

Example RNA-seq report

February 8, 2018

Issue number: #9999
Request by: Jan User <jan.user@ki.se>
Principal Investigator: Maria Investigator <maria.investigator@ki.se>
Organisation: Karolinska Insitutet
NBIS staff: Olga Dethlefsen <olga.dethlefsen@nbis.se>

Contents

1	Support request	3
2	Important practical information	3
2.1	Data responsibility	3
2.2	Acknowledgments	3
3	Work log	3
4	Materials and Methods	4
4.1	Available data	4
4.2	Data processing	4
4.3	Differential expression	4
4.4	Exon usage	5
4.5	HOMER Motif Analysis	5
5	Results	6
5.1	Data exploration	6
5.2	Differential expression	8
5.2.1	Estimating BCVs	8
5.2.2	GR1: time2 vs. time1	9
5.2.3	GR2: time2 vs. time1	11
6	Deliverable	13
7	R session info	14
8	Where to go next	14
9	Support project closing procedures	14

1 Support request

To answer our question we have performed RNA-seq, total RNA ribosomal RNA depleted. We had 4 groups and would like help to finding differentially expressed genes

2 Important practical information

2.1 Data responsibility

NBIS & Uppnexus Unfortunately, we do not have resources to keep any files associated with the support request. We kindly suggest that you store safely the results delivered by us. In addition, we kindly ask that you remove the files from UPPMAX/UPPNEX. The main storage at UPPNEX is optimized for high-speed and parallel access, which makes it expensive and not the right place for longer time archiving. Please consider others by not taking up the expensive space.

Long-term backup The responsibility for data archiving lies with universities and we recommend asking your local IT for support with long-term data archiving. Also a newly established **Data Office** at SciLifeLab may be of help to discuss other options.

Data submission NBIS can help with raw sequencing data submission to ENA. Please contact our data manager, Niclas Jareborg, niclas.jareborg@nbis.se for more details.

2.2 Acknowledgments

If you are presenting the results in a paper, at a workshop or conference, we kindly ask you to acknowledge us.

NBIS staff are encouraged to be co-authors when this is merited in accordance to the ethical recommendations for authorship, e.g. **ICMJE recommendations**. If applicable, please include **Olga Dethlefsen, National Bioinformatics Infrastructure Sweden, Science for Life Laboratory, Stockholm University** as co-author. In other cases, NBIS would be grateful if support by us is acknowledged in publications according to this example: "Support by NBIS (National Bioinformatics Infrastructure Sweden) is gratefully acknowledged."

Uppmax kindly asks you to acknowledge UPPMAX and SNIC. If applicable, please add: **The computations were performed on resources provided by SNIC through Uppsala Multidisciplinary Center for Advanced Computational Science (UPPMAX) under Project b29999.**

NGI Stockholm In any and all publications based on data from NGI Sweden, the authors must acknowledge SciLifeLab, NGI and Uppmax, like so: **The authors would like to acknowledge support from Science for Life Laboratory, the National Genomics Infrastructure, NGI, and Uppmax for providing assistance in massive parallel sequencing and computational infrastructure.**

3 Work log

A brief project history containing key points

2015-09-15 first meeting with Jan to discuss experimental design, available data and desired results. As first results, Jan would like to receive lists of differentially expressed (DE) genes between the two time points for the two groups

2015-10-09 meeting with Jan to go over the DE results. Agreed that Jan will go over the DE results and try running gene set enrichment analyses using DAVID website.

2015-11-06 I have run and emailed Jan the exon usage results for the 4 comparisons

2015-12-01 I have run and emailed Jan the motif discovery Homer results

2015-12-10 meeting with Jan to discuss the additional

2015-12-20 final results and report delivered, project is closed

4 Materials and Methods

4.1 Available data

Data were delivered to Inbox on Uppnexus b29999 in fastq format using Illumina 1.8 quality scores. Data were from a paired-end run, with one file for the forward reads and one file for the reverse reads following a naming convention: [LANE-][DATE-][FLOWCELL-][SCILIFE NAME-][READ].fastq.gz.

4.2 Data processing

Raw sequencing reads were processed to obtain counts per genes for each samples. This included:

1. FastQC/0.11.2 [Andrews 2010] quality check on raw sequencing reads
2. trimmomatic/0.32 [Bolger, Lohse, and Usadel 2014] reads filtering for quality score and read length. Reads with average quality below 20 (within 4-base wide sliding window) and/or shorter than 36 bases were removed
3. star/2.4.1c [Dobin et al. 2013] was used to align the reads to the reference genome `Mus_musculus.GRCm38.dna.primary_assembly.fa` using the annotation `Mus_musculus.GRCm38.81.gtf`, with reference genome and annotation downloaded from <http://www.ensembl.org/index.html>
4. featureCounts/1.5.0 from subread/1.4.5 Liao, Smyth, and Shi 2013] was used to count the fragments in the exon regions as defined in the `Mus_musculus.GRCm38.81.gtf` file, using default parameters. Specifically, for paired-end reads, a fragment is said to overlap a feature if at least one read base is found to overlap the feature. Fragments overlapping with more than one feature and multi-mapping reads are not counted
5. Counts from multiple lanes were added, if applicable
6. samtools/0.1.19 [Li et al. 2009] were used to sort and index the BAM files containing the aligned reads, e.g. for visualization in IGV genome browser
7. MultiQC/0.3.1 [Ewels et al. 2016] was used to aggregate results from FastQC/0.11.2, star/2.4.1c and featureCounts/1.5.0 across many samples into a single report

4.3 Differential expression

All analyses were performed under R, a programming language and software environment for statistical computing and graphics. Details on the R version and packages used can found at the end of this document in [R session info](#)

1. biomaRt package was used to annotate Ensembl gene identifiers with chromosome name, official gene symbol and description.

2. low count reads were filtered by keeping reads with at least 1 read per million in at least 2 samples
3. **edgeR** package was used to normalize for the RNA composition by finding a set of scaling factors for the library sizes that minimize the log-fold changes between the samples for most genes, using a trimmed mean of M values (TMM) between each pair of samples.
4. the normalized counts were used to examine the samples for outliers and relationships, using Multidimensional Scaling and heatmap based on the Pearson correlation coefficient between every sample pair
5. the normalized counts were used to examine the samples for outliers and relationships, using Multidimensional Scaling and heatmap based on the Pearson correlation coefficient between every sample pair
6. **edgeR** package was to define design matrix based on the experimental design, fitting gene-wise glms model and conducting likelihood ratio tests for the selected group comparisons

4.4 Exon usage

1. **DEXseq** package in R was used to infer differential exon usage
2. the provided with the **DEXseq** package **Python** scripts were used to prepare a flattened GTF file based on the `Mus_musculus.GRCm38.81.gtf` and to obtain counts per each exon given the aligned BAM files
3. size factors measuring the relative sequencing depth were estimated to adjust for coverage biases
4. variability of the data was then estimated to be able to distinguish technical and biological variation from real effects on exon usage due to the different conditions. Briefly, per-exon dispersions are calculated using a Cox-Reid adjusted profile likelihood estimation, then a dispersion-mean relation is fitted to this individual dispersion values and finally, the fitted values are taken as a prior in order to shrink the per-exon estimates towards the fitted values
5. having the dispersion estimates and the size factors, differential exon usage was tested. For each gene, **DEXSeq** fits a generalized linear model with the formula $\sim \text{samples} + \text{exon} + \text{condition} : \text{exon}$ and compares it to the null model $\sim \text{samples} + \text{exon}$. The deviances of both fits are compared using a χ^2 -distribution, giving rise to a p value, indicative whether the null model is sufficient to explain the data or whether it may be rejected in favor of the alternative containing an interaction coefficient for condition:exon. The latter means that the fraction of the gene's reads that fall onto the exon under the test differs significantly between the experimental conditions.
6. the obtained p-values were BH adjusted for multiple comparison

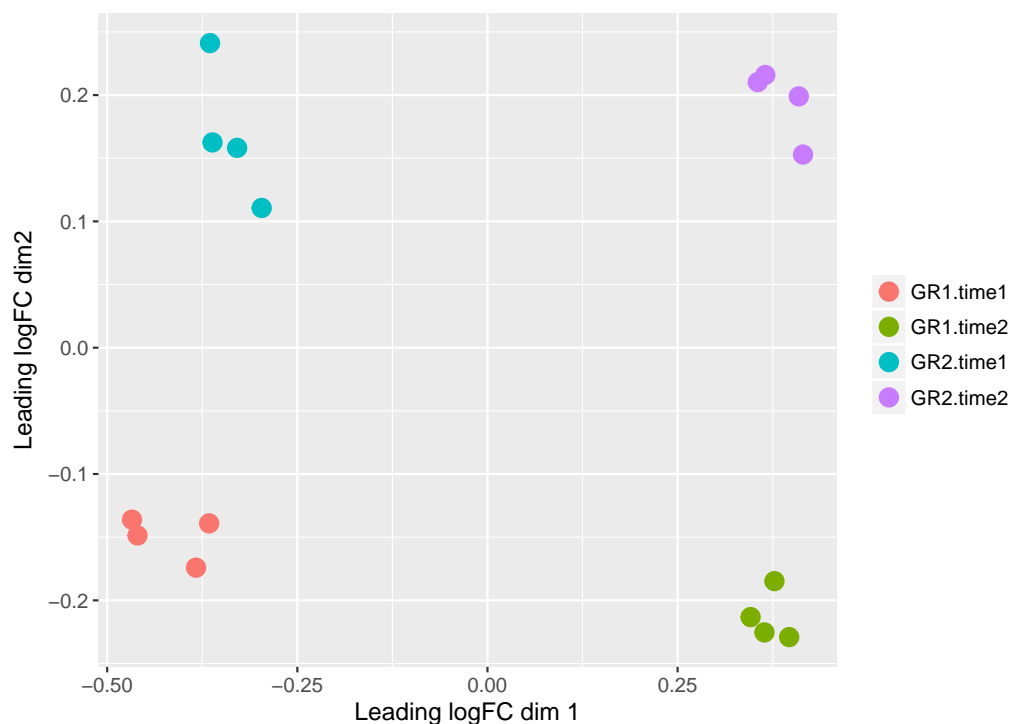
4.5 HOMER Motif Analysis

1. **homer/4.7.2** was used to analyze the promoters of genes and look for motifs that are enriched in the target gene promoters relative to other promoters. The analyses followed the '[Analyzing lists of genes with promoter motif analysis](#)' tutorial
2. briefly, for each comparison, the analyses were run for the differentially expressed genes, separately for down- and up-regulated genes, were differentially expressed genes were defined at 5% FDR and absolute minimum log 2 fold change of 1.
3. the analyses included Gene Ontology enrichment calculations, *de novo* motif analysis and known motif enrichment analysis

5 Results

5.1 Data exploration

Table count, containing counts measured across many genes and samples, is a typical example of a multidimensional dataset, where N objects (samples) were measured on p numeric variables. Hence, to examine the samples for outliers and other relationship one can look at the several multivariate techniques that aim to reduce the dataset dimensions and to reveal the data structure by plotting samples in one or two dimensions. Multi-dimensional Scaling provides a visual representation of the pattern of proximities among a set of objects, here samples. Another method of visualization is a heatmap based on the Pearson correlation coefficient, calculated between every sample pair.



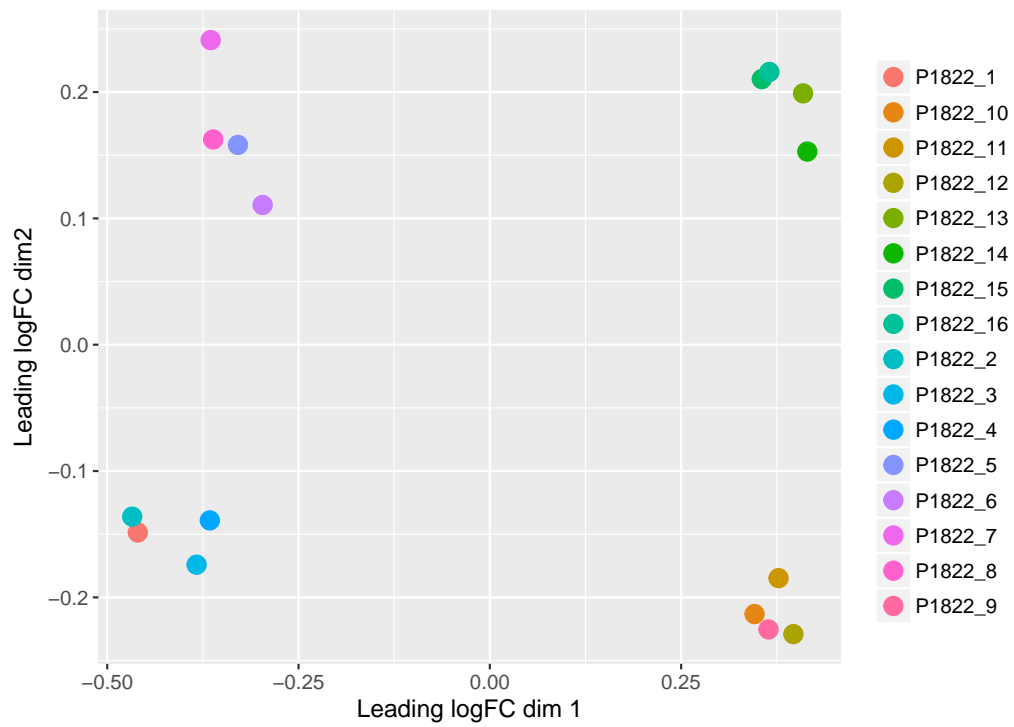


Figure 2: MDS plot on the normalized and filtered counts colour-coded by individual samples

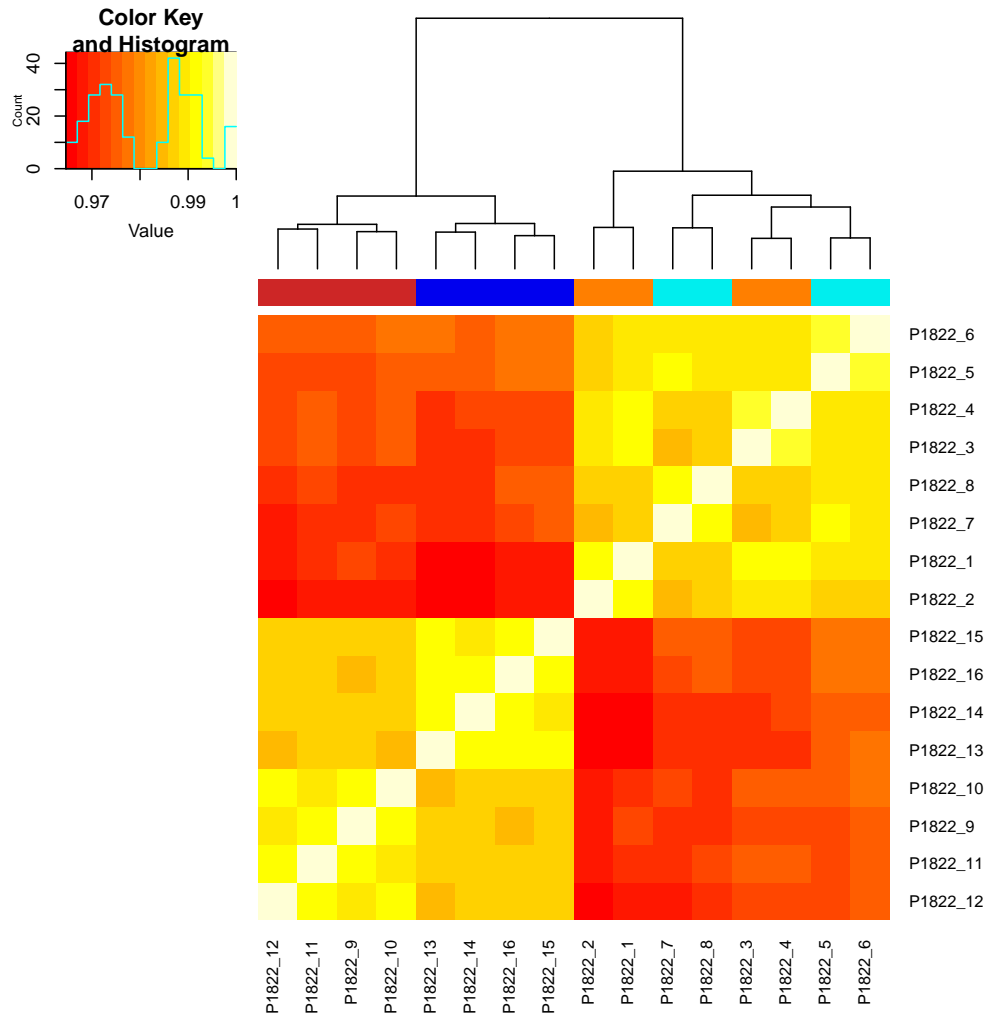


Figure 3: Heatmap based on the pair-wise Pearson correlation coefficient between samples

5.2 Differential expression

5.2.1 Estimating BCVs

Two levels of variation, technical and biological, can be distinguished in any RNA-Seq experiment. Biological coefficient of variation (BCV) is the coefficient of variation with which the (unknown) true abundance of the gene varies between replicate RNA samples. It represents the CV that would remain between biological replicates if sequencing depth could be increased indefinitely. The technical CV decreases as the size of the counts increases. BCV on the other hand does not. BCV is therefore likely to be the dominant source of uncertainty for high-count genes, so reliable estimation of BCV is crucial for realistic assessment of differential expression in RNA-Seq experiments. **edgeR** uses empirical Bayes methods that permit the estimation of gene-specific biological variation, even for experiments with minimal levels of biological replication.

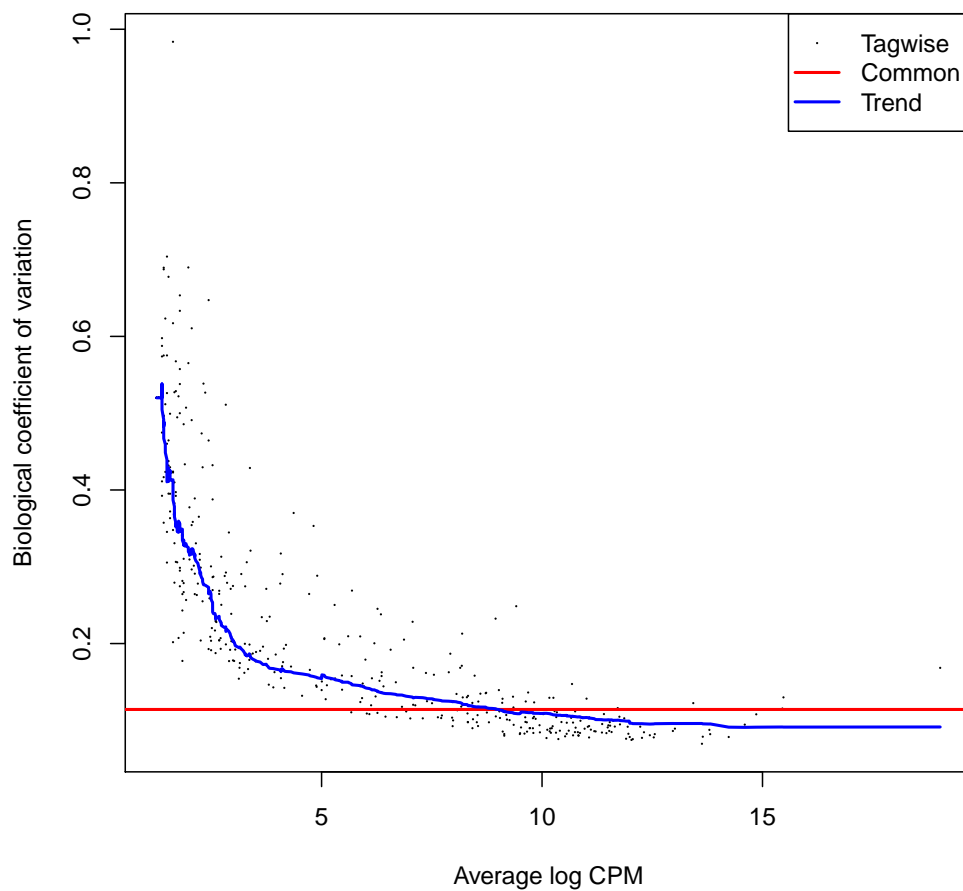


Figure 4: Biological coefficient of variation plot showing the dispersion estimates

5.2.2 GR1: time2 vs. time1

	Down-regulated	Non-significant	Up-regulated
Genes	40	363	36

Table 1: Number of down- and up-regulated differentially expressed genes given 5% FDR and absolute minimum log-fold-change of 1

ensembl_gene_id	mgi_symbol	GR1.logFC	GR1.FDR
ENSMUSG00000056870	Gulp1	2.33	2.62E-109
ENSMUSG00000025911	Adhfe1	2.87	1.03E-81
ENSMUSG00000026064	Ptp4a1	-1.77	6.51E-75
ENSMUSG00000041859	Mcm3	-1.88	1.18E-63
ENSMUSG00000064294	Aox3	4.29	6.81E-53
ENSMUSG00000026023	Cdk15	3.67	5.07E-46
ENSMUSG00000038305	Spats2l	-1.61	5.58E-44
ENSMUSG00000026069	Il1rl1	-1.63	5.86E-40
ENSMUSG00000026070	Il18r1	-2.63	4.20E-37
ENSMUSG00000051951	Xkr4	3.28	1.16E-31

Table 2: 10 genes with the smallest FDR values

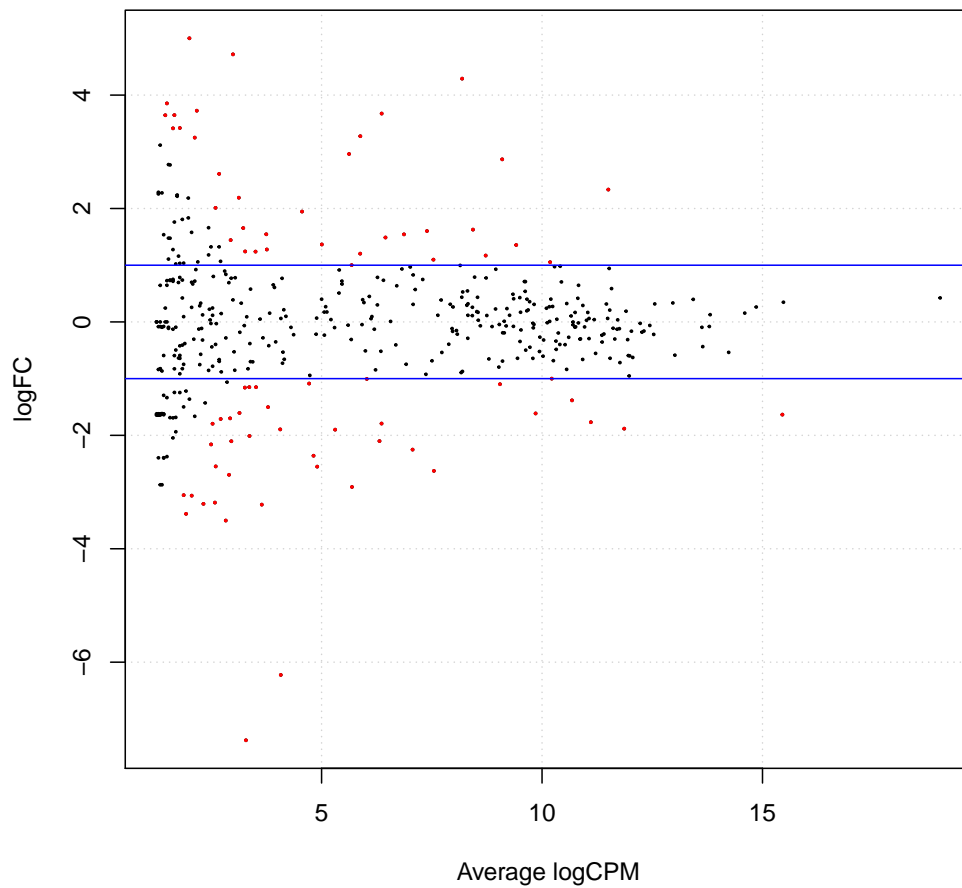


Figure 5: Plot log-fold change against log-counts per million, with DE genes highlighted

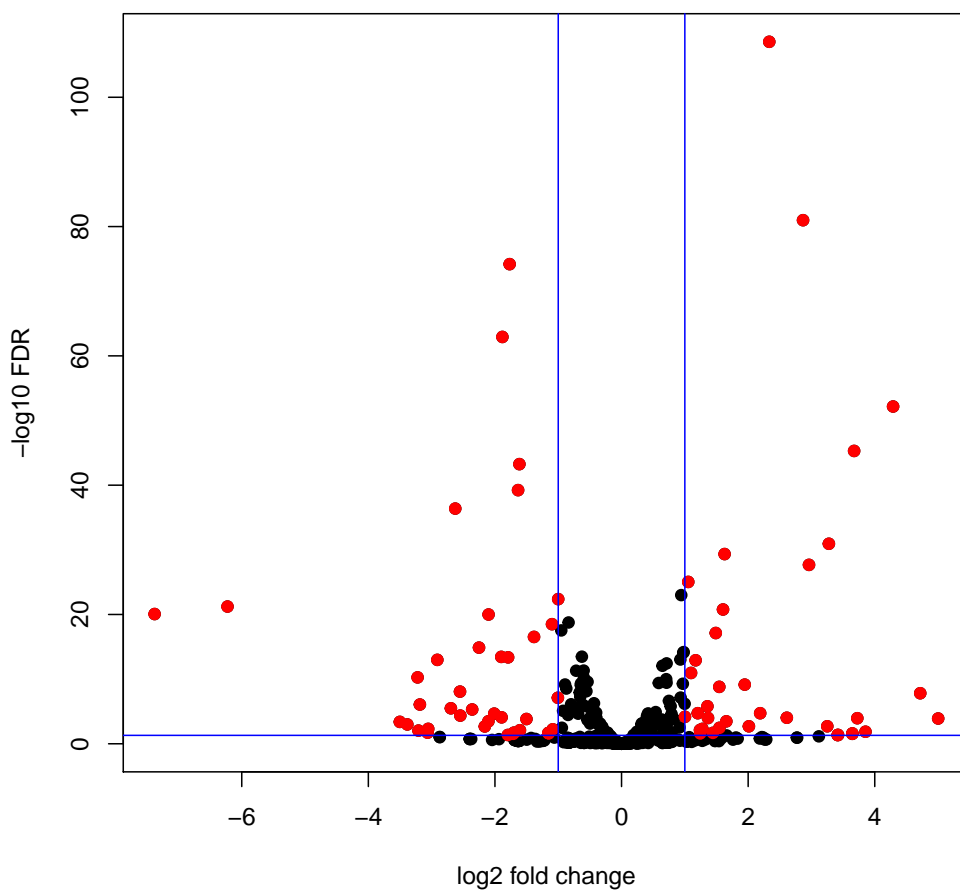


Figure 6: Volcano plot with DE genes highlighted in red. Horizontal blue line corresponds to $FDR=0.05$ and vertical blue lines correspond to absolute \log_2 fold change of 1

5.2.3 GR2: time2 vs. time1

	Down-regulated	Non-significant	Up-regulated
Genes	36	368	35

Table 3: Number of down- and up-regulated differentially expressed genes given 5% FDR and absolute minimum log-fold-change of 1

ensembl_gene_id	mgi_symbol	GR2.logFC	GR2.FDR
ENSMUSG00000025993	Slc40a1	3.20	2.12E-128
ENSMUSG00000056870	Gulp1	2.34	2.36E-112
ENSMUSG00000026069	Il1rl1	-2.14	1.49E-65
ENSMUSG00000051951	Xkr4	3.88	1.72E-58
ENSMUSG00000025911	Adhfe1	2.31	4.48E-56
ENSMUSG00000026023	Cdk15	4.07	4.82E-47
ENSMUSG00000064294	Aox3	3.65	5.34E-43
ENSMUSG00000026070	Il18r1	-2.62	1.14E-38
ENSMUSG00000041859	Mcm3	-1.35	5.18E-34
ENSMUSG00000026064	Ptp4a1	-0.90	4.66E-20

Table 4: 10 genes with the smallest FDR values

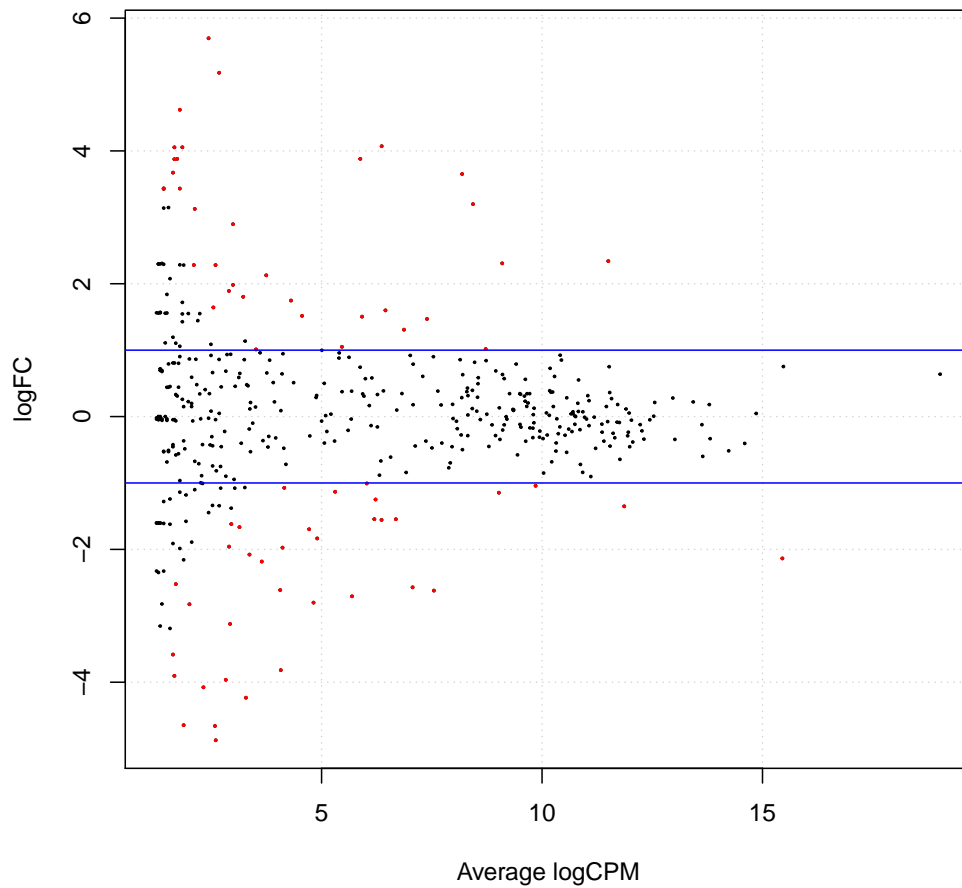


Figure 7: Plot log-fold change against log-counts per million, with DE genes highlighted

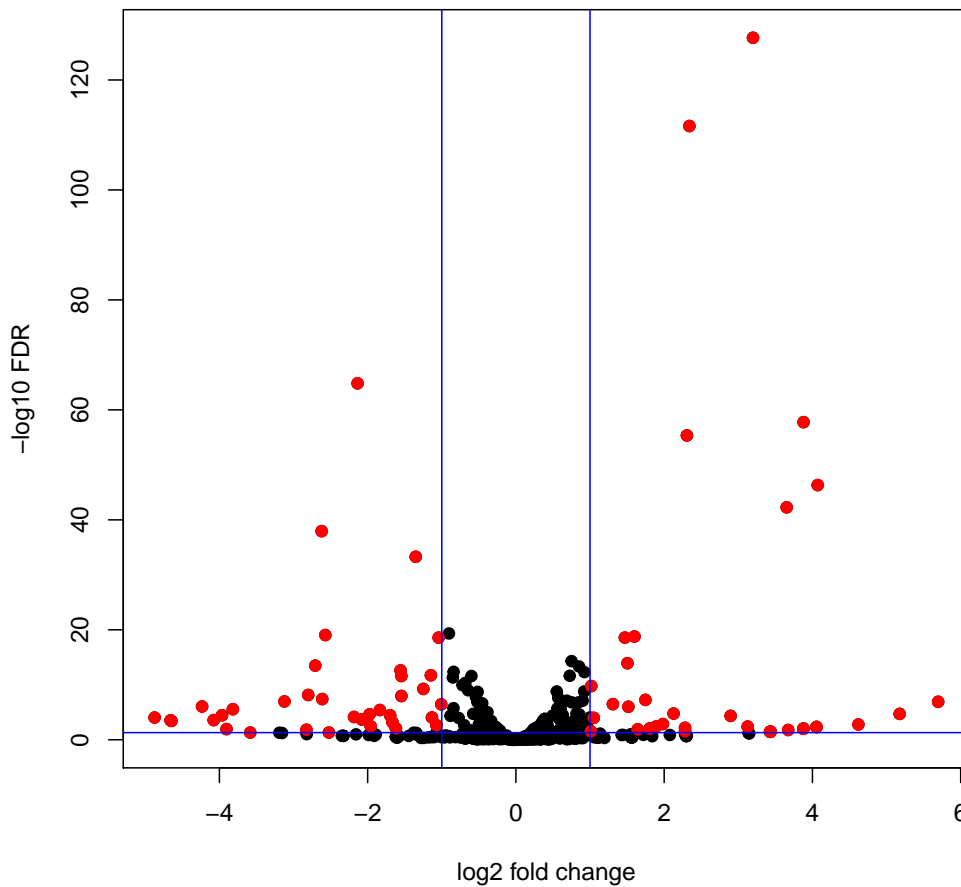


Figure 8: Volcano plot with DE genes highlighted in red. Horizontal blue line corresponds to $FDR=0.05$ and vertical blue lines correspond to absolute \log_2 fold change of 1

6 Deliverable

Below is the list of key tab-delimited text files containing the key results from the described data analyses. All results incl. computational scripts can be found on Uppnex under project **b29999**.

```
[1] "DE.txt" "count_table.txt"
[3] "count_table_annotations.txt" "norm_table.txt"
[5] "norm_table_annotations.txt"
```

where,

DE.txt contains the differential expression results for all the analyses

count_table.txt contains the raw genes counts

count_table_annotations.txt contains the annotations for the genes in the count_table.txt

norm_table.txt contains filtered and normalized genes expression values (TMM)

norm_table_annotations.txt contains the annotations for the norm_table.txt

7 R session info

```
R version 3.3.3 (2017-03-06)
Platform: x86_64-apple-darwin13.4.0 (64-bit)
Running under: macOS 10.13.3

locale:
[1] C

attached base packages:
[1] stats      graphics  grDevices  utils      datasets  methods   base

other attached packages:
[1] xtable_1.8-2   biomaRt_2.30.0 gplots_3.0.1  ggplot2_2.2.1 edgeR_3.16.5
[6] limma_3.30.13

loaded via a namespace (and not attached):
 [1] Rcpp_0.12.14      pillar_1.1.0      plyr_1.8.4
 [4] bitops_1.0-6      tools_3.3.3       digest_0.6.14
 [7] bit_1.1-12        RSQLite_2.0       memoise_1.1.0
[10] tibble_1.4.1      gtable_0.2.0      lattice_0.20-35
[13] rlang_0.1.6       DBI_0.7           parallel_3.3.3
[16] S4Vectors_0.12.2 gtools_3.5.0      caTools_1.17.1
[19] IRanges_2.8.2     locfit_1.5-9.1    stats4_3.3.3
[22] bit64_0.9-7       grid_3.3.3        Biobase_2.34.0
[25] AnnotationDbi_1.36.2 XML_3.98-1.9      gdata_2.18.0
[28] blob_1.1.0        scales_0.5.0      BiocGenerics_0.20.0
[31] splines_3.3.3     colorspace_1.3-2  labeling_0.3
[34] KernSmooth_2.23-15 RCurl_1.95-4.10   lazyeval_0.2.1
[37] munsell_0.4.3
```

8 Where to go next

There is a wide selection of online user-friendly tools available for investigating the interesting genes or list of DE genes. Few recommended below

ClustVist for creating Principal Component Analysis plots and heatmaps

Venny helps to prepare Venn diagram showing relations between a finite collection of different sets, e.g. between list of DE genes from different comparisons

DAVID for a comprehensive set of functional annotation tools for investigators to understand biological meaning behind large list of genes, including identification of enriched Gene Ontology terms, discovering enriched functional-related gene groups, visualizing genes on BioCarta & KEGG pathway and many more

REVIGO for summarizing long list of Gene Ontology terms and visualization in semantic similarity-based scatter plots, interactive graphs or tag clouds

9 Support project closing procedures

You should soon be contacted by one of our managers, Jessica Lindvall <jessica.lindvall@nbis.se> or Henrik Lantz <henrik.lantz@nbis.se>, with a request to close down the project in our internal system and for invoicing matters. If we do not hear from you within **30 days** the project will be automatically closed and invoice sent.

Again, we would like to remind you about data responsibility and acknowledgements, see [Data responsibility](#) and [Acknowledgments](#).

You are naturally more than welcome to come back to us with further data analysis request at any time via <http://nbis.se/support/support.html>. Thank you for using NBIS and all the best for future research.

References

- Andrews, Simon. 2010. *FastQC*. <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>.
- Bolger, Anthony M, Marc Lohse, and Bjoern Usadel. 2014. “Trimmomatic: A flexible trimmer for Illumina sequence data”. *Bioinformatics* 30:2114.
- Dobin, Alexander, et al. 2013. “STAR: ultrafast universal RNA-seq aligner.” *Bioinformatics (Oxford, England)* 29 (1): 15–21. ISSN: 1367-4811 (Electronic). doi:[10.1093/bioinformatics/bts635](https://doi.org/10.1093/bioinformatics/bts635).
- Liao, Yang, Gordon K Smyth, and Wei Shi. 2013. “The Subread aligner: fast, accurate and scalable read mapping by seed-and-vote”. *Nucleic Acids Research* 41 (10): e108–e108. ISSN: 0305-1048. doi:[10.1093/nar/gkt214](https://doi.org/10.1093/nar/gkt214). <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3664803/>.
- Li, Heng, et al. 2009. “The Sequence Alignment/Map format and SAMtools”. *Bioinformatics* 25 (16): 2078–2079.
- Ewels, Philip, et al. 2016. “MultiQC: summarize analysis results for multiple tools and samples in a single report”. *Bioinformatics* 32 (19): 3047–3048. ISSN: 1367-4803. doi:[10.1093/bioinformatics/btw354](https://doi.org/10.1093/bioinformatics/btw354). <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC5039924/>.