

# RNA-seq data analysis practical - San Michele all'Adige (Trento), Italy

Mar González-Porta

2013/07/02

## Contents

<b>Introduction</b>	<b>2</b>
<b>Dealing with raw data</b>	<b>2</b>
The FASTQ format . . . . .	2
Quality assessment (QA) . . . . .	3
Filtering FASTQ files . . . . .	4
De-multiplexing samples . . . . .	6
Aligning reads to the genome . . . . .	6
<b>Dealing with aligned data</b>	<b>8</b>
The SAM/BAM format . . . . .	8
Visualising aligned reads . . . . .	9
Filtering BAM files . . . . .	10
Counting reads overlapping annotated genes . . . . .	10
With htseq-count . . . . .	10
With R . . . . .	11
Alternative approaches . . . . .	13
Normalising counts . . . . .	14
With RPKMs . . . . .	14
With DESeq . . . . .	14
Differential gene expression . . . . .	15
Count normalisation . . . . .	15
Dispersion estimation . . . . .	16
Differential expression test . . . . .	16

<b>Extra information</b>	<b>17</b>
Software requirements . . . . .	17
Other resources . . . . .	17
Course data . . . . .	17
Tutorials . . . . .	17
Cheat sheets . . . . .	18
<b>Aknowledgments</b>	<b>18</b>

## Introduction

This tutorial will illustrate how to use standalone tools, together with R and Bioconductor for the analysis of RNA-seq data. Keep in mind that this is a rapidly evolving field and that this document is not intended as a review of the many tools available to perform each step; instead, we will cover one of the many existing workflows to analyse this type of data.

We will be working with a subset of a publicly available dataset from *Drosophila melanogaster*, which is available both in the Short Read archive (SRP001537<sup>1</sup> - raw data) and in Bioconductor (pasilla package<sup>2</sup> - processed data). For more information about this dataset please refer to the original publication (Brooks et al. 2010<sup>3</sup>).

The tools and R packages that we will be using during the practical are listed below (see Software requirements<sup>4</sup>) and the necessary data files can be found here<sup>5</sup>. After downloading and uncompressing the `tar.gz` file, you should have the following directory structure in your computer:

```

RNAseq
|-- reference          # reference info (e.g. genome sequence and annotation)
`-- data
    |-- raw            # raw data: fastq files
    |-- mapped         # mapped data: BAM files
    `-- demultiplexing # extra fastq files for the demultiplexing section

```

## Dealing with raw data

### The FASTQ format

The nucleotide sequences and qualities of the short reads produced in a sequencing experiment are commonly stored in a plain text file using the FASTQ format. In the `data/raw` directory, you will find two fastq files, which contain information about the short reads obtained from one of the samples in the *Drosophila melanogaster* experiment.

**Exercise:** Why do we have two fastq files for this given sample? Solution<sup>6</sup>

<sup>1</sup><http://www.ebi.ac.uk/ena/data/view/SRP001537>

<sup>2</sup><http://www.bioconductor.org/packages/release/data/experiment/html/pasilla.html>

<sup>3</sup><http://genome.cshlp.org/content/early/2010/10/04/gr.108662.110>

<sup>4</sup><https://github.com/mgonzalezporta/TeachingMaterial#software-requirements>

<sup>5</sup><http://www.ebi.ac.uk/~mar/courses/RNAseq.tar.gz>

<sup>6</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_fastq\\_ex1.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_fastq_ex1.md)

To confirm that we are working with a fastq file and to get an idea of how this format looks like we can print the first lines of our files by typing this into the terminal:

```
zcat SRR031714_1.fastq.gz | head
zcat SRR031714_2.fastq.gz | head
```

**Exercise:** How many lines are used to represent a read in the fastq file? Which information do they contain? Solution<sup>7</sup>

**Exercise:** How many reads are there in each file? Do both files contain the same number of reads? Is that what we would expect? Solution<sup>8</sup>

## Quality assessment (QA)

We will be using FastQC<sup>9</sup> to generate our first QA report. This software can be executed in two different modes: either using the graphical user interface (if we just type `fastqc` on the terminal) or as a command itself (if we add extra parameters). For example, we can print the help documentation by typing the following: `bash fastqc -h`

To generate a report for our files we only have to provide the file names as an argument: `bash # might take a while fastqc SRR031714_1.fastq.gz SRR031714_2.fastq.gz`

As a result we will obtain the file `filename_fastqc.zip`, which will be automatically unzipped in the `filename_fastqc` directory. There we will find the QA report (`fastqc_report.html`), which provides summary statistics about the numbers of reads, base calls and qualities, as well as other information (you will find a detailed explanation of all the plots in the report in the project website<sup>10</sup>).

**Exercise:** The information provided by the QA report will be very useful when deciding on the options we want to use in the filtering step. After checking it, can you come up with some criteria for the filtering of our file (i.e. keeping/discarding reads based on a specific quality threshold)? Solution<sup>11</sup>

**Exercise:** As we have seen in the previous section, fastq files contain information on the quality of the read sequence. The reliability of each nucleotide in the read is measured using the Phred quality score, which represents the probability of an incorrect base call:

$$Q = -10 \cdot \log_{10} P$$

Figure 1: Phred quality score formula

where  $Q$  is the Phred quality value and  $P$  the probability of error. For example, a Phred quality score of 20 would indicate a probability of error in the base call of 1 in 100 (i.e. 99% accuracy). If you inspect the fastq file again though, you will see that this information is not displayed in number format, but is encoded in a set of characters. During the filtering step, we will be using tools that read these characters and transform them into quality values, so we need to be sure first about the encoding format used in our data (either phred 33 or phred 64). Using the information provided in the QA report (under the *per base sequence quality* section) and in the Wikipedia entry for the FASTQ format<sup>12</sup>, can you guess which encoding format was used? Solution<sup>13</sup>

<sup>7</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_fastq\\_ex2.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_fastq_ex2.md)

<sup>8</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_fastq\\_ex3.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_fastq_ex3.md)

<sup>9</sup><http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

<sup>10</sup><http://www.bioinformatics.babraham.ac.uk/projects/fastqc/Help/3%20Analysis%20Modules/>

<sup>11</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_qa\\_ex1.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_qa_ex1.md)

<sup>12</sup>[http://en.wikipedia.org/wiki/FASTQ\\_format](http://en.wikipedia.org/wiki/FASTQ_format)

<sup>13</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_qa\\_ex2.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_qa_ex2.md)

**Exercise:** As we have seen, a visual interpretation of the QA report is a very useful practice when dealing with HTS data. However, it becomes a very tedious task if we are working with huge volumes of data (imagine we have 1000 fastq files to inspect!). Thankfully, the developers of FastQC have thought of that. Can you spot any alternative output of this software that we could use in this situation? Solution<sup>14</sup>

## Filtering FASTQ files

After analysing the QA report, one might want to discard some of the reads based on several criteria, such as quality and nucleotide composition. We will use two different tools to perform these filtering steps: PRINSEQ<sup>15</sup> and fastqc-mcf<sup>16</sup>.

PRINSEQ offers a wide range of options for filtering and we can learn about them in the manual:

```
prinseq-lite -h
```

Based on what we have learned from the QA report, we could decide to apply the following filters:

```
zcat SRR031714_1.fastq.gz | prinseq-lite \
    -fastq stdin \
    -out_good SRR031714_1_filt1 \
    -out_bad null \
    -trim_qual_right 30 -min_len 32 \
    -ns_max_p 5 \
    -lc_method dust -lc_threshold 10

zcat SRR031714_2.fastq.gz | prinseq-lite \
    -fastq stdin \
    -out_good SRR031714_2_filt1 \
    -out_bad null \
    -trim_qual_right 30 -min_len 32 \
    -ns_max_p 5 \
    -lc_method dust -lc_threshold 10
```

**Exercise:** Which are the criteria that we are using to discard reads? Solution<sup>17</sup>

**Exercise:** Which extra option should we have had to use if our files had been in phred 64 format? Solution<sup>18</sup>

**Exercise:** After filtering the fastq file it is not a bad idea to obtain a QA report again to decide whether we are happy with the results. Do you think we got rid of the main issues spotted initially? Solution<sup>19</sup>

**Exercise:** Usually it is also useful to keep track of the number of reads available in the fastq file both before and after the filtering step. Can you gather this information from the FastQC reports? Given

---

<sup>14</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_qa\\_ex3.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_qa_ex3.md)

<sup>15</sup><http://prinseq.sourceforge.net/>

<sup>16</sup><https://code.google.com/p/ea-utils/>

<sup>17</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_filtering\\_fastq\\_ex1.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_filtering_fastq_ex1.md)

<sup>18</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_filtering\\_fastq\\_ex2.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_filtering_fastq_ex2.md)

<sup>19</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_filtering\\_fastq\\_ex3.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_filtering_fastq_ex3.md)

that we are working with paired-end data, do you see any limitation? Now imagine we were working with a bigger number of files; can you come up with a more automated way to check that? Solution<sup>20</sup>

The filtering step might become complicated if you are working with paired-end data, since you have to be sure that both fastq files (one for each read pair) contain the same number of reads in the same order. There are some tools available that simplify this task, for example, fastqc-mcf. We can print a list of the available options just by typing in the name of the tool: `bash fastqc-mcf`

We observe that a main functionality of fastqc-mcf is to remove adapters from the fastq file. For this reason we need to provide a fasta file with the adapter sequences. In our case, we have decided to check only for the standard Illumina paired-end adapter: `bash cat adapters.fa`

Now that we have a good understanding of the required input we can proceed to execute the tool, trying to match the options that we used previously with PRINSEQ: `bash fastqc-mcf adapters.fa SRR031714_1.fastq.gz SRR031714_2.fastq.gz \ -o SRR031714_1_filt2.fastq -o SRR031714_2_filt2.fastq \ -q 30 -P 33 -l 32 --max-ns 1`

**Exercise:** How does the QA report look like this time? Do we have the same number of reads in each file? Solution<sup>21</sup>

**Exercise:** Something else that one might want to check is the read length. How long were our reads before we started with the filtering step? How long are they now? Solution<sup>22</sup>

Depending on the filtering options used, we might end up with a set of reads with different lengths. *A priori*, this should not be a limitation, but we might encounter some difficulties in the downstream analyses, depending on the tools we want to use. For this reason, if we have a clear idea of what we want to do with the data, it is always a good practice to check the requirements for the tools that we are planning to use before taking any steps. If that is not the case, in order to be on the safe side, we can use filtering options that affect all reads equally, eg:

```
zcat SRR031714_1.fastq.gz | prinseq-lite \
  -fastq stdin \
  -out_good SRR031714_1_filt3 \
  -out_bad null \
  -trim_right 5
```

```
zcat SRR031714_1.fastq.gz | prinseq-lite \
  -fastq stdin \
  -out_good SRR031714_1_filt3 \
  -out_bad null \
  -trim_right 5
```

**Exercise:** Let us generate the QA reports one last time. How do they compare to the previous ones? Solution<sup>23</sup>

In conclusion, there are many filtering combinations that you can apply, and the specific options will largely depend on the type of data and the posterior analyses. We recommend to check the PRINSEQ manual<sup>24</sup> for a nice overview on the topic.

---

<sup>20</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_filtering\\_fastq\\_ex4.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_filtering_fastq_ex4.md)

<sup>21</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_filtering\\_fastq\\_ex5.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_filtering_fastq_ex5.md)

<sup>22</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_filtering\\_fastq\\_ex6.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_filtering_fastq_ex6.md)

<sup>23</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_filtering\\_fastq\\_ex7.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_filtering_fastq_ex7.md)

<sup>24</sup><http://prinseq.sourceforge.net/manual.html>

## De-multiplexing samples

Nowadays, NGS machines produce so many reads that the coverage obtained per lane for the transcriptome of organisms with small genomes is very high. Sometimes it is more valuable to sequence more samples with lower coverage than sequencing only one with very high coverage. With this end, multiplexing techniques have been optimised to sequence several samples in a single lane using 4-6 bp barcodes to uniquely identify the sample within the library (e.g. Lefrançois et al. 2009<sup>25</sup>). This approach is very advantageous for researchers, especially in terms of cost, but it adds an additional layer of pre-processing that is not as trivial as one would think, given the average 0.1-1% sequencing error rate that introduces a lot of multiplicity in the actual barcodes. Most commonly the data is de-multiplexed immediately after sequencing, and only FASTQ files that are ready for analyses are distributed. However, you might encounter the necessity to perform the de-multiplexing step yourself, or, given de-multiplexed FASTQ files, to remove the adaptors manually; thus, it is important to learn how to deal with such data.

The data which we were working on in the previous section was not multiplexed, so we will now work with a different fastq file that can be found under the `data/demultiplexing` directory. In this directory you will also find information on the barcodes used:

```
cat barcodes.txt
```

**Exercise:** Imagine we do not know whether the barcode was introduced in the 5' or the 3' end of the reads. How can we figure that out? Solution<sup>26</sup>

In order to separate the reads in 4 different fastq files (one for each barcode/sample) we will use `fastq-multx`<sup>27</sup>. We can learn more about this tool by typing its name in the terminal:

```
fastq-multx
```

Although we already know where our barcodes are located within the read, from the documentation we observe that `fastq-multx` will attempt to guess the position for us. Let us try the automatic guessing with the following command:

```
fastq-multx barcodes.txt barcoded.fastq -o %.barcoded.fastq
```

After executing the command above you should have five new fastq files: one corresponding to each sample and one for the reads that did not match any of the barcodes.

**Exercise:** Try generating a QA report for one of the samples. How does it compare to the report for the initial multiplexed fastq file? What happened to the read length? Solution<sup>28</sup>

## Aligning reads to the genome

So far we have been working with fastq files, which contain the reads that were generated during the sequencing experiment. *A priori* we do not know from which transcripts those reads were originated, and that is precisely what will be addressed in following steps, starting with the mapping. There are essentially two ways of approaching this: one is to align the reads to a known transcriptome or

---

<sup>25</sup><http://www.biomedcentral.com/1471-2164/10/37>

<sup>26</sup>[../solutions/\\_demultiplexing\\_ex1.md](#)

<sup>27</sup><https://code.google.com/p/ea-utils/>

<sup>28</sup>[../solutions/\\_demultiplexing\\_ex2.md](#)

genome, and the other is to assemble these reads *de novo* into a transcriptome without the need for any reference.

**Exercise:** Can you think of any advantages/disadvantages for the above mentioned approaches?  
Solution<sup>29</sup>

Here we decided to align our reads to a known reference genome. To achieve this task, you could use any aligner capable of exporting its results in the SAM/BAM format or in a format that can be easily converted to this one. In the case of RNA-seq data, we also want to be able to retain information about split reads (i.e. reads with a gap) and spliced reads (i.e. those that span multiple exons). There are many aligners available, some of them optimised for working with RNA-seq data (e.g. TopHat, GSNAP... - see Fonseca et al. 2012<sup>30</sup> for a review). In this practical we decided to use TopHat<sup>31</sup>, and these are the commands we would use to map the fastq files that we have generated during the filtering step:

```
# do not run - very time consuming
# output already provided in the data/mapped directory

cd reference

# obtain the reference genome from Ensembl
wget ftp://ftp.ensembl.org/pub/release-62/fasta/drosophila_melanogaster/dna/\
    Drosophila_melanogaster.BDGP5.25.62.dna_rm.toplevel.fa.gz
gunzip Drosophila_melanogaster.BDGP5.25.62.dna_rm.toplevel.fa.gz

# index the reference genome
bowtie-build Drosophila_melanogaster.BDGP5.25.62.dna_rm.toplevel.fa \
    Drosophila_melanogaster.BDGP5.25.62.dna_rm.toplevel

# align the reads
tophat Drosophila_melanogaster.BDGP5.25.62.dna_rm.toplevel \
    ../data/raw/SRR031714_1_filt3.fastq \
    ../data/raw/SRR031714_2_filt3.fastq
```

The last command runs TopHat with the default options. A detailed description of those, as well as information on other commands (e.g. for single-end reads), can be found in the manual<sup>32</sup> or just by typing the name of the tool in the console (i.e. type `tophat` on its own).

Notice the `reference` directory, which contains, amongst other files, the genomic sequence for *Drosophila melanogaster* (`Drosophila_melanogaster.BDGP5.25.62.dna_rm.toplevel.fa`). We can inspect which chromosomes are present in this fasta file with the following command:

```
grep '^>' Drosophila_melanogaster.BDGP5.25.62.dna_rm.toplevel.fa | head
```

**Exercise:** Suppose we had decided to align to the transcriptome instead. Similarly to what we did with the genome sequence, the transcriptome sequence for *Drosophila melanogaster* can be obtained from Ensembl with the following command:

---

<sup>29</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_aligning\\_ex1.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_aligning_ex1.md)

<sup>30</sup><http://bioinformatics.oxfordjournals.org/content/28/24/3169>

<sup>31</sup><http://tophat.cbcb.umd.edu/>

<sup>32</sup><http://tophat.cbcb.umd.edu/manual.html>

```
# do not run
wget ftp://ftp.ensembl.org/pub/release-62/fasta/drosophila_melanogaster/cdna/\
Drosophila_melanogaster.BDGP5.25.62.cdna.all.fa.gz
gunzip Drosophila_melanogaster.BDGP5.25.62.cdna.all.fa.gz
```

This file has been provided on the **reference** directory. What do the “chromosome” names correspond to in this case? Solution<sup>33</sup>

We are now familiarised with the input required to align reads to a reference genome or transcriptome. In both cases, the output produced by the mapping tool is going to be stored in SAM/BAM format, which we will inspect in the next section.

## Dealing with aligned data

### The SAM/BAM format

The SAM/BAM format is the standard way of representing the results from the alignment step. It contains the same information as in the fastq file, plus some extra fields providing mapping information, for example, the coordinates where each of the reads was aligned. A SAM file is a plain text file with the information spread across different columns, and a BAM file is just its compressed version in binary format. In order to save disk space, we will typically work with BAM files; however, we can easily transform a BAM file into SAM format using samtools<sup>34</sup>:

```
# do not run
# output already provided in data/mapped
samtools view -h -o untreated3.sam untreated3.bam
```

We can now inspect the first lines of the file with standard Unix commands:

```
head -n20 untreated3.sam
```

Alternatively, we can directly inspect the contents of a BAM file with the following samtools command:

```
samtools view untreated3.bam | head
```

**Exercise:** Why do we get a different output from the two previous commands? How can we obtain information about the header from the BAM file? Hint: try typing **samtools view** into the terminal. Solution<sup>35</sup>

**Exercise:** The first column in the BAM file contains the read name. Take a closer look at the first alignments: why do you think some of the names appear twice, while others seem to be present only once? Solution<sup>36</sup>

**Exercise:** A description of the SAM format can be found in the samtools website<sup>37</sup>, under the section *SAM format*. With the combination of samtools and Unix commands, try to answer the following questions:

<sup>33</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_aligning\\_ex2.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_aligning_ex2.md)

<sup>34</sup><http://samtools.sourceforge.net/samtools.shtml>

<sup>35</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_bam\\_ex1.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_bam_ex1.md)

<sup>36</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_bam\\_ex2.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_bam_ex2.md)

<sup>37</sup><http://samtools.sourceforge.net/samtools.shtml>



- How many reads are mapped in total?
- How many reads map to each chromosome?
- How many different mapping qualities are represented in the BAM file, and how many reads have each of them assigned?
- How many different alignment flags can you find in the BAM file? What do they represent? *Hint:* <http://picard.sourceforge.net/explain-flags.html>
- Try to print the unique CIGAR strings for the first 300 reads. What is their meaning? *Hint:* <http://genome.sph.umich.edu/wiki/SAM>

Solution<sup>38</sup>

## Visualising aligned reads

Several genome browsers exist to visualise the files generated during the analysis of high-throughput sequencing data, including BAM files. Two of the most popular tools are IGV<sup>39</sup> and Tablet<sup>40</sup>. In this practical we will be using IGV to visualise our BAM file. We can launch this tool from the *Download* section in the project website<sup>41</sup>.

Once the interface is loaded, we can proceed to load the necessary files. In our case, we will load the following information:

- the reference genome:

```
Genomes > Load genome from file >
    reference/Drosophila_melanogaster.BDGP5.25.62.dna_rm.toplevel.fa
```

- the BAM file:

```
File > Load from file > data/mapped/untreated3.bam

# IGV requires the BAM file to be indexed, which can be achieved with samtools
# (i.e. `samtools index bam_file`)
# For this practical the index is already provided, so there is no need to run this command
```

- the annotation:

```
File > Load from file > reference/Drosophila_melanogaster.BDGP5.25.62.gtf
```

**Exercise:** Spend some time exploring the loaded BAM file and how the reads overlap with the known annotation. Can you find examples of split and spliced reads? For a subset of the reads, some nucleotides are highlighted in a different color. What do you think the explanation for this is? *Hint:* <http://www.broadinstitute.org/igv/AlignmentData>

Solution<sup>42</sup>

<sup>38</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_bam\\_ex3.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_bam_ex3.md)

<sup>39</sup><http://www.broadinstitute.org/igv/>

<sup>40</sup><http://bioinf.scri.ac.uk/tablet/>

<sup>41</sup><http://www.broadinstitute.org/igv/download>

<sup>42</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_visualising\\_ex1.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_visualising_ex1.md)

## Filtering BAM files

Samtools can also be used to further modify and/or subset BAM files. For example, some tools will require that the reads are sorted by coordinate or by name. In addition, and similarly to what we did with the fastq files, one might consider to discard reads with a low alignment quality (including reads that align to several locations in the genome), or in the case of paired-end data, discard reads that are not properly paired.

**Exercise:** Try to answer to the following questions using samtools and the information provided in the documentation<sup>43</sup>:

- How many reads are properly paired?
- By default, TopHat creates BAM files where the reads are sorted by coordinate. How would you sort the properly paired reads by name instead? Save the output in a new BAM file called `untreated3_paired.bam`, which we will use later on during the practical.
- Which percentage of those properly paired reads map uniquely? *Hint:* Have a look at the options for `samtools sort` and `samtools view`.

Solution<sup>44</sup>

## Counting reads overlapping annotated genes

Following the read mapping step, we can proceed working with BAM files with standalone tools or load them directly in R. These two workflows are not exclusive and we will cover both of them for illustrative purposes.

### With htseq-count

htseq-count<sup>45</sup> is a simple but yet powerful tool to overlap a BAM file with the genome annotation and thus obtain the number of reads that overlap with our features of interest. As usual, we can obtain information on the tool by typing `htseq-count -h` and by referring to its website.

**Exercise:** One of the input files required by htseq-count is a GTF file. For this practical, you will find this file under the directory `reference`. Which information does it contain? Solution<sup>46</sup>

**Exercise:** As we have already mentioned, the other required input file is a BAM file. Can you spot any specific requirement regarding this file? *Hint*<sup>47</sup> - Solution<sup>48</sup>

In addition to the input file requirements, special care must be taken in dealing with reads that overlap more than one feature (e.g. overlapping genes), and thus might be counted several times in different features. To deal with this, htseq-count offers three different counting modes: union, intersection-strict and intersection-nonempty.

**Exercise:** What are the differences between these three counting modes? *Hint*<sup>49</sup> - Solution<sup>50</sup>

Now that we have a good understanding of the input files and options, we can proceed to execute htseq-count:

---

<sup>43</sup><http://samtools.sourceforge.net/samtools.shtml>

<sup>44</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_filtering\\_bam.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_filtering_bam.md)

<sup>45</sup><http://www-huber.embl.de/users/anders/HTSeq/doc/count.html>

<sup>46</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_counting\\_ex1.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_counting_ex1.md)

<sup>47</sup><http://www-huber.embl.de/users/anders/HTSeq/doc/count.html>

<sup>48</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_counting\\_ex2.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_counting_ex2.md)

<sup>49</sup><http://www-huber.embl.de/users/anders/HTSeq/doc/count.html>

<sup>50</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_counting\\_ex3.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_counting_ex3.md)

```
# the following command takes a while to execute
# the output file is already provided in the data/mapped directory
samtools view untreated3_paired.bam | htseq-count \
    --mode=intersection-nonempty \
    --stranded=no \
    --type=exon \
    --idattr=gene_id - \
    ../../reference/Drosophila_melanogaster.BDGP5.25.62.gtf > htseq_count.out
```

**Exercise:** In addition to the counts that overlap known genes, the output file also contains some extra information on reads that could not be assigned to any of those; can you find it? Solution<sup>51</sup>

## With R

Computing gene counts in R is very similar to what we have done so far with htseq-count. However, it requires some extra steps, since we first need to load the necessary files (i.e. BAM files and annotation).

*Note: All the commands provided in this section have to be executed in R. Make sure to specify the working directory properly before starting (e.g. `setwd("./data")`).*

**Importing BAM files** There are three main functions to load BAM files into R:

- *scanBam*: this function is part of the *Rsamtools* package and is the low level function used by the other two. It potentially reads *all* fields (including CIGAR strings and user defined tags) of a BAM file into a list structure, but allows you to select specific fields and records to import.
- *readAligned*: a higher-level function defined in the *ShortRead* package which imports some of the data (query names, sequences, quality, strand, reference name, position, mapping quality and flag) into an *AlignedRead* object. *ShortRead* was the first package developed to read in NGS data and is able to read almost sequencer every manufacturer proprietary formats, so you could for example also use it to read an Illumina export file produced by a GenomeAnalyzer GAIIX.
- *readGappedAlignments* and *readGappedAlignmentPairs*: two functions from the *GenomicRanges* package that create an object intended for operations such as searching for overlaps or coverage. Each alignment is described by its position and strand on the reference and read ids, sequences and base qualities are discarded for the sake of memory usage and speed.

In this section we will import the data using the *readGappedAlignmentPairs* function, intended for paired-end data. In order to speed up the process of importing the data, we will use the function *ScanBamParam* to load only the reads that map to chromosome 4:

```
library(GenomicRanges)
library(Rsamtools)

# define a filter
which=RangesList(IRanges(1, 1351857))
names(which)="chr4"
which
param=ScanBamParam(which=which)
```

---

<sup>51</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_counting\\_ex4.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_counting_ex4.md)

```
# import the data
aln_chr4=readGappedAlignmentPairs("untreated3.bam", use.names=T, param=param)
aln_chr4
```

We have now stored our data in an object of the class *GappedAlignmentPairs*, which has been defined in the *GenomicRanges* package and does not correspond to the standard R classes. For this reason, it is useful to check the documentation for this package<sup>52</sup> in order to learn how to access our data.

**Exercise:** After having a look at the *GenomicRanges* documentation, try to answer the following questions:

- How many reads have been loaded?
- How can we access the read names? What about the strand information?
- How can we access the information for the first reads in the pair? Try to print a vector with their start coordinates.
- How many reads are properly paired?
- What is the percentage of reads that map to multiple locations?
- What information does the command `seqlevels(aln_chr4)` provide? *Hint:* look for the *GappedAlignmentPairs* class

Solution<sup>53</sup>

Since we have loaded only the reads that map to chromosome 4, we can proceed to modify the `aln_chr4` object accordingly:

```
seqlevels(aln_chr4)="chr4"
aln_chr4
```

**Importing the annotation** To link the alignments to their respective features, we need access to the genome annotation for the studied organism, in our case *Drosophila melanogaster*, which contains information on the coordinates of known exons, genes and transcripts. Similarly to what we encountered when loading BAM files, there is more than one way to load the annotation in R (see the Bioconductor resources<sup>54</sup> for further details). It is extremely important to pay attention to overlapping features (e.g. exons shared by multiple transcripts within the same gene), since they might end up complicating the downstream analysis (e.g. we need to make sure not to count the same read multiple times). In order to circumvent this limitation, in this practical we will use the *biomaRt* package to query Ensembl directly from R and retrieve only the necessary information:

```
library(biomaRt)

ensembl62=useMart(host="apr2011.archive.ensembl.org",
  biomart="ENSEMBL_MART_ENSEMBL",
  dataset="dmelanogaster_gene_ensembl")

fields=c("chromosome_name",
  "strand",
  "ensembl_gene_id",
```

<sup>52</sup><http://www.bioconductor.org/packages/release/bioc/manuals/GenomicRanges/man/GenomicRanges.pdf>

<sup>53</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_counting\\_ex5.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_counting_ex5.md)

<sup>54</sup><http://www.bioconductor.org/help/course-materials/>

```

"ensembl_exon_id",
"start_position",
"end_position",
"exon_chrom_start",
"exon_chrom_end")
annot=getBM( fields, mart=ensembl62)

```

**Exercise:** Have a look at the newly created `annot` object. What type of object is it? *Hint:* use the function `class`. Solution<sup>55</sup>

**Exercise:** In the next subsection we will calculate the overlap between the loaded BAM file and the annotation with the function `summarizeOverlaps` from the *GenomicRanges* package. What is the input required? Do we have all the necessary objects ready? *Hint:* type `?summarizeOverlaps`. Solution<sup>56</sup>

Before we proceed to calculate the counts, we need to store the annotation information in an object of the proper class: `rconsole annot=GRanges( seqnames = Rle(paste("chr", annot$chromosome_name, sep="")), ranges = IRanges(start=annot$exon_chrom_start, end=annot$exon_chrom_end), strand = Rle(annot$strand), exon=annot$ensembl_exon_id, gene=annot$ensembl_gene_id )`  
`annot class(annot)`

**Counting reads over known genes in R** Now that we have the alignment locations (`aln_unique` object) and the genome annotation (`annot` object), we can quantify gene expression by counting reads over all exons of a gene and summing them together. Similarly to what we encountered with `htseq-count`, we need to pay attention to those reads that overlap with several features.

```

counts_chr4=summarizeOverlaps(
  annot_chr4, aln_chr4, ignore.strand=T, mode="IntersectionNotEmpty")
exon_counts_chr4=assays(counts_chr4)$counts[,1]
names(exon_counts_chr4)=elementMetadata(annot_chr4)$gene

head(exon_counts_chr4, n=15)

```

**Exercise:** So far we have obtained the number of reads overlapping each exon. How can we combine this information to obtain gene counts? *Hint:* use the functions `split` and `sapply`. Solution<sup>57</sup>

## Alternative approaches

In this section of the practical we have seen how to calculate the number of reads that overlap known gene models. In the two approaches evaluated here, those reads that mapped to multiple features were not considered. This is a simplification we may not want to pursue, and alternatively, there are several methods to probabilistically estimate the expression of overlapping features (e.g. Trapnell et al. 2010<sup>58</sup>, Turró et al. 2011<sup>59</sup>, Li et al. 2010<sup>60</sup>...).

<sup>55</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_counting\\_ex6.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_counting_ex6.md)

<sup>56</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_counting\\_ex7.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_counting_ex7.md)

<sup>57</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_counting\\_ex8.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_counting_ex8.md)

<sup>58</sup><http://www.nature.com/nbt/journal/v28/n5/abs/nbt.1621.html>

<sup>59</sup><http://genomebiology.com/content/12/2/R13>

<sup>60</sup><http://bioinformatics.oxfordjournals.org/content/26/4/493.long>

## Normalising counts

### With RPKMs

While in the previous sections the data was derived from a single sample, in this exercise we will work with the precomputed counts for all the samples in our experiment<sup>61</sup>:

```
library("pasilla")

data("pasillaGenes")
counts=counts(pasillaGenes)
head(counts)
```

A common way to normalise reads is to convert them to RPKMs (or FPKMs in the case of paired-end data). This implies normalising the read counts depending on the feature size (exon, transcript, gene model...) and on the total number of reads sequenced for that library:

$$RPKMs = \frac{\text{gene counts}}{\text{gene length} \cdot \text{library size}} \cdot 10^9$$

Figure 2: RPKM formula

**Exercise:** Let us obtain RPKMs for the table `counts` following these steps:

- Calculate the length of the exons in the `annot` object and store the result in a vector, with the name of each element set to the corresponding gene. *Hint:* check the `width` and `elementMetadata` accessors.
- Obtain gene lengths by adding up the lengths of all the exons in each gene. *Hint:* check the functions `split` and `sapply`.
- Normalise the counts by the library size. *Hint:* check the function `colSums`.
- Continue normalising by gene length, but be aware that the object that contains the gene lengths and the one that contains the normalised counts might have a different number of genes.
- Finally, obtain RPKMs by multiplying by a factor of  $10^9$ .

Solution<sup>62</sup>

Such a count normalisation is suited for visualisation, but sub-optimal for further analyses. A better way of normalising the data is to use either the `edgeR` or `DESeq` packages.

### With DESeq

RPKM normalisation is not the most adequate for certain types of downstream analysis (e.g. differential gene expression), given that it is susceptible to library composition biases. There are many other normalisation methods that should be considered with that goal in mind (see Dillies et al. 2012<sup>63</sup> for a comparison). In this section we are going to explore the one offered within the `DESeq` package:

<sup>61</sup><http://bioconductor.org/packages/2.11/data/experiment/html/pasilla.html>

<sup>62</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_normalising\\_ex1.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_normalising_ex1.md)

<sup>63</sup><http://bib.oxfordjournals.org/content/early/2012/09/15/bib.bbs046.long>

```
library(DESeq)

# create an object of class CountDataSet, which is the data container used by DESeq
# remember we have already loaded the pasilla dataset in the previous section
cds = newCountDataSet(counts(pasillaGenes), phenoData(pasillaGenes)$condition)
cds

# the CountDataSet class is a container for high-throughput assays and experimental metadata
# the counts can be accessed with the counts function
head(counts(cds))
# and the phenotypic data with the pData function
pData(cds)
```

In order to normalise the raw counts we will start by determining the relative library sizes, or *size factors* of each library. For example, if the counts of the expressed genes in one sample are, on average, twice as high as in another, the size factor for the first sample should be twice as large as the one for the other sample. These size factors can be obtained with the function *estimateSizeFactors*:

```
cds=estimateSizeFactors(cds)
sizeFactors(cds)
```

Once we have this information, the normalised data is obtained by dividing each column of the count table by the corresponding size factor. We can perform this calculation by calling the function *counts* with a specific argument as follows:

```
deseq_ncounts=counts(cds, normalized=TRUE)
```

**Exercise:** We have now accumulated three different versions of the same dataset: the raw counts (*counts*), the RPKM quantifications (*rpkm*) and the DESeq normalised counts (*deseq\_ncounts*). How would you visualise the performance of each normalisation method in getting rid of the variation that does not associate to the experimental conditions that are being evaluated? Solution<sup>64</sup>

## Differential gene expression

A basic task in the analysis of expression data is the detection of differentially expressed genes. The *DESeq* package provides a method to test for differential expression by using the negative binomial distribution. It expects a matrix of count values where each column corresponds to a sample and each line to a feature (e.g. a gene). Typically, a *DESeq* analysis is performed in three steps: count normalisation, dispersion estimation and differential expression test.

### Count normalisation

Since we have already generated a matrix with the normalised counts in the previous section, we will use it directly as input for the next step.

---

<sup>64</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_normalising\\_ex2.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_normalising_ex2.md)

## Dispersion estimation

An important step in differential expression analysis is to figure out how much variability we can expect in the expression measurements within the same condition. Unless this is known, we cannot make inferences about whether the change in expression observed for a given gene is big enough to be considered significant, or whether it corresponds to the variability that we would expect by chance. This is why it is so important to have replicates: they show how much variation occurs without a difference in the condition.

In *DESeq*, in order to estimate the dispersion for each gene, we can use the function *estimateDispersions*:

```
cds=estimateDispersions(cds)
str(fitInfo(cds))
head(fData(cds))
```

The list returned by the *fitInfo* function contains the empirical dispersion values for each gene, the curve fitted through the dispersions, and the fitted values. This can be visualised by plotting the per-gene estimates against the normalised mean expressions per gene, and overlaying the fitted curve:

```
plot(rowMeans(counts(cds, normalized=TRUE)), fitInfo(cds)$perGeneDispEsts, log="xy")
xg=10^seq(-.5, 5, length.out=300)
lines(xg, fitInfo(cds)$dispFun(xg), col="red")
```

## Differential expression test

Finally, we will use the function *nbinomTest* to contrast the two studied conditions:

```
de=nbinomTest(cds, "treated", "untreated")
head(de)
```

The *padj* column in the table **de** contains the p-values adjusted for multiple testing with the Benjamini-Hochberg procedure (i.e. FDR). This is the information that we will use to decide whether the expression of a given gene differs significantly across conditions (e.g. we can arbitrarily decide that genes with an  $FDR < 0.10$  are differentially expressed).

**Exercise:** How would you select those genes that pass a given FDR threshold (e.g.  $FDR < 0.10$ )? Which are the most significant? Solution<sup>65</sup>

**Exercise:** Let us generate an MA plot to evaluate the results of the differential expression analysis. Follow this steps:

- Obtain the mean expression values from the **de** object.
- Obtain the log2 fold change values from the **de** object.
- Create a scatterplot with the mean in the *x* axis and the log2 fold change in the *y* axis.
- How would we highlight those genes that are differentially expressed with a different color? *Hint:* check the *col* argument for the *plot* function and the function *ifelse*.

Solution<sup>66</sup>

---

<sup>65</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_de\\_ex1.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_de_ex1.md)

<sup>66</sup>[https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/\\_de\\_ex2.md](https://github.com/mgonzalezporta/TeachingMaterial/blob/master/solutions/_de_ex2.md)



## Extra information

### Software requirements

*Note: depending on the topics covered in the course some of these tools might not be used.*

- Standalone tools:
  - FastQC<sup>67</sup>
  - PRINSEQ<sup>68</sup>
  - eautils<sup>69</sup>
  - samtools<sup>70</sup>
  - IGV<sup>71</sup>
  - htseq-count<sup>72</sup>
- Bioconductor packages:
  - GenomicRanges<sup>73</sup>
  - Rsamtools<sup>74</sup>
  - biomaRt<sup>75</sup>
  - pasilla<sup>76</sup>
  - DESeq<sup>77</sup>

### Other resources

#### Course data

- Complete course data, including command outputs and R sessions<sup>78</sup>

#### Tutorials

- Course materials available at the Bioconductor website<sup>79</sup>
- Online training resources at the EBI website<sup>80</sup>
- R and Bioconductor tutorial by Thomas Girke<sup>81</sup>
- Do not forget to check the documentation for the packages used in the practical!

---

<sup>67</sup><http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

<sup>68</sup><http://prinseq.sourceforge.net/>

<sup>69</sup><https://code.google.com/p/ea-utils/>

<sup>70</sup><http://sourceforge.net/projects/samtools/>

<sup>71</sup><http://www.broadinstitute.org/software/igv/download>

<sup>72</sup><http://www-huber.embl.de/users/anders/HTSeq/doc/count.html>

<sup>73</sup><http://www.bioconductor.org/packages/release/bioc/html/GenomicRanges.html>

<sup>74</sup><http://www.bioconductor.org/packages/release/bioc/html/Rsamtools.html>

<sup>75</sup><http://www.bioconductor.org/packages/release/bioc/html/biomaRt.html>

<sup>76</sup><http://www.bioconductor.org/packages/release/data/experiment/html/pasilla.html>

<sup>77</sup><http://www.bioconductor.org/packages/release/bioc/html/DESeq.html>

<sup>78</sup>[http://www.ebi.ac.uk/~mar/courses/RNAseq\\_all.tar.gz](http://www.ebi.ac.uk/~mar/courses/RNAseq_all.tar.gz)

<sup>79</sup><http://www.bioconductor.org/help/course-materials/>

<sup>80</sup>[http://www.ebi.ac.uk/training/online/course-list?topic%5B%5D=13&views\\_\\_exposed\\_\\_form\\_\\_focused\\_\\_field=](http://www.ebi.ac.uk/training/online/course-list?topic%5B%5D=13&views__exposed__form__focused__field=)

<sup>81</sup>[http://manuals.bioinformatics.ucr.edu/home/R\\_BioCondManual](http://manuals.bioinformatics.ucr.edu/home/R_BioCondManual)

## Cheat sheets

- R reference card<sup>82</sup>
- Unix comand line cheat sheet<sup>83</sup>

## Aknowledgments

This tutorial has been inspired on material developed by Ângela Gonçalves, Nicolas Delhomme, Simon Anders and Martin Morgan, who I would like to thank and acknowledge. Special thanks must go to Ângela Gonçalves, with whom I started teaching, and Gabriella Rustici, for always finding a way to organise a new course.

---

<sup>82</sup><http://cran.r-project.org/doc/contrib/Short-refcard.pdf>

<sup>83</sup>[http://sites.tufts.edu/cbi/files/2013/01/linux\\_cheat\\_sheet.pdf](http://sites.tufts.edu/cbi/files/2013/01/linux_cheat_sheet.pdf)