character(reads))

## [1] "NACCACTCAGCTCTGGCCAATTATTGCCGTGCAGGAGTGTGGGCTCCTAGTGGCAGGGGGTCTGGAACCTGTGGAAGAAGCAGGCAAACGC"
## [2] "NTCCCCAAAGCACAGGGCTCAGCTCCAGAGGGAGACGGGCTGGGCTGTGACGGGGCCAGGGGCACGCCACTGTTCAGAACAACTGGTTG"
## [3] "NCTATGGACTGTGGTAAAGCTAGGATTAGTAACCAGACATTACTTACCTTGGCTCCGATCTGGTTGCCACACTGGCCAATCTGAATATGG"
## [4] "NACTGAAAAAGGATGCTTTGGAAAAAGAAAGTGGGTCTGGCAACACTGACTCAACCTTGAATTCCCCGCACGATGACACGGATGACAGGG"
## [5] "GTTATTTAAGCCACCCAGTCTGTGTTTGTATGGCAGGCTGAGGAGACTATGACAGGAACCAACACAAAAAAACAAAAACTCTGGAGG"
## [6] "ATTTTGTAACCTCTTATCCTTAGATCCAAAAGATATGTTTCATCTAGGCTTGATAAGCACATGTGCATTATACCACACTCTATAGTTCT"

ids <- id(fastq1.1)
class(ids)

## [1] "BStringSet"
## attr(,"package")
## [1] "Biostrings"

head(as.character(ids))

## [1] "A81CH8ABXX:4:1101:1524:1813#NGACCAAT/1" "A81CH8ABXX:4:1101:1583:1836#NGACCAAT/1"
## [3] "A81CH8ABXX:4:1101:1729:1852#TGACCAAT/1" "A81CH8ABXX:4:1101:1642:1867#TGACCAAT/1"
## [5] "A81CH8ABXX:4:1101:1715:1891#TGACCAAT/1" "A81CH8ABXX:4:1101:1624:1941#TGACCAAT/1"
```

It is also possible to use this package to look at quality scores. For example, following up on the example above:

```
quals <- quality(fastq1.1)
class(quals)

## [1] "SFastqQuality"
## attr(,"package")
## [1] "ShortRead"

quals

## class: SFastqQuality
## quality:
## A BStringSet instance of length 2500
## width seq
```

```
## [1] 90 BQXXQY[V[Yccc_c__V\_ccc__\[X^~^~^BBBBBBB...BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
## [2] 90 BBBB...BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB...BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
## [3] 90 BWSSQVUVUTTQVSUUUW__BBBBBBBBBBBBBBBBBBBB...BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
## [4] 90 BBBB...BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB...BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
## [5] 90 gefggggdgegfgggdfegee`eeddbdfZeffde`dBBB...BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
## ... ..
## [2496] 90 gfggggggggggggggfgggefegfgcgaafdebfgggad...eee_ddeabdb_cedeed0Z`^IZXYW`^YW]\`ccX_adbdb
## [2497] 90 gggggggggggggfggggggggdgfgcgggfgfgggggfgg...bggfgefggggegceggedfdeeg`dbdee`bNddee`a
## [2498] 90 ggggggggggggggggggggggggggggggggggggggggggg...egfgggggaeeede[~d[_Y`^b]a[bZXZ[Y`U^BBBBBB
## [2499] 90 eedeee]decdbdbacca`bdbbdac_adWdb__dcee\cee...ee_ede`ddbdaUWbbd__aa[dee`e\bb_cad`dZcZc\b
## [2500] 90 gggggfgggggggggggdgdgfgggggggfgfggggdggggg...fbggggggggggfgggggcgg^gggeffedfgggfeggeeg
```

The **ShortRead** package will attempt to guess what these quality scores mean, for example see:

```
encoding(quality(fastq1.1))

## ; < = > ? @ 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 [
## -5 -4 -3 -2 -1 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27
## \ \ ] ^ _ ` a b c d e f g h i
## 28 29 30 31 32 33 34 35 36 37 38 39 40 41

fastq2.1 <- readFastq(' ../data/fastq_files/fastq2_1.txt')
encoding(quality(fastq2.1))

## ! " # $ % & ' ( ) * + , - . / 0 1 2 3 4 5 6 7 8 9 : ; < = > ? @ A
## 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32
## B C D E F G H I J
## 33 34 35 36 37 38 39 40 41
```

(Optional, and somewhat difficult) Exercise: Generate a plot showing the average Phred score as a function of the position in the read.

1.3 Using the Galaxy server (30 minutes)

Start by identifying a working instance of the Galaxy server from this [location](#). There are multiple choices here, so feel free to experiment. In case of doubt, I propose that you use the [main Galaxy instance](#) which I know works for the purpose of this exercise. However, you can try different options but be ready for small and subtle differences and be flexible associated with different versions of the software.

Exercise: Create an account on one of these Galaxy servers and upload the pair of fastq *fastq2.1.txt* and *fastq2.2.txt*. Run the **fastqc** code on these files (note that a standalone version also exists and can be found [here](#)).

(Optional) Exercise: Run **fastqc** locally (after downloading it from the [web](#)) and compare the output with what you get from Galaxy (it should be identical!).

2 Reading and manipulating BAM files (about 1.5h)

BAM files are compressed files that contain the information from the FASTQ files, plus additional information about the location where the reads map. A key feature of the BAM files is that they can be indexed, i.e. an associated file contains information about where each of the reads are located in the file. It allows very quick retrieval of reads that map to a given genomic location, which is the typical way one wants to use BAM files (for example, to extract all the reads that map to a gene of interest in order to find rare variants).

2.1 Using shell scripts and samtools (40 minutes)

`samtools` is the key piece of software that is used to read, write and index BAM files. The [manual page](#) is the first place to go to find information about how to use `samtools`. A very useful feature of `samtools` is that it can work over a ftp site, by downloading the index locally and only pulling the reads that are relevant. That allows to access rich datasets online without having to download very large files. The exercise below illustrates some of these features.

Exercise: The following should be useful as an introduction to the sort of things one may want to do with `samtools`. The manual page should have all the commands and ideas to go through these, so best to have a go and try.

1. Download the index of all the aligned file from the 1,000 Genomes. You can for example do it using:

```
wget -O results/1KG_index.tab ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/exome.alignment.index
```

2. Using `samtools view` over the web, download the *BRCA1* reads for the sample HG00118. You will first need to find where that file is located on the web, using the index you just downloaded. Make sure that your output is in BAM format.

Note: This step may seem difficult at first so the resulting BAM file should be available by default as part of the github set of files. You will also need the address of the 1KG ftp site: `ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp`.

3. Index this BAM file using `samtools index`.
4. Still using `samtools view`, subset the BAM file for the first coding exon (not the UTR) of *BRCA1* (use transcript ENST00000357654 from Ensembl) and output in SAM format.
5. Exercise: What do these lines do?

```
samtools view -f 0x0002 results/BRCA1_HG00118.bam | wc -l
samtools view -F 0x0002 results/BRCA1_HG00118.bam | wc -l
samtools view -F 0x0040 results/BRCA1_HG00118.bam | wc -l

## 6519
## 159
## 3347
```

2.2 Understanding the SAM/BAM headers (10 minutes)

Key information about a BAM file can be obtained from the headers. `Samtools` is once again the key tool to read these data. You can view these headers using:

```
samtools view -H results/BRCA1_HG00118.bam > results/headers.txt
```

```
less -S results/headers.txt
```

Exercise: Go through the output of `samtools view -H` and make sure you understand what all the fields mean. This will be discussed in class. In the case of the downloaded 1KG file, much processing has happened so not all fields will make sense, but the important thing is to understand the general philosophy of the file. The ID, LB, SM tags are particularly useful.

2.3 Using R and Rsamtools (20 minutes)

The **Rsamtools** package in R is very effective to parse BAM files, and extremely memory efficient, making full use of BAM indexes. Look at the example below for example. Inspect the output object called **bam.reads**. Can you understand its structure? See what it contains and how the data are organised? We will be using these tools later on in the CNV analysis section.

```
library(Rsamtools)
library(GenomicRanges)

which <- GRanges(seqnames=Rle("21"),
                 IRanges(start = 43000000, end = 45000000))

what <- c("rname", "strand", "pos", "qwidth", "seq")
param <- ScanBamParam(which=which, what=what)

bam.reads <- scanBam(file = "../data/BAM_files/HG00251.mapped.ILLUMINA.bwa.GBR.exome.20121211.bam",
                    param=param)
names(bam.reads[[1]])

## [1] "rname" "strand" "pos" "qwidth" "seq"

head(bam.reads[[1]]$pos)

## [1] 43823883 43823889 43823889 43823894 43823897 43823904
```

3 Merging overlapping paired-end reads (20 minutes)

Another useful thing to be aware of is the possibility to merge read pairs that overlap in the middle. This happens when the combined length of both reads exceeds the length of the DNA fragment they originate from. A review of tools to perform this task is available [here](#). I propose here to have a look at the program **PEAR**, which has been useful for the applications I had in mind. Start by downloading the pre-compiled binary to your own computer (this program is not preinstalled). You should be able to run something like the command below, after changing the path to the appropriate location of the executable:

```
bin/pear-0.9.4-64 -e -v 10 \  
-f ../data/fastq_for_merging/reads1.fq -r ../data/fastq_for_merging/reads2.fq -o results/merged -c 40
```

Check what the output looks like, in particular the read length of the assembled reads. Have a look at the options that PEAR uses, and see how the choice of the `v` argument affects the final results.

Warning: Note that PEAR requires that you know which of the reads are forward and reverse, which is probably not the case for a fastq file that is not stranded or that has been processed. So be a bit careful before running this code on a generic fastq.

Session info

```
sessionInfo()

## R version 3.1.0 (2014-04-10)
## Platform: x86_64-unknown-linux-gnu (64-bit)
##
## locale:
##  [1] LC_CTYPE=en_US.iso885915      LC_NUMERIC=C                LC_TIME=en_US.iso885915
##  [4] LC_COLLATE=en_US.iso885915    LC_MONETARY=en_US.iso885915 LC_MESSAGES=en_US.iso885915
##  [7] LC_PAPER=en_US.iso885915      LC_NAME=C                   LC_ADDRESS=C
## [10] LC_TELEPHONE=C                LC_MEASUREMENT=en_US.iso885915 LC_IDENTIFICATION=C
##
## attached base packages:
## [1] parallel  methods  stats      graphics  grDevices  utils      datasets  base
##
## other attached packages:
## [1] ShortRead_1.20.0      Rsamtools_1.14.3      lattice_0.20-29       Biostrings_2.30.1
## [5] GenomicRanges_1.14.4 XVector_0.2.0         IRanges_1.20.7        BiocGenerics_0.8.0
## [9] knitr_1.6
##
## loaded via a namespace (and not attached):
##  [1] Biobase_2.22.0        bitops_1.0-6          evaluate_0.5.5        formatR_0.10
##  [5] grid_3.1.0            highr_0.3             hwriter_1.3           latticeExtra_0.6-26
##  [9] RColorBrewer_1.0-5    stats4_3.1.0          stringr_0.6.2         tools_3.1.0
## [13] zlibbioc_1.8.0
```