

Introduction to RNA-Seq data analysis

Cancer Genomics Course, EMBL-EBI

Alexey Larionov

20 June 2019

Plan of the practical session

1. Introduction
 - Environment
2. Alignment with *minimap2*
3. Explore RNA-Seq BAM
 - *IGV*
 - *samtools*
 - *BAMSeek*
4. Calculating transcript counts with *Salmon*
5. Differential gene expression with *DESeq2*
 - Prepare R environment and accessory data
 - Import Salmon counts into R using *tximport*
 - Read data to *DESeq-DataSet* object
 - Calculate differentially expressed genes
 - Visualising results
6. *STAR-fusion*
 - Detecting fusion transcripts
 - Visualising fusion transcripts using *Chimeraviz*

1. Introduction

RNA-Seq data can be analysed in many different ways (see lecture slides). Numerous good tools and algorithms have been developed to process RNA-Seq data. This session will provide practical experience with a number of selected tasks and tools, aiming to illustrate the diversity of RNA-Seq data analysis.

This handout shows only the key elements of the code. The full code is provided in the accompanying scripts, which should be studied during the practical.

To follow the **reproducible research** principles, the output of Bash scripts may be captured to logs like this:

```
./script.sh &> log.txt & # ./script.sh : start script
                        # & > log.txt : direct output and errors to log.txt,
                        # & : return control to the terminal (optional)

tail -f log.txt # Watch the log updates (optional)
                # Exit with ctr-C
```

For the analysis in R, the output of Rmd scripts can be “knitted” to nice-looking html logs, as will be shown during the session.

Environment

VM specifications

Check the number of cores and amount of memory available on your machine (hint: you may use **htop** command in terminal). Some of the tools allow for multi-processing: the number of requested threads should not exceed the number of available cores. Some of the tools will require at least 20-30GB of RAM.

Working folder

Explore ... /**RNA-Seq** folder *on your machine*. It should contain subfolders for scripts, data, resources and results. There is no actual results or resources within these folders yet. It is expected that you will create results and some additional resources during this session.

There is a similar folder on **penelopeprime** drive: ... /**Cancer_Genomics_June19/Day_4/RNA-Seq**.

It contains copies of the source data and scripts. You may copy files from **penelopeprime** to your machine, if you need. However, don't change, remove or modify any files directly on **penelopeprime**.

Shared resources

In addition to the copies of source data and scripts, RNA-Seq folder on **penelopeprime** contains some shared resources that will be used during the practical session:

- b38 transcriptome fastq
- b38 transcriptome annotations
- b38 Cancer Transcriptome Analysis Toolkit (**CTAT**) library from the Trinity project: https://github.com/NCIP/Trinity_CTAT/wiki

You *do not need* to copy these shared resources from **penelopeprime** to your home folder on VM. The source scripts will look for the resources directly on **penelopeprime**.

CTAT resources library

Amongst other files, the CTAT library includes

- reference genome in FASTA format: **ref_genome.fa**
- reference annotations in GTF format: **ref_annot.gtf**
- index and other resources for STAR-Fusion: **ref_genome.fa.star.idx** etc

Make yourself familiar with the content of **ref_annot.gtf** file (hint: you may use **head** command in a terminal)

BAM files to explore in IGV

Finally, **penelopeprime** contains some BAM files in folder ... /**Cancer_Genomics_June19/Day_4/RNA-Seq/data/bams**. These files will be used to explore different types of RNA-Seq data (short-reads, nanopore and pacbio) in IGV.

2. Alignment with minimap2

minimap2 is a splice-aware aligner, which could be used to align RNA-Seq data to reference genome.

It is written by an extremely reputable author, it is able to deal with both short and long reads, and it, supposedly, outperforms many other aligners, which can be used for RNA-Seq:

- <https://github.com/lh3/minimap2>
- <https://doi.org/10.1093/bioinformatics/bty191>
- <http://bioinfo.zesoi.fer.hr/index.php/hr/blog-en/56-gmap-vs-minimap2>
- <https://doi.org/10.1093/bioinformatics/btx668>

Unlike most of other aligners, **minimap2** does not require a pre-build index before the alignment. In fact, it builds the index on-the-fly: it only takes a couple of minutes for minimap2 to generate an index for human genome. Although it is possible to build and save the index, it still may be a good idea to generate index on-the-fly because the preferred indexing settings might slightly change for different types data.

Source data

For the **minimap2** alignment exercise you will use nanopore RNA-Seq FASTQ file located in `.../RNA-Seq/data/nanopore` folders on your machine. This dataset has been taken from the paper, which explored 2D ONT sequencing of to study EGFR-related events in H1975 cell line:

- <https://doi.org/10.1093/dnares/dsx027>

The FASTQ file presents sequencing of EGFR-containing cDNA amplicon.

Alignment script

The alignment script is located on your machine in folder `.../RNA-Seq/scripts/s01_alignment`. It illustrates alignment of nanopore data:

```
minimap2 -t 10 -ax splice -L "${ref_genome}" "${fastq}" > "${sam}"
```

`-t 10` instructs to use 10 threads (to speed up the calculations). You may increase it to 18 (your machine should have 20 cores). You may watch how computer uses resources during the calculations using **htop** command in a separate terminal window.

`-ax splice` defines output to BAM (a) and a set of parameters for splice-aware alignment (x splice). Please note that PacBio RNA-Seq alignment might require slightly different settings (`-ax splice:hq`)

`-L` is a safeguard against a bug within BAM file format specification:

<https://github.com/lh3/minimap2#working-with-65535-cigar-operations>

For the same reason we output data to SAM, not to BAM format.

Results

Alignment of the nanopore data should take less than 10 min. If you run the provided script as intended, the resulting SAM file should appear in the `.../RNA-Seq/results/s01_alignment` folder on your machine. After the alignment is complete, you may explore the log (if you use logging) and the SAM file (e.g. `samtools flagstat`, `samtools view | head` etc).

STAR alignment of short-read RNA-Seq data (optional)

You are not advised to run this step because of the short time of the practical session. However, just in case, I placed a script that shows how to use STAR for alignment of Illumina PE75 RNA-Seq FASTQ files.

The used Illumina PE75 RNA-Seq FASTQ files have been taken from the paper, which studied mechanisms of drug resistance in lung cancer:

- <https://doi.org/10.1158/0008-5472.CAN-17-3146>

A small subset of the whole data had been extracted to allow alignment within 10 min. The extracted slice includes reads for EGFR, KRAS, ALK and some other cancer-related genes.

Please note that running two alignment scripts at the same time (minimap2 and STAR) might not be a good idea because they would compete for the cores and memory on your machine.

3. Explore RNA-Seq BAM

You may start exploring RNA-Seq BAMs before the alignment is complete.

IGV

IGV (Integrated Genome Viewer) is the *de-facto* standard for exploring and visualising BAM files. Most likely you have already used it during the course to explore BAM files generated in DNA-sequencing. However, the default IGV settings may need to be adjusted for comfortable viewing of *spliced* RNA alignments, as described here:

- https://software.broadinstitute.org/software/igv/splice_junctions

Also, IGV provides special tool for viewing *Sashimi plots*:

- <http://software.broadinstitute.org/software/igv/Sashimi>

Adjusting IGV settings

First, go to **Menu > View > Preferences > Alignments**:

- Set **Visibility range 500** : this is necessary to see the alignment track for entire genes
- Note that the latest IGV version contains panels to configure settings for RNA-Seq and Long-Read (Third Gen) data. Because IGV documentation has not yet been updated to describe these panels, we will not use these panels.
- Note options related to **Splice Junction Track**. We will configure **Junction Track** using right-click menu later, when viewing the data.

RNA-Seq tutorial available on IGV server

You may use RNA-Seq data from IGV Server to explore RNA-Seq data, while **minimap2** is still performing our own alignment:

- Set genome to **hg19** (this is an old version of human reference genome)
- Load data: **Menu > File > Load from Server > Tutorials > RNA-Seq (Body Map)**
- Navigate to **SLC25A3** gene
- Use right-click menu to show **Splice Junction Track**
- Switch “on” **Autoscale** and apply **Expand** option to the junction truck(s)
- Show **Sashimi plots** for both heart and liver (right click on a junction truck -> Sashimi plot)

- Using context menu (right click) on Sashimi plot:
 - Show/Hide **exon coverage data**
 - Set **minimal junctions coverage** to 10 (or pick the threshold that you like :)
 - Find whether the transcripts came from **Forward or Reverse strand**

Is there any evidence for alternative splicing of SLC25A3 gene in these two tissues?

Compare different types of RNA-Seq data in IGV

For this section you will use the BAM files provided in folder `/Day_4/RNA-Seq/data/bams` on **penelopeprime** drive. Please **copy** this folder from **penelopeprime** to any suitable location on **your machine** before you start exploring the files.

The folder contains 3 bam files with indices (bai). The type of data is indicated by the suffixes in file names: “ont” for Oxford-Nanopore-Technology, “sr” for Illumina PE75 short reads, “pb” for PacBio.

All these BAMs were aligned to b38. So **set genome to hg38** before exploring the files. To open the files in IGV use menu: **Menu -> File -> Load from file -> ...path-to/bams/...**

SRR3534924 short-reads sample

- Explore **CCND1** and **CCNB1** genes: try to show Autoscale, Expand/Collapse, Sashimi plot. Which of these genes does not show evidence of alternative splicing? Are they expressed on forward or reverse strand?
- Another gene that may be interesting to explore is **CSDE1**. Some of its transcripts skip Exon 2. Is it a common event? Is exon 3 used by any transcript? Is this gene expressed on Forward or Reverse strand?
- Look at **ALK** gene: can you see any irregularity in expression of this gene?
- You may also explore **NRAS**, **KRAS**, **HRAS**, **EGFR** and **CDKN3** (if you have time :)
- Why there is no expression of **ACTB**? (hint: I made a slice of the data :)

DRR059318-ONT and SRR7346977-PacBio samples

DRR059318-ONT contains only cDNA amplicon of EGFR gene.

SRR7346977-PacBio contains chromosomes 7,8,14 and 17 only.

Explore and compare results from 3 sequencing technologies (e.g. **EGFR** gene is present in all 3 BAM files). Which technology has the best base qualities?

samtools (optional)

Samtools is a powerful toolset, which can be used for variant calling and for many other tasks concerning BAM/SAM files (see <http://www.htslib.org/doc/samtools.html> for details).

In a terminal go to the folder containing the BAM files and execute the following commands:

```
samtools flagstat SRR3534924_selected_sr.bam
samtools view SRR3534924_selected_sr.bam | head
samtools view SRR3534924_selected_sr.bam | head -n 300 | awk '$6 ~ "N"'
```

What the last command does?

What is special about the reads shown by the last command?

BAMSeek (optional)

BAMSeek is a handy GUI tool to explore BAM content at low level. You can find it in `.../RNA-Seq/tools`. Launch it by double-clicking on the icon; then use menu to navigate and open a BAM file. When opening a BAM file for the first time, BAMSeek creates its own index for the file.

Use BAMSeek to open the short-reads bam file. Can you see the content of the BAM file header? Is BAM file sorted? What version of BAM-format is used?

Use BAMseek to open a nanopore bam. Why CIGAR string may exceed the size initially reserved in BAM specification?

4. Calculating transcript counts with *Salmon*

Arguably, **gene expression analysis** could be perceived as the main application of RNA-Seq (of course, people focused on fusion detection will disagree :) The first step in the expression measurement is counting reads overlapping the genes (discussed in more details during the lecture). Traditionally, read count was based on genomic alignments, as implemented in *Cufflinks* or *RSEM*. However, alignment-free methods, such as *Kallisto* or *Salmon*, are becoming increasingly popular.

In this session we will use *Salmon*. In fact, instead of relying on another tool for genome-alignment, *Salmon* performs its own fast and accurate “quasi-mapping” to transcriptome:

<https://combine-lab.github.io/salmon/>

Preparing transcriptome file for quasi-mapping

Human transcriptome could be obtained/compiled from publicly available sources, such as Ensembl. For this tutorial a combination of cDNA and non-coding RNA sequences was chosen as the transcriptome:

```
wget ftp://ftp.ensembl.org/pub/path/to/cdna/Homo_sapiens.GRCh38.cdna.all.fa.gz
wget ftp://ftp.ensembl.org/pub/current_fasta/homo_sapiens/ncrna/Homo_sapiens.GRCh38.ncrna.fa.gz

zcat Homo_sapiens.GRCh38.cdna.all.fa.gz Homo_sapiens.GRCh38.ncrna.fa.gz >
  Homo_sapiens.GRCh38.rna.fa.gz
```

This has been done already, so you do not need to repeat this step. You can find the transcriptome file in the **penelopeprime** drive `.../Day_4/RNA-Seq/resources` folder.

Prepare Salmon index

Like most of the other tools performing “alignment-like” tasks, Salmon prepares index for the reference. In this case it prepares index for the reference transcriptome:

```
salmon index \
-t "/path/to/Homo_sapiens.GRCh38.rna.fa.gz" \
-i "/path/to/salmon_index"
```

Again, there is a script **salmon_index.sh**, which you can use to create the index. The intended location of the index is `.../RNA-Seq/resources/salmon_index` on *your machine*.

Run Salmon

The quantitation step could be performed using **salmon quant** command:

```
salmon quant \  
--index "/path/to/salmon_index" \  
--libType A \  
--mates1 "fastq1" \  
--mates2 "${fastq2}" \  
--output "${out_folder}" \  
--threads 10
```

Use `salmon quant --help-reads` for information about the options. Script `salmon_quant.sh` is provided to run the quantification. The intended output folder is `.../RNA-Seq/results/s02_transcripts_count` on *your machine*.

Have a look at the content of the output folder after the run. Amongst the other files, there will be a `quant.sf` file. This is the result of Salmon quantification. Have a quick look at the content of this file (hint: you may use `head` command in the terminal).

5. Differential gene expression with *DESeq2*

Detecting differentially expressed genes is an important step in the gene expression analysis. The top differentially expressed genes are often used for gene expression signatures for molecular classification of cancer, to evaluate prognosis or to predict response to treatment.

There are many statistical approaches for detecting differential expression. We will use the approach implemented in **DESeq2** R package. It assumes Negative Binomial distribution of the counts. This assumption is based on the fact that empirical counts data are closer to the Negative Binomial than, for instance, to Normal or Poisson distributions.

In addition to the sophisticated mathematical methods, *DESeq2* package provides an implementation of these methods and a number of convenience functions for data processing and assessment: <https://bioconductor.org/packages/release/bioc/vignettes/DESeq2/inst/doc/DESeq2.html>

Source data

`quant.sf` files were pre-calculated for 3 Tumour-Normal pairs of renal cancer. The data was generated by CAGEKID consortium (<https://www.cng.fr/cagekid/index.html>). We will use these pre-calculated Salmon counts to detect genes differentially expressed between renal carcinoma and normal tissue.

The data are located in `.../RNA-Seq/data/kidney_cancer/salmon_counts` folder on *your machine*. The patients IDs are *LT344*, *LR364* and *LR365*. Suffixes *N* and *T* indicate to Normal and Tumour sample respectively.

R-markdown script

For this part of the workshop we switch from Shell scripts to R environment. The code for our entire Differential Expression analysis is provided in `DESeq2_analysis.Rmd` script. The key fragments of this script will be discussed in the text below. However, plenty of additional information is provided within the script. You may review and run the corresponding parts of the script along with reading this text.

The script is written using **R-markdown**. This allows organising code into convenient **chunks**, adding formatted comments between the chunks and, most importantly, it allows to **knit** a detailed report of the performed analysis. R-markdown is a widely used tool facilitating **reproducible research** practices in R environment.

Prepare R environment and accessory data

The **Start chunk** of the script cleans the environment (just in case) and shows how the required packages were installed. The installation lines are commented because it needed to be done just once. Then we load **dplyr** and **ggplot2**: these are general libraries commonly used for data handling and plotting. The specialised libraries are loaded later: in the chunks, where they are used.

The following accessory data should be prepared for **DESeq2** analysis:

List of the “.sf” files, which contain *Salmon* counts

```
# Get list of folders with salmon counts
counts_folder=paste(base_folder, "data/kidney_cancer/salmon_counts", sep="/")
salmon_folders <- list.dirs(counts_folder, recursive = FALSE)

# Make list of sf files
salmon_sf_files <- paste(salmon_folders, "quant.sf", sep="/")

# Assign names to the vector of files
names(salmon_sf_files) <- basename(salmon_folders)
```

Data frame with samples description

```
# Prepare vectors with data
run <- c("LR344N", "LR344T", "LR364N", "LR364T", "LR365N", "LR365T")
condition <- c("N", "T", "N", "T", "N", "T")
patient <- c("LR344", "LR344", "LR364", "LR364", "LR365", "LR365")

# Combine vectors to data frame
samples.df <- data.frame(run, condition, patient)

# Convert condition column to factor
samples.df$condition <- as.factor(samples.df$condition)

# Add rownames
rownames(samples.df) <- samples.df$run

# Check result
samples.df
```

```
##           run condition patient
## LR344N LR344N         N   LR344
## LR344T LR344T         T   LR344
## LR364N LR364N         N   LR364
## LR364T LR364T         T   LR364
## LR365N LR365N         N   LR365
## LR365T LR365T         T   LR365
```


Table linking transcript-ID to gene-ID

```
# Load required library
require(rtracklayer) # contains readGFF function: see ?readGFF

# Read annotations file
annotation_file <- paste(base_folder, "resources/annotations/Homo_sapiens.GRCh38.89.gtf", sep="/")
gene_annotation.df <- readGFF(annotation_file)

# Make table linking transcripts to genes
trans2genes.df <- gene_annotation.df %>%
  filter(type == "transcript") %>%
  select(transcript_id, transcript_version, gene_id)

# Combine transcript ID and version
trans2genes.df <- trans2genes.df %>%
  mutate(transcript_id_version=paste(transcript_id, transcript_version, sep=".") %>%
  select(transcript_id_version, gene_id)

# Check result
head(trans2genes.df)
```

```
## transcript_id_version      gene_id
## 1 ENST00000456328.2 ENSG00000223972
## 2 ENST00000450305.2 ENSG00000223972
## 3 ENST00000488147.1 ENSG00000227232
## 4 ENST00000619216.1 ENSG00000278267
## 5 ENST00000473358.1 ENSG00000243485
## 6 ENST00000469289.1 ENSG00000243485
```

Table linking gene-ID to gene-name

```
# Make data frame
geneId2Name.df <- gene_annotation.df %>%
  select(gene_id, gene_name) %>%
  arrange(gene_name) %>%
  distinct()

# Copy gene IDs to rownames
rownames(geneId2Name.df) <- geneId2Name.df$gene_id

# Check result
head(geneId2Name.df) # NB: gene names are not unique !
tail(geneId2Name.df)
```

```
## gene_id gene_name
## ENSG00000252830 ENSG00000252830 5S_rRNA
## ENSG00000276442 ENSG00000276442 5S_rRNA
## ENSG00000274408 ENSG00000274408 5S_rRNA
## ENSG00000274059 ENSG00000274059 5S_rRNA
## ENSG00000277313 ENSG00000277313 7SK
## ENSG00000275933 ENSG00000275933 7SK
```

```
## ...
##           gene_id  gene_name
## ENSG00000229956 ENSG00000229956 ZRANB2-AS2
## ENSG00000121903 ENSG00000121903  ZSCAN20
## ENSG00000162415 ENSG00000162415  ZSWIM5
## ENSG00000203995 ENSG00000203995  ZYG11A
## ENSG00000162378 ENSG00000162378  ZYG11B
## ENSG00000036549 ENSG00000036549   ZZZ3
```

At this point all the accessory tables have been prepared.

Import Salmon counts into R

tximport R package provides a convenient way of importing transcript abundance data from multiple upstream applications, including *Salmon*, *Kallisto*, *RSEM* and others:

<https://bioconductor.org/packages/release/bioc/vignettes/tximport/inst/doc/tximport.html>

Loading *Salmon* data can be preformed in the following way:

```
# Load required libraries
require(tximport) # ?tximport
require(rjson) # required for soem tximport features

# Read and convert salmon counts
gene_counts.ls <- tximport(salmon_sf_files, type="salmon", tx2gene=trans2genes.df)

# Check results
head(gene_counts.ls$counts)
```

```
##           LR344N  LR344T  LR364N  LR364T  LR365N  LR365T
## ENSG00000000457 454.3549 621.4804 429.68625 642.2974 654.4376 280.6368
## ENSG00000000460 128.3026 277.5019 105.65094 334.8424 160.7700 2711.3892
## ENSG00000000938 140.0000 992.0003 71.99995 674.9995 206.0005 349.0004
## ENSG00000000971 2255.8532 3603.0045 1536.92162 9826.9728 3812.2033 25005.4770
## ENSG00000001460 541.3224 380.3891 274.47184 496.8335 357.9931 264.4049
## ENSG00000001461 1121.4062 1220.1628 754.08937 2521.9820 1296.2460 945.2384
```

Note that **tximport** places data into a *list* object. Also it aggregates the transcript-level counts into genes.

Later the **tximport list** will be transformed into **DESeq-DataSet** object that is needed to perform the differential expression tests.

Explore imported data

Gene expression data assessment and QC could be done using **hierarchical clustering** or **PCA** analysis. Note that this text illustrates only selected steps of data QC: review and run the accompanying script for more details.

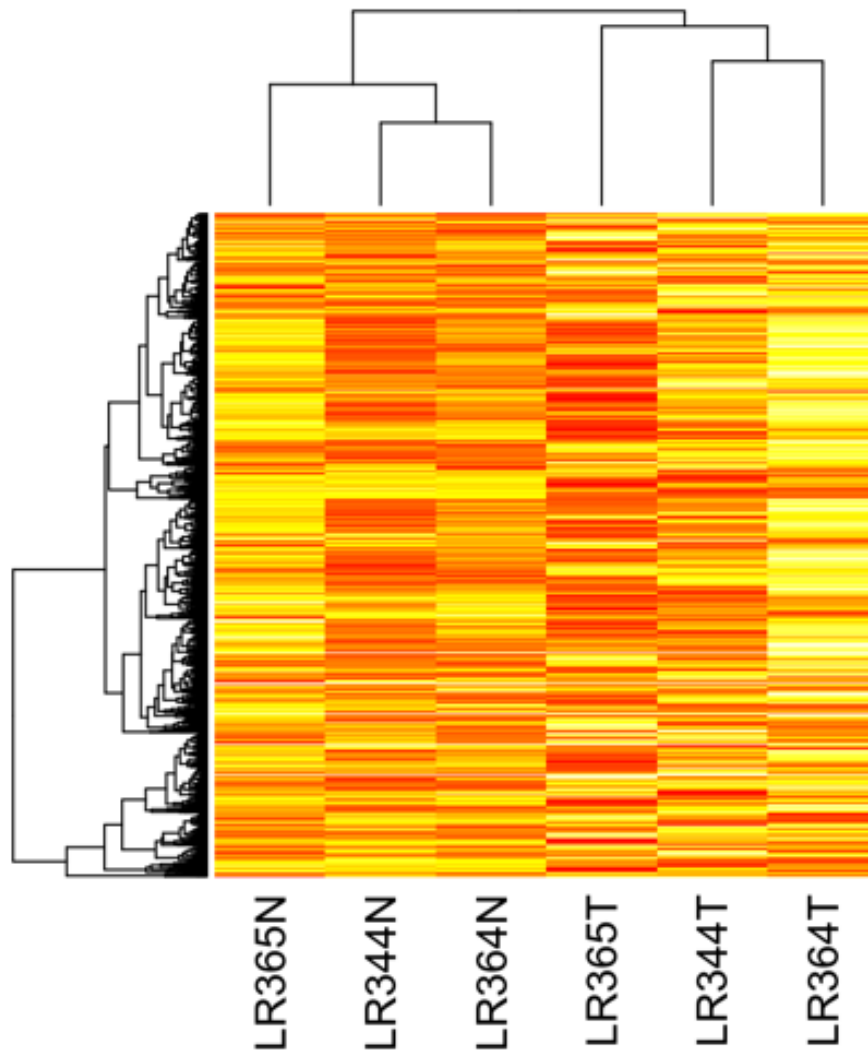
Within our **DESeq2** workflow the data assessment and QC could be performed at two points:

- immediately after the data import into the **tximport list**
- at a later stage: after importing data into **DESeq-DataSet** object

For teaching purposes, we will do data assessment at both of these steps and compare the results.

Hierarchical clustering with heatmap

```
# Extract and log-transform normalised expression values  
expr.mx <- gene_counts.ls$abundance  
expr.mx <- log2(expr.mx + 1) # +1 in case if some expressions are 0  
  
# Filter low expressed genes  
low_expressed_genes <- mean_expr < 1  
expr.mx <- expr.mx[! low_expressed_genes, ]  
  
# Plot heatmap  
heatmap(expr.mx, labRow=NA)
```

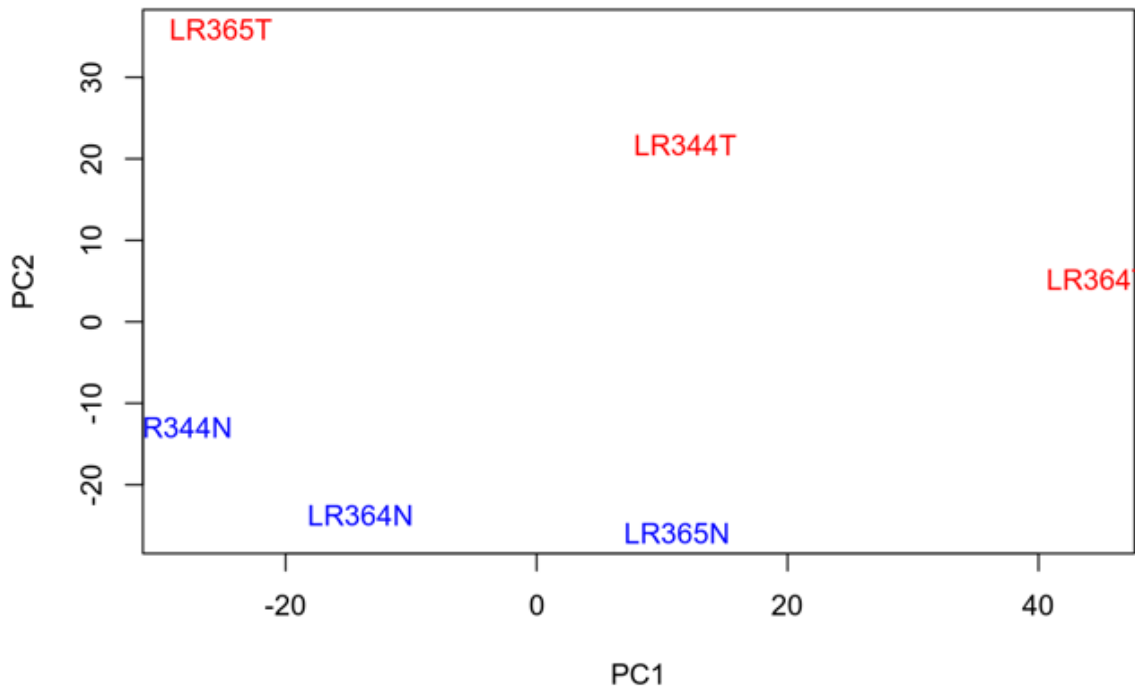


PCA plot

```
# Calculate PCA
expr.pca <- prcomp(t(expr.mx),
                  center = TRUE,
                  scale = TRUE)

# Extract matrix of PCs
pc.mx <- expr.pca$x

# Plot PC1 vs PC2
PC1 <- pc.mx[,1]
PC2 <- pc.mx[,2]
plot(PC1, PC2, type="n")
text(PC1, PC2, labels = row.names(pc.mx),
     col=c("blue","red","blue","red","blue","red"))
```



Both techniques (**hierarchical clustering** and **PCA** analysis) showed separation of tumour and normal samples, even using the entire set of expressed genes. This confirms a good quality of data and suggests presence of strong differentially expressed genes, which could be detected in downstream analysis.

Read data to *DESeq-DataSet* object

At this step we will read data into the **DESeq-DataSet** object and filter data, re-using the list of low-expressed genes from the previous step (more details about selecting the low-expressed genes are shown in the accompanying script).

```
# Load required library
require(DESeq2)
##?DESeqDataSetFromTximport

# Read data into DESeq2 dataset
dds <- DESeqDataSetFromTximport(gene_counts.ls, colData = samples.df, design = ~condition)

# Remove low-expressed genes from DESeq2 dataset object
dds <- dds[! low_expressed_genes, ]

# Check result
dds

## class: DESeqDataSet
## dim: 2065 6
## metadata(1): version
## assays(6): counts avgTxLength ... H cooks
## rownames(2065): ENSG00000000457 ENSG00000000460 ... ENSG00000283761 ENSG00000283773
## rowData names(22): baseMean baseVar ... deviance maxCooks
## colnames(6): LR344N LR344T ... LR365N LR365T
## colData names(3): run condition patient
```

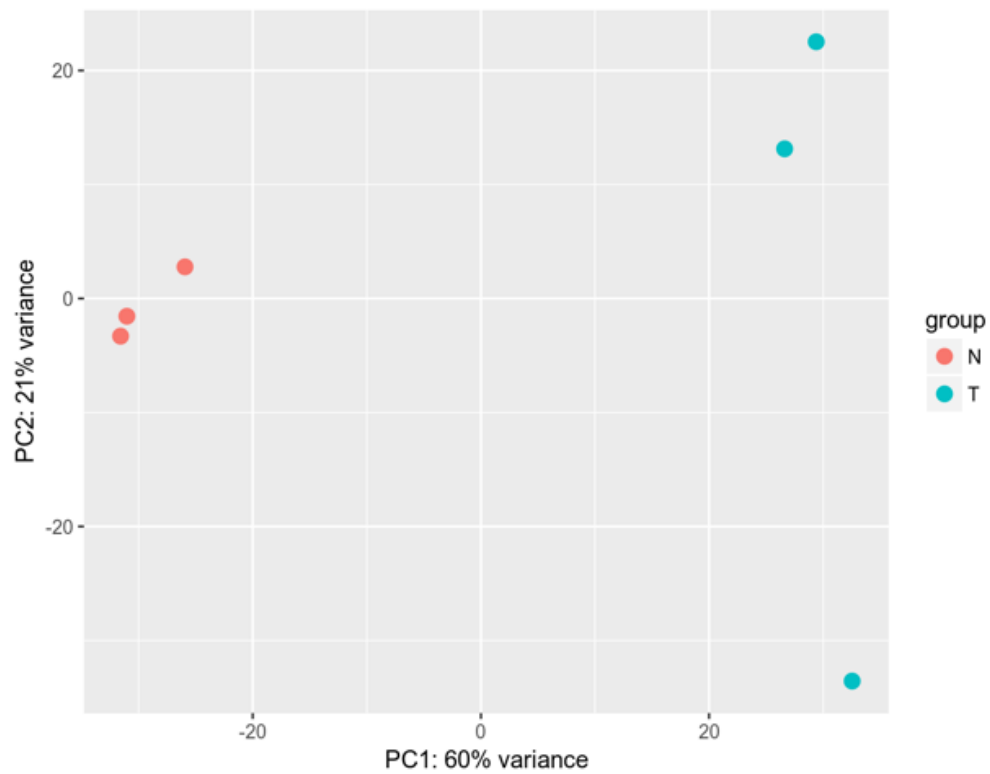
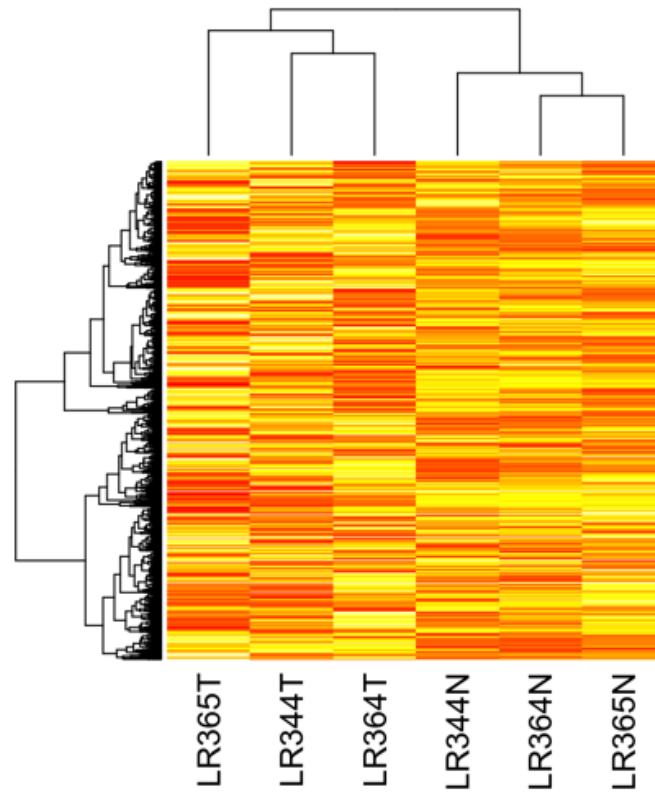
Explore data within *DESeq-DataSet* object

The data have just been placed into **DESeq-DataSet** object. No tests for differential expression have yet been done. However, we may already use some convenience tools provided by *DESeq2* package for data assessment within **DESeq-DataSet** object. These tools include advanced normalisation functions (such as **varianceStabilizingTransformation**) and an internal function to make PCA plot on the vst-transformed data.

```
# Apply VST (for visualising)
vsd <- varianceStabilizingTransformation(dds)

# Heatmap after VST
heatmap(assay(vsd), labRow=NA)

# PCA after VST (built-in DESeq2 function)
plotPCA(vsd)
```



Consistent with the previous assessment, both methods show a good separation of normal and tumour samples after VST. This confirms good quality of data and suggests presence of differentially expressed genes.

Calculate differentially expressed genes

Finally, we are in a position to make tests for differential expression.

The testing algorithm includes estimation of effective library size, genes dispersion and including these parameters into Negative Binomial Generalized Linear Model:

<https://genomebiology.biomedcentral.com/articles/10.1186/s13059-014-0550-8>

<https://bioconductor.org/packages/release/bioc/vignettes/DESeq2/inst/doc/DESeq2.html#theory>

Luckily, all this mathematical complexity is wrapped into a single simple function call:

```
dds <- DESeq(dds)
```

The test results have now been added to the **DESeq-DataSet** object.

The table with Fold-Change and adjusted P-values can now be extracted from the **DESeq-DataSet** object using *results()* function:

```
result.df <- as.data.frame(results(dds))
head(result.df)
```

##		baseMean	log2FoldChange	lfcSE	stat	pvalue	padj
##	ENSG00000000457	516.8618	-0.1215062	0.5610834	-0.2165564	0.8285540686	0.945858318
##	ENSG00000000460	722.2178	3.3731988	0.8638281	3.9049422	NA	NA
##	ENSG00000000938	370.1845	2.0667952	0.5567747	3.7120853	0.0002055587	0.002680444
##	ENSG00000000971	8052.6020	2.3867098	0.7483717	3.1892038	0.0014266527	0.013044657
##	ENSG00000001460	377.8124	-0.1443089	0.5027592	-0.2870338	0.7740864756	0.933888418
##	ENSG00000001461	1265.1619	0.5165444	0.5161500	1.0007642	0.3169408247	0.657948099

Err ... It's exciting to know that ENSG00000000938 has adjusted P-value of 0.002 ...

Still, it seems that some finishing touches might be helpful:

```
# Add gene_id to colimms
result.df <- cbind(gene_id=rownames(result.df), result.df)

# Remove genes with no p-value
uninformative_genes <- is.na(result.df$padj)
result.df <- result.df[!uninformative_genes,]

# Add gene names and order by p-value
result.df <- left_join(result.df, geneId2Name.df, by="gene_id") %>%
  select(gene_id, gene_name, baseMean, log2FoldChange, padj) %>%
  arrange(padj, desc(log2FoldChange))

# Copy gene names to rownames
rownames(result.df) <- result.df$gene_name

# Check result
head(result.df)
```

##		gene_id	gene_name	baseMean	log2FoldChange	padj
##	NPHS2	ENSG00000116218	NPHS2	3498.4932	-9.798823	1.979669e-59
##	FAM151A	ENSG00000162391	FAM151A	14583.8887	-8.725494	1.313300e-44
##	DIO1	ENSG00000211452	DIO1	5840.0166	-8.609517	1.012498e-40
##	MCOLN3	ENSG00000055732	MCOLN3	389.4464	-5.605388	2.719366e-23
##	VTCN1	ENSG00000134258	VTCN1	397.4958	-6.727450	3.816724e-22
##	KCNJ10	ENSG00000177807	KCNJ10	784.2833	-6.443739	6.568008e-22

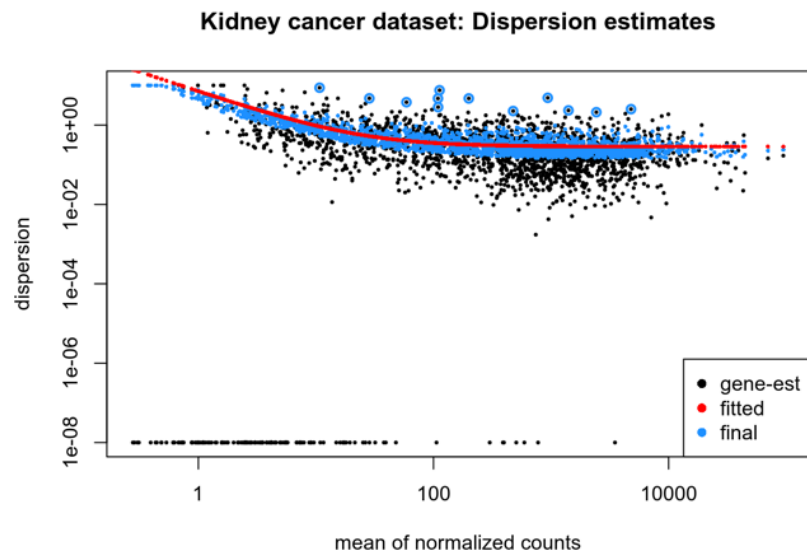
It looks much better now !

Visualising results of differential expression analysis

In addition to the table with p-values and fold-changes, **DESeq2** package provides several functions to visualise results.

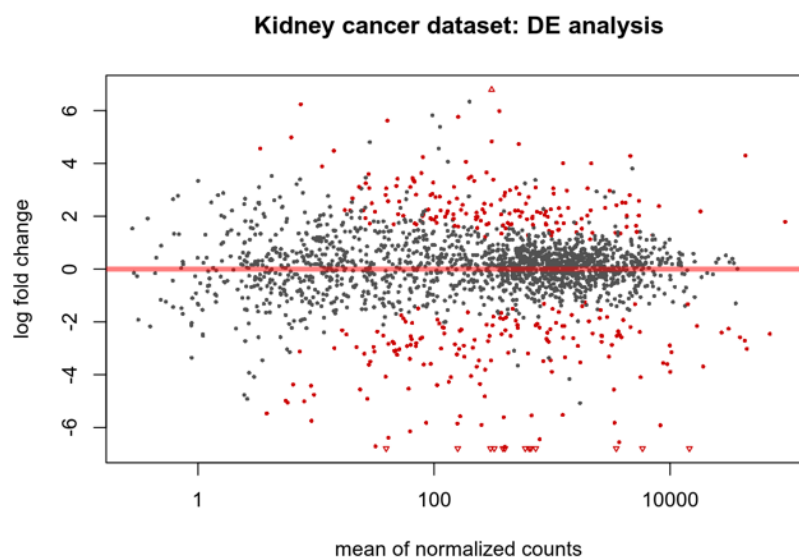
Dispersion estimates

```
plotDispEsts(dds, main="Kidney cancer dataset: Dispersion estimates")
```



MA plot

```
plotMA(dds, main="Kidney cancer dataset: DE analysis")
```



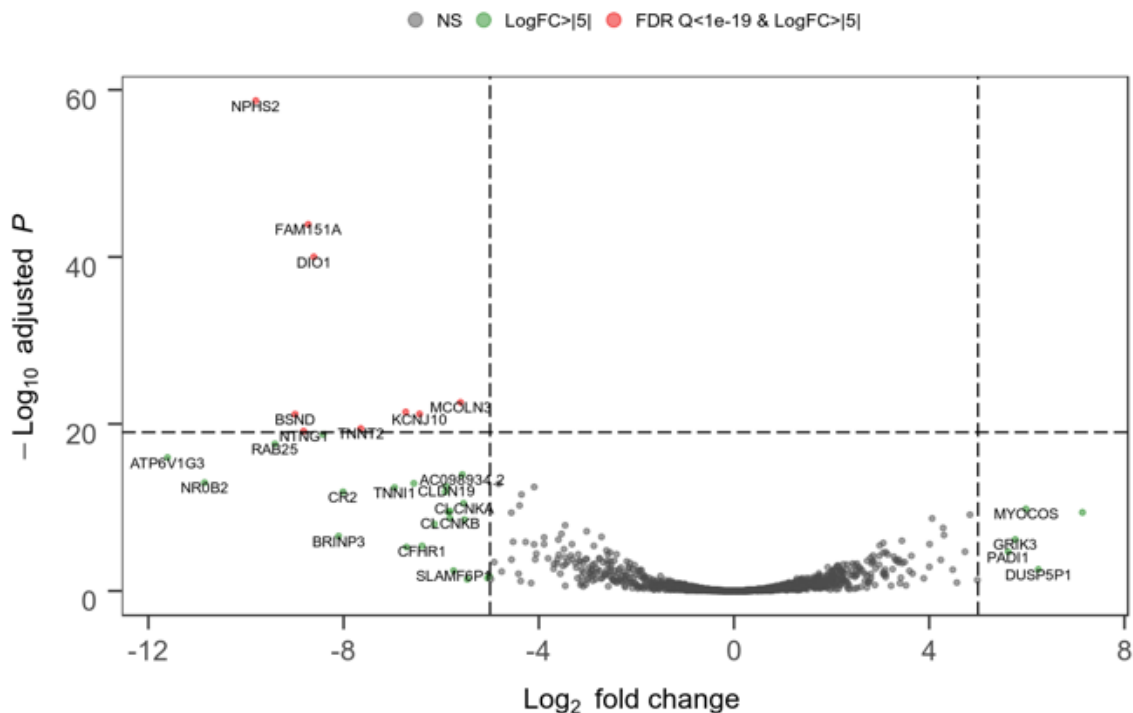
Volcano plot

Surprisingly, it seems that **DESeq2** does not provide function for a *Volcano plot*. However, a quick and nice *Volcano plot* can be generated from **DESeq2** results using a function, available at this URL:

<https://github.com/kevinblighe/EnhancedVolcano>

```
# Load plotting function
source("EnhancedVolcanoDESeq2.R")

# Make plot
EnhancedVolcanoDESeq2(result.df, AdjustedCutoff=10E-20, LabellingCutoff=0.05, FCCutoff=5.0)
```



Focused analysis of top significant genes

Finally, a focused analysis of the most differentially expressed genes may give clues for selecting candidates for downstream analyses, e.g. for inclusion into the gene expression signatures.

Select genes by p-value and fold-change

```
# Set cut-offs for top genes
log_fc_cutoff <- 5
adj_p_cutoff <- 0.001

# Extract results for top genes
top_genes.df <- result.df %>%
  filter(abs(log2FoldChange) > log_fc_cutoff, padj < adj_p_cutoff)
```

```
# Print top genes
```

```
top_genes.df
```

```
##          gene_id  gene_name  baseMean log2FoldChange      padj
## 1  ENSG00000116218      NPHS2  3498.49324      -9.798823  1.979669e-59
## 2  ENSG00000162391  FAM151A 14583.88874      -8.725494  1.313300e-44
## 3  ENSG00000211452      DIO1  5840.01659      -8.609517  1.012498e-40
## 4  ENSG00000055732  MCOLN3   389.44640      -5.605388  2.719366e-23
## 5  ENSG00000134258  VTCN1   397.49580      -6.727450  3.816724e-22
## 6  ENSG00000177807  KCNJ10   784.28326      -6.443739  6.568008e-22
## 7  ENSG00000162399      BSND   383.66461      -8.993231  7.019311e-22
## 8  ENSG00000118194  TNNT2   727.26329      -7.649031  3.646336e-20
## 9  ENSG00000162631  NTNG1   301.16639      -8.819620  7.587532e-20
## 10 ENSG00000116183  PAPP2   661.52114      -8.420240  1.850836e-19
## 11 ENSG00000132698  RAB25   323.74983      -9.406230  2.472867e-18
## 12 ENSG00000151418  ATP6V1G3 592.94906     -11.606765  1.084382e-16
## 13 ENSG00000234996 AC098934.2 166.42232      -5.568136  1.154303e-14
## 14 ENSG00000131910  NROB2   158.85090     -10.848179  1.082700e-13
## 15 ENSG00000132855  ANGPTL3 3708.79307      -6.560491  1.321762e-13
## 16 ENSG00000164007  CLDN19   252.88165      -5.898905  2.915436e-13
## 17 ENSG00000159173  TNNI1   640.02290      -6.955767  4.109309e-13
## 18 ENSG00000162896  PIGR   8287.39239      -5.916263  1.370312e-12
## 19 ENSG00000117322      CR2   398.51388      -8.011660  1.370312e-12
## 20 ENSG00000186510  CLCNKA   670.13466      -5.542289  3.204980e-11
## 21 ENSG00000283683  MYOCOS   357.52627       5.980955  1.540559e-10
## 22 ENSG00000143001  TMEM61    85.71740      -5.819244  2.645604e-10
## 23 ENSG00000182901  RGS7   157.89409      -5.851889  3.778038e-10
## 24 ENSG00000117598  PLPPR5   307.45508       7.138759  4.043817e-10
## 25 ENSG00000184908  CLCNKB  3390.11996      -5.820860  2.077361e-09
## 26 ENSG00000158014  SLC30A2 1230.27754      -5.522122  2.932829e-09
## 27 ENSG00000117601  SERPINC1  62.64627      -6.137652  9.784109e-09
## 28 ENSG00000162670  BRINP3   39.27424      -8.101568  2.771871e-07
## 29 ENSG00000163873  GRIK3   160.06775       5.762889  7.044498e-07
## 30 ENSG00000244414  CFHR1   41.03793      -6.392451  4.024608e-06
## 31 ENSG00000236136  ADORA2BP1 31.90155      -6.710285  5.748792e-06
## 32 ENSG00000142623  PADI1   40.22486       5.619188  2.062910e-05
## 33 ENSG00000159166  LAD1  1725.02072      -5.077473  2.594881e-05
```

Heatmap of top significant genes

```
# Get expression values for the top genes
```

```
top_genes_exprs.mx <- expr.mx[top_genes.df$gene_id, ]
```

```
# Change gene IDs to gene names
```

```
rownames(top_genes_exprs.mx) <- top_genes.df$gene_name
```

```
# Plot heatmap
```

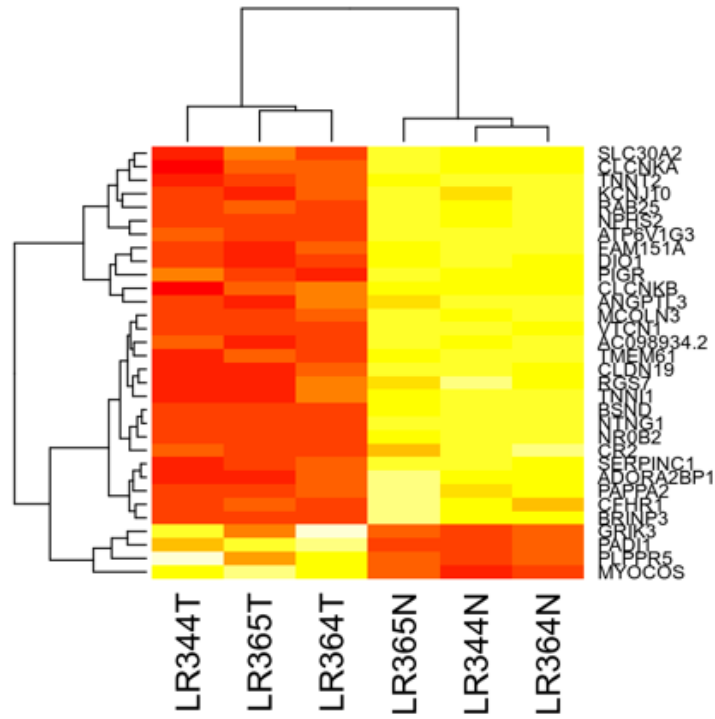
```
num_genes <- nrow(top_genes.df)
```

```
main <- paste(num_genes, "genes with log2(FC) >", log_fc_cutoff, "and adj.P <", adj_p_cutoff, "\n")
```

```
par(oma=c(0,0,3,0)) # make space for the main header
```

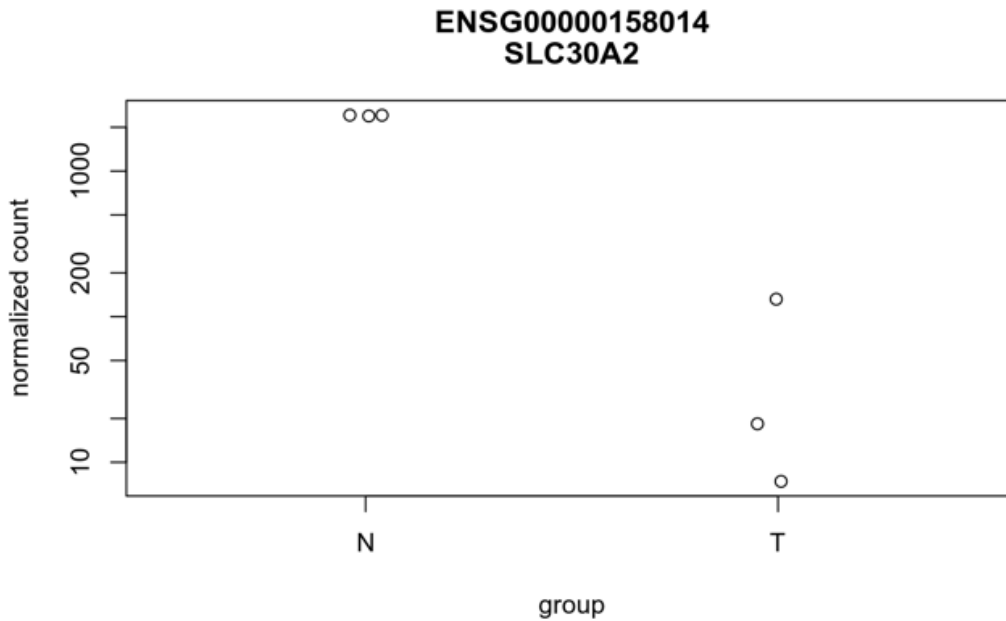
```
heatmap(top_genes_exprs.mx, main=main)
```

32 genes with $\log_2(\text{FC}) > 5$ and $\text{adj.P} < 0.001$



Plot normalised counts for a selected gene

```
gene_id <- "ENSG00000158014"
gene_name <- geneId2Name.df[gene_id,]
plotCounts(dds, gene=gene_id, main=gene_name)
```



5. STAR-fusion

Detecting fusion transcripts

Fusion transcripts are emerging as clinically important targets in oncology. STAR-Fusion allows fast and accurate fusion detection from RNA-Seq data:

<https://www.biorxiv.org/content/early/2017/03/24/120295>

<https://github.com/STAR-Fusion/STAR-Fusion/wiki>

STAR-fusion may use BAM files generated by STAR. However, in most cases it is preferable to run STAR-fusion starting from the original FASTQ files. In this case, STAR-fusion runs STAR-alignment before calling fusions anyway. However, it runs STAR-alignment with the parameters optimised for fusion detection.

STAR fusion launcher is implemented as a Python script. It requires STAR, samtools and CTAT genome library. A STAR-fusion analysis, starting from FASTQ files may be launched like this:

```
STAR-Fusion \  
--genome_lib_dir "${ctat_lib_folder}" \  
--left_fq "${fastq1}" \  
--right_fq "${fastq2}" \  
--CPU 12 \  
--output_dir "${out_folder}"
```

CPU parameter describes the number of threads, requested for STAR-aligner. Again, the full script is provided in the scripts folder (`star_fus.sh`). Review this script and run STAR-fusion. The run may take 5-10 min, depending on the available RAM and number of requested threads.

After the run, explore the content of the output folder. Amongst the other files it should contain `star-fusion.fusion_predictions.tsv` and `star-fusion.fusion_predictions.abridged.tsv` files. Open the last file in LibreOffice Calc or Excel (dont make or save any changes - this file will be used in downstream analysis!).

What fusions are supported mainly by split reads, rather than by spanning fragments? What fusion has clinically relevant annotations? What fusion drives cancer growth in SRR3534924 sample?

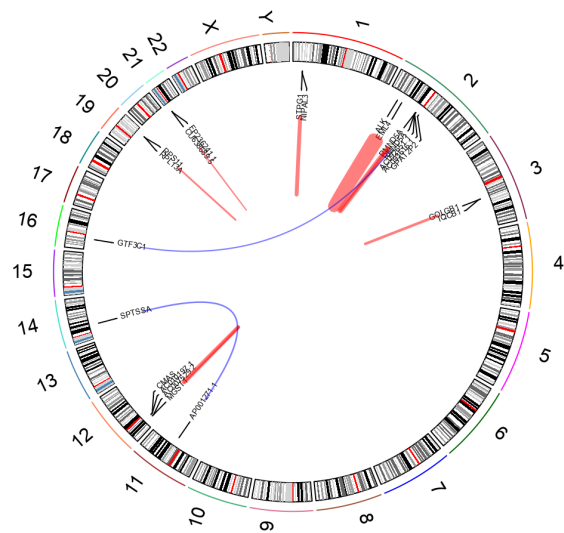
Visualise fusion transcripts using *Chimeraviz*

Chimeraviz is a simple R tool to visualise results of STAR-Fusion:

<https://academic.oup.com/bioinformatics/article/33/18/2954/3835381>

Five lines of code produce simple report, which includes a circus-plot and a table of fusions:

```
# Install chimeraviz library (done once)  
source("https://bioconductor.org/biocLite.R")  
biocLite("chimeraviz")  
  
# Load chimeraviz library  
library(chimeraviz)  
  
# Read STAR-fusion output  
fusions <- importStarfusion("star-fusion.fusion_predictions.abridged.tsv", "hg38", 20)  
  
# Generate quick report with circus plot  
createFusionReport(fusions, "chimeraviz_report.html")
```



If you have free time at the end of this session: try making *Chimeraviz* gene-pair plot for the driver-fusion.