

# BCGES short courses, session 7, transcriptome sequencing (RNA-Seq)

Vincent Plagnol

## Contents

<b>1</b>	<b>A last (non-RNA-Seq) point about VCF format and R</b>	<b>2</b>
<b>2</b>	<b>GTF format to store information gene-centric information (20 minutes)</b>	<b>2</b>
2.1	Ensembl data . . . . .	2
2.2	UCSC data . . . . .	2
2.3	Why aren't my reads aligning? . . . . .	2
<b>3</b>	<b>Library normalization choices (40 minutes)</b>	<b>6</b>
3.1	A toy example . . . . .	6
3.2	Further reading on normalization of RNA-Seq data . . . . .	7
<b>4</b>	<b>Aligning RNA-Seq data and estimating gene expression levels (45 minutes)</b>	<b>8</b>
4.1	Aligning with tophat and bowtie . . . . .	8
4.2	Expression level estimation using Cufflinks . . . . .	8
4.3	What if we do not have a GTF file? . . . . .	9
<b>5</b>	<b>Differential expression analysis (30 minutes)</b>	<b>10</b>
5.1	Using DESeq . . . . .	10
5.2	Using DESeq2 . . . . .	11
<b>6</b>	<b>Galaxy server</b>	<b>12</b>

One thing before we start: please unzip the file `chr12_short.fa.gz` in the `data/RNASeq` folder and go back to session 7.

## 1 A last (non-RNA-Seq) point about VCF format and R

R is actually much better than I anticipated at reading VCF files, and I am currently learning what is available. I came across the excellent package `VariantAnnotation`. Installation instructions are available [here](#). I wanted to highlight quickly what this package does with the code below. I find it very elegant and practical, so have a look, this may be quite useful.

```
library(VariantAnnotation)
input.vcf <- "../session4_multi_sample_calling_annotation/results/UBASH31A_samtools.vcf"
my.vcf <- readVcf(input.vcf, "hg19")
genotypes <- geno(my.vcf)
head(genotypes$DP)
head(genotypes$GT)
```

## 2 GTF format to store information gene-centric information (20 minutes)

### 2.1 Ensembl data

If one works with genes and exons, it is important to have a format that captures this information. The file format that does this is the GTF format. A good place to download GTF file is the <http://www.ensembl.org/info/data/ftp/index.html>. One can start by using the `curl` function (which is a combination of `cat` and `url`) to obtain the first few lines of an example GTF file.

```
curl --silent ftp://ftp.ensembl.org/pub/release-76/gtf/homo_sapiens/Homo_sapiens.GRCh38.76.gtf.gz | \
  zcat | head -100 > results/human_gtf_example.gtf
```

**Exercise:** Go over the GTF format and understand what the fields mean, and how the data are organised.

You can now download the full ensembl file to get an idea of the size of the file. We will use the `wget` function that was used before in these practicals (note that the code below is not executed, because too long to go through).

```
wget -O results/ensembl_human_GRCh38.gtf.gz \
  ftp://ftp.ensembl.org/pub/release-76/gtf/homo_sapiens/Homo_sapiens.GRCh38.76.gtf.gz
```

### 2.2 UCSC data

UCSC is the other obvious place to obtain genome-scale data. The webpage you want to become familiar with is [this one](#).

**Exercise:** Look for a human GTF file generally equivalent to the one you just downloaded from UCSC. Compare the sizes of both files, look for differences and similarities.

```
#You want to set the group option \texttt{mRNA and EST}.
#Use the Human mRNA track
#In output format select \texttt{GTF - gene transfer format}
#Specify a name in the output file
#Maybe request a compressed file to limit transfer time
#Compressed the UCSC GTF is ... and the compressed Ensembl one is 16M.
```

### 2.3 Why aren't my reads aligning?

There are many reasons why you may get a very low alignment rate with sequence data, in particular RNA-Seq. **Exercise:** The files below are fastq files for RNA-Seq post-mortem brain data. Can you understand why the aligned proportion is extremely low? This is a complicated problem without a hint so here is one: think about the

role of adapters when DNA fragments are short. Also note that the Illumina forward and reverse adapters start with AGATCGGAAGAG. How would you go about fixing it? There should be a program in your PATH that is designed exactly for this issue.

```
ls ../data/fastq_for_RNASeq/file1.fq ../data/fastq_for_RNASeq/file2.fq
```

```
## ../data/fastq_for_RNASeq/file1.fq
## ../data/fastq_for_RNASeq/file2.fq
```

```
## for reference, this is how it runs on my mac
```

```
python2.7 /Users/vplagnol/Library/Python/2.7/lib/python/site-packages/cutadapt \
  -a AGATCGGAAGAG -A AGATCGGAAGAG \
  -o results/trimmed1.fastq -p results/trimmed2.fastq \
  ../data/fastq_for_RNASeq/file1.fq ../data/fastq_for_RNASeq/file2.fq
```

```
## This is cutadapt 1.8.3 with Python 2.7.10
```

```
## Command line parameters: -a AGATCGGAAGAG -A AGATCGGAAGAG -o results/trimmed1.fastq -p results/trimmed2.
```

```
## Trimming 2 adapters with at most 10.0% errors in paired-end mode ...
```

```
## Finished in 0.01 s (100 us/read; 0.60 M reads/minute).
```

```
##
```

```
## === Summary ===
```

```
##
```

```
## Total read pairs processed: 100
```

```
## Read 1 with adapter: 97 (97.0%)
```

```
## Read 2 with adapter: 93 (93.0%)
```

```
## Pairs written (passing filters): 100 (100.0%)
```

```
##
```

```
## Total basepairs processed: 20,000 bp
```

```
## Read 1: 10,000 bp
```

```
## Read 2: 10,000 bp
```

```
## Total written (filtered): 15,743 bp (78.7%)
```

```
## Read 1: 7,838 bp
```

```
## Read 2: 7,905 bp
```

```
##
```

```
## === First read: Adapter 1 ===
```

```
##
```

```
## Sequence: AGATCGGAAGAG; Type: regular 3'; Length: 12; Trimmed: 97 times.
```

```
##
```

```
## No. of allowed errors:
```

```
## 0-9 bp: 0; 10-12 bp: 1
```

```
##
```

```
## Bases preceding removed adapters:
```

```
## A: 12.4%
```

```
## C: 41.2%
```

```
## G: 32.0%
```

```
## T: 14.4%
```

```
## none/other: 0.0%
```

```
##
```

```
## Overview of removed sequences
```

```
## length count expect max.err error counts
```

```
## 10 1 0.0 1 1
```

```
## 12 1 0.0 1 1
```

```
## 13 4 0.0 1 4
```

```
## 15 4 0.0 1 4
```

```
## 16 9 0.0 1 8 1
```

```
## 17 4 0.0 1 4
```

```
## 18 8 0.0 1 8
```

```
## 19 4 0.0 1 4
```

```
## 20 5 0.0 1 5
```

```
## 21 9 0.0 1 7 2
```

```

## 22 10 0.0 1 9 1
## 23 4 0.0 1 4
## 24 4 0.0 1 4
## 25 6 0.0 1 5 1
## 26 4 0.0 1 4
## 27 5 0.0 1 5
## 28 5 0.0 1 5
## 29 2 0.0 1 2
## 30 2 0.0 1 2
## 34 1 0.0 1 1
## 36 1 0.0 1 1
## 39 1 0.0 1 1
## 41 1 0.0 1 0 1
## 47 1 0.0 1 1
## 55 1 0.0 1 1
##
## === Second read: Adapter 2 ===
##
## Sequence: AGATCGGAAGAG; Type: regular 3'; Length: 12; Trimmed: 93 times.
##
## No. of allowed errors:
## 0-9 bp: 0; 10-12 bp: 1
##
## Bases preceding removed adapters:
##   A: 16.1%
##   C: 38.7%
##   G: 36.6%
##   T:  8.6%
##  none/other: 0.0%
##
## Overview of removed sequences
## length count expect max.err error counts
## 10 1 0.0 1 1
## 12 1 0.0 1 1
## 13 4 0.0 1 3 1
## 15 4 0.0 1 3 1
## 16 8 0.0 1 7 1
## 17 3 0.0 1 3
## 18 8 0.0 1 8
## 19 3 0.0 1 3
## 20 4 0.0 1 3 1
## 21 9 0.0 1 9
## 22 9 0.0 1 9
## 23 4 0.0 1 4
## 24 4 0.0 1 4
## 25 6 0.0 1 5 1
## 26 4 0.0 1 4
## 27 6 0.0 1 6
## 28 5 0.0 1 5
## 29 2 0.0 1 2
## 30 2 0.0 1 2
## 34 1 0.0 1 1
## 36 1 0.0 1 1
## 39 1 0.0 1 1
## 41 1 0.0 1 1
## 47 1 0.0 1 1
## 55 1 0.0 1 0 1

```

```
## and this is what should work on your linux box
cutadapt -a AGATCGGAAGAG -A AGATCGGAAGAG \
  -o results/trimmed1.fastq -p results/trimmed2.fastq \
  ../data/fastq_for_RNASeq/file1.fq ../data/fastq_for_RNASeq/file2.fq
```

## 3 Library normalization choices (40 minutes)

### 3.1 A toy example

A key issue for RNA-Seq data is normalization: how do you compare two RNA-Seq datasets generated at two different time points? One popular choice is RPKM which stands for reads per kb and per Million reads. The recipe is simple: divide by the total number of reads, divide by the length of the gene, and multiply by a number (one million, which is what M stands for) to get a number that is not too small. However, this may not be quite what you want. Indeed, we will now look at an example that illustrates the limitations of this measurement. Our dataset will be very basic (the code below is R code):

```
##let us assume that all genes have the same length, to put this problem aside for now
read.data <- data.frame (sample1 = 100, sample2 = c(rep(143, 70), rep(0, 30)))
## so one sample has all genes at level 100
##another has more reads for 70 genes and no read for the other 30
```

**Exercise:** Setting aside the gene length question, how would you compute RPKM values in this case? To get numbers easier to manipulate, let us rather use the number of reads per 1,000. How do you look into these RPKM values? Do they match what you would see as the intuitive explanation for these data?

```
RPKM.sample1 <- read.data$sample1/sum(read.data$sample1) * 1000
RPKM.sample2 <- read.data$sample2/sum(read.data$sample2) * 1000
table(RPKM.sample1)

## RPKM.sample1
## 10
## 100

table(RPKM.sample2)

## RPKM.sample2
##          0 14.2857142857143
##          30          70

## the tables above indicate that for sample1, all genes have the same level of
## expression (which happens to be 10). However for sample 2, 70 genes have a 14,
## and 30 genes have a 0 value. I don't think this matches the intuition we
## should have in this case. It seems to me that in that case one would rather
## conclude that 30 genes are not expressed in sample2, but that the other genes
## are pretty much on the same level.
```

So now what to do, and how to interpret the data. Here is an [interesting forum answer](#) that recapitulates the problem. I copy paste below the explanation:

To estimate the library size, simply taking the total number of (mapped or unmapped) reads is, in our experience, not a good idea. Sometimes, a few very strongly expressed genes are differentially expressed, and as they make up a good part of the total counts, they skew this number. After you divide by total counts, these few strongly expressed genes become equal, and the whole rest looks differentially expressed. The following simple alternative works much better:

- Construct a "reference sample" by taking, for each gene, the geometric mean of the counts in all samples.
- To get the sequencing depth of a sample relative to the reference, calculate for each gene the quotient of the counts in your sample divided by the counts of the reference sample. Now you have, for each gene, an estimate of the depth ratio.
- Simply take the median of all the quotients to get the relative depth of the library.

This is what the 'estimateSizeFactors' function of our DESeq package does.

This answer summarizes the problem well. One can see that in our case the 30 genes with 0 read have a large effect on the overall expression measurements that is probably not warranted. We will use the DESeq package to address this issue.

**Exercise:** There are three functions that you need to estimate the size factors in DESeq based on the dataset above. One is `newCountDataSet` to create a new object that DESeq can manipulate. The other two functions are: `estimateSizeFactors` and `sizeFactors` (the latter extracts the size factors from a DESeq object). Using these

two functions (and starting with loading the `DESeq` library), can you compute the size factors, normalize the data using these, and get new gene level estimates of expression? Is the result now more consistent with your intuitive interpretation of the data?

```
### Here is my answer to the question above, in R:
library(DESeq)
CDS <- newCountDataSet(read.data[, 1:2], condition = c('sample1', 'sample2'))
CDS <- estimateSizeFactors(CDS)
size.factors <- sizeFactors(CDS)
print(table(read.data$sample1 / size.factors[ 1 ]))

##
## 119.582607431014
##                100

print(table(read.data$sample2 / size.factors[ 2 ]))

##
##          0 119.582607431014
##        30                70
```

### 3.2 Further reading on normalization of RNA-Seq data

The issue of normalization of RNA-Seq data has been of interest to a lot of people. You can for example have a read of Lior Pachter's [blog post on the matter](#). The blog links to a talk that is probably interesting (though I have not yet seen it).

This [other blog post](#) is a good read for the list of available methods, even though I do not think I agree with all the details. In particular the sentence: "Again, the methods in this section allow for comparison of features with different length WITHIN a sample but not BETWEEN samples" does not make sense to me. If we normalize, it is exactly to compare data across samples (there is no point otherwise). So while there are caveats, as always, I don't think this (rather crucial) statement makes sense.

## 4 Aligning RNA-Seq data and estimating gene expression levels (45 minutes)

Aligning short-read RNA-Seq data is not fundamentally different from aligning DNA sequencing data. It is however made more complex by the presence of introns, which can create reads or paired-reads spanning large distances. A popular aligner for RNA-Seq data is **tophat** and we will go over some basic commands.

### 4.1 Aligning with tophat and bowtie

It is important to note that the underlying alignment engine for **tophat** is **bowtie**, hence many commands are shared with standard calls to **bowtie**. We start by building a **bowtie** index for a short portion of chromosome 12, which we will use as an example for this class. Go through the steps below. Alternatively you should be able to execute the script **scripts/tophat\_bowtie\_scripts.sh**. It will generate all the output files we want to look into, and the following goes through these commands in more details. But the step by step walkthrough is probably a better way to learn.

```
bowtie2-build -f ../data/RNASeq/chr12_short.fa ../data/RNASeq/chr12_short
```

With this, we can now perform the alignment step. But we first create some output folders to store all the output files:

```
mkdir results/tophat_output_with_gtf
## mkdir: results/tophat_output_with_gtf: File exists
```

Now we can start working with the fastq files:

```
f1=../data/RNASeq/reads_1.fq.gz
f2=../data/RNASeq/reads_2.fq.gz

tophat --no-coverage-search -o results/tophat_output_with_gtf -r 220 --library-type fr-unstranded \
--segment-length 30 -G ../data/RNASeq/chr12_short.gtf ../data/RNASeq/chr12_short ${f1} ${f2} \
--rg-sample test --rg-id test ## these sets the tags in the header of the BAM file
```

**Exercise:** For now, simply make sure you can run these scripts and that the output makes sense. Look at the BAM file, make sure you can see the appropriate tags in it. Also look at the general output of **tophat**.

**Exercise:** Can you see a subtle difference between a RNA-Seq BAM file and a DNA sequencing BAM file? The read extracted below should illustrate this.

```
samtools index results/tophat_output_with_gtf/accepted_hits.bam
samtools view results/tophat_output_with_gtf/accepted_hits.bam | grep 16M2150N34M > results/odd_read.sam

## open: No such file or directory
## [bam_index_build2] fail to open the BAM file.
## open: No such file or directory
## [main_samview] fail to open "results/tophat_output_with_gtf/accepted_hits.bam" for reading.
```

*#The Ns in the CIGAR string indicate introns, identified by split reads.*

### 4.2 Expression level estimation using Cufflinks

A popular software often associated with **tophat** is **cufflinks**. This piece of software is designed to estimate the abundance of each gene (and potentially isoforms). A call to **cufflinks** is pretty straightforward:

```
cufflinks -o results/cufflinks_output --GTF ../data/RNASeq/chr12_short.gtf \
results/tophat_output_with_gtf/accepted_hits.bam
```

**Exercise:** Go through the output, make sure you understand what all columns mean. Can you see the distinction between the isoform level and gene based estimates?



### 4.3 What if we do not have a GTF file?

A quick way to identify introns (and therefore exons) is to use these split reads to see where the gaps in the sequence are. I propose for example to pick a highly expressed gene (*DCP1B*) for example and to look at the introns in that gene. Because how to do this did not seem obvious to me, I wrote a small perl script that does the parsing (something very basic). See the example below, and make sure that you can run that code.

```
samtools index results/tophat_output_with_gtf/accepted_hits.bam
samtools view results/tophat_output_with_gtf/accepted_hits.bam 12:2055213-2113677 | \
  awk '{if ($6 ~ /N/ ) {print;}}' | ./scripts/get_introns.pl | \
  awk '{print $4"_"$5}' | sort | uniq -c > results/DCP1B_with_gtf.tab

## open: No such file or directory
## [bam_index_build2] fail to open the BAM file.
## open: No such file or directory
## [main_samview] fail to open "results/tophat_output_with_gtf/accepted_hits.bam" for reading.
```

There are plenty of situations where we do not have a GTF file. It could be because the species is not well annotated, or because for some reason you do not trust a published GTF file (low quality, unusual tissue type...).

**Exercise:** What does the output of `tophat` look like in the absence of a GTF? Start by creating a folder that will contain the output of tophat and run tophat without a GTF file. Do you identify the same introns and at the same frequency in the gene *DCP1B*? I suggest to start by creating a folder to store the data

```
mkdir results/tophat_output_no_gtf

## mkdir: results/tophat_output_no_gtf: File exists


f1=../data/RNASeq/reads_1.fq.gz
f2=../data/RNASeq/reads_2.fq.gz

tophat --no-coverage-search -o results/tophat_output_no_gtf -r 220 --library-type fr-unstranded \
  --segment-length 30 ../data/RNASeq/chr12_short ${f1} ${f2} \
  --rg-sample test --rg-id test ## these sets the tags in the header of the BAM file

## Error in running command bash
```

## 5 Differential expression analysis (30 minutes)

### 5.1 Using DESeq

We start by loading the DESeq package as well as an example dataset from a mouse brain RNA-Seq experiment.

```
library(DESeq)
load("../data/RNASeq/deseq_counts_TDP43.RData")
head(genes.counts)
```

##	control_rep1_dexseq_counts.txt		
## ENSMUSG00000000001	208		
## ENSMUSG00000000003	0		
## ENSMUSG00000000028	15		
## ENSMUSG00000000037	9		
## ENSMUSG00000000049	4		
## ENSMUSG00000000056	233		
##	control_rep2_dexseq_counts.txt		
## ENSMUSG00000000001	295		
## ENSMUSG00000000003	0		
## ENSMUSG00000000028	26		
## ENSMUSG00000000037	20		
## ENSMUSG00000000049	1		
## ENSMUSG00000000056	390		
##	control_rep3_dexseq_counts.txt		
## ENSMUSG00000000001	239		
## ENSMUSG00000000003	0		
## ENSMUSG00000000028	13		
## ENSMUSG00000000037	13		
## ENSMUSG00000000049	3		
## ENSMUSG00000000056	346		
##	control_rep4_dexseq_counts.txt	KD_rep1_dexseq_counts.txt	
## ENSMUSG00000000001	292	326	
## ENSMUSG00000000003	0	1	
## ENSMUSG00000000028	13	21	
## ENSMUSG00000000037	21	11	
## ENSMUSG00000000049	2	2	
## ENSMUSG00000000056	381	339	
##	KD_rep2_dexseq_counts.txt	KD_rep3_dexseq_counts.txt	
## ENSMUSG00000000001	371	316	
## ENSMUSG00000000003	0	0	
## ENSMUSG00000000028	22	18	
## ENSMUSG00000000037	12	18	
## ENSMUSG00000000049	1	1	
## ENSMUSG00000000056	359	317	
##	KD_rep4_dexseq_counts.txt		
## ENSMUSG00000000001	339		
## ENSMUSG00000000003	0		
## ENSMUSG00000000028	18		
## ENSMUSG00000000037	30		
## ENSMUSG00000000049	2		
## ENSMUSG00000000056	379		

We can now define the model for the differential expression analysis:

```
formula1 <- count ~ condition
formula0 <- count ~ 1
design.deseq <- c('control', 'control', 'control', 'control', 'KD', 'KD', 'KD', 'KD')
```

And now the computations can properly start. Note that these steps are very long, and therefore the code is not executed as part of this file (to be more precise, it is executed once, and the output is saved).

```

CDS <- newCountDataSet(genes.counts, condition = design.deseq)

CDS <- estimateSizeFactors(CDS)
CDS <- estimateDispersions(CDS, method = 'pooled')

fit0 <- fitNbinomGLMs( CDS, formula0 )
fit1 <- fitNbinomGLMs( CDS, formula1 )

deseq.pval <- fit1
deseq.pval$EnsemblID <- row.names( deseq.pval)
deseq.pval$basic.pval <- signif(nbinomGLMTest( fit1, fit0 ), 4)
save(list = 'deseq.pval', file = 'results/DE_pvalues_ranked.RData')

```

See below some polishing: a multiple testing/false discovery rate Bonferroni-Hochberg analysis, and the ordering of the results by significance of P-values.

```

load('results/DE_pvalues_ranked.RData')
deseq.pval$adj.pval <- signif(p.adjust( deseq.pval$basic.pval, method="BH" ), 4)

deseq.pval <- deseq.pval[ order(deseq.pval$basic.pval, decreasing = FALSE), ]
head(deseq.pval)

##                (Intercept) conditionKD  deviance converged
## ENSMUSG00000023224      5.788135      1.732927  4.110022      TRUE
## ENSMUSG00000023826      6.559152     -2.002398  7.912218      TRUE
## ENSMUSG00000026547      6.032272      1.688365  3.296005      TRUE
## ENSMUSG00000039419      9.697903     -1.184795 12.093286      TRUE
## ENSMUSG00000040424      8.449948     -1.373362  6.262590      TRUE
## ENSMUSG00000041459     10.080164     -1.691099  5.525777      TRUE
##                EnsemblID basic.pval adj.pval
## ENSMUSG00000023224 ENSMUSG00000023224          0          0
## ENSMUSG00000023826 ENSMUSG00000023826          0          0
## ENSMUSG00000026547 ENSMUSG00000026547          0          0
## ENSMUSG00000039419 ENSMUSG00000039419          0          0
## ENSMUSG00000040424 ENSMUSG00000040424          0          0
## ENSMUSG00000041459 ENSMUSG00000041459          0          0

```

## 5.2 Using DESeq2

Relatively recently, the authors of DESeq have released a new version of the DESeq package and called it DESeq2. The commands are similar, but there are also differences. The best way to learn about a R package is to work through the vignette, which highlights the main capabilities of the package. The vignette for DESeq2 is located [here](#).

**Exercise:** Perform a differential expression analysis using the newer DESeq2 package. This is essentially about finding the right commands in the vignette and transposing them to your situation.

```

library(DESeq2)
load("../data/RNASeq/deseq_counts_TDP43.RData")
design.deseq <- c('control', 'control', 'control', 'control', 'KD', 'KD', 'KD', 'KD')

dds <- DESeqDataSetFromMatrix(countData = genes.counts,
                              colData = data.frame(KD = design.deseq),
                              design = ~ KD)

dds <- DESeq(dds)
res <- results(dds)
resOrdered <- res[order(res$padj),]
head(resOrdered)

```

## 6 Galaxy server

Have a look at the Galaxy server tools for RNA-Seq analysis. Can you replicate the alignment steps? And can you run cufflinks as well? Get a feeling for the tools that Galaxy offers and decide whether you much rather the (more constrained) web interface, or whether you are OK with the command lind tools. Both solutions are absolutely acceptable.