

Task 3: RNA-seq Differential Gene Expression & RNA-seq Functional Profiling

Simon Safron [Mn: 11923407], Alexander Veith [Mn: 12122739], Miriam Überbacher [Mn:01627576]

2025-12-06

Contents

1	Introduction	2
2	Raw Reads and Mapping QC	2
2.1	FastQC:	2
2.2	HTSeq-count:	3
3	Dataset	4
3.1	Raw read counts	4
3.2	Sample information	4
4	Preprocessing of the data	5
4.1	Filtering of the data	5
5	DGE Analysis	5
5.1	Differential Expression Analysis	5
5.2	Extracting results	5
6	Vizualising data	8
6.1	Experimental QC - Clustering of samples (PCA)	8
6.2	Viewing counts for a single geneID	9
	References	11

1 Introduction

2 Raw Reads and Mapping QC

2.1 FastQC:

- a) According to FastQC: What was the minimum and the maximum number of read pairs sequenced per sample?

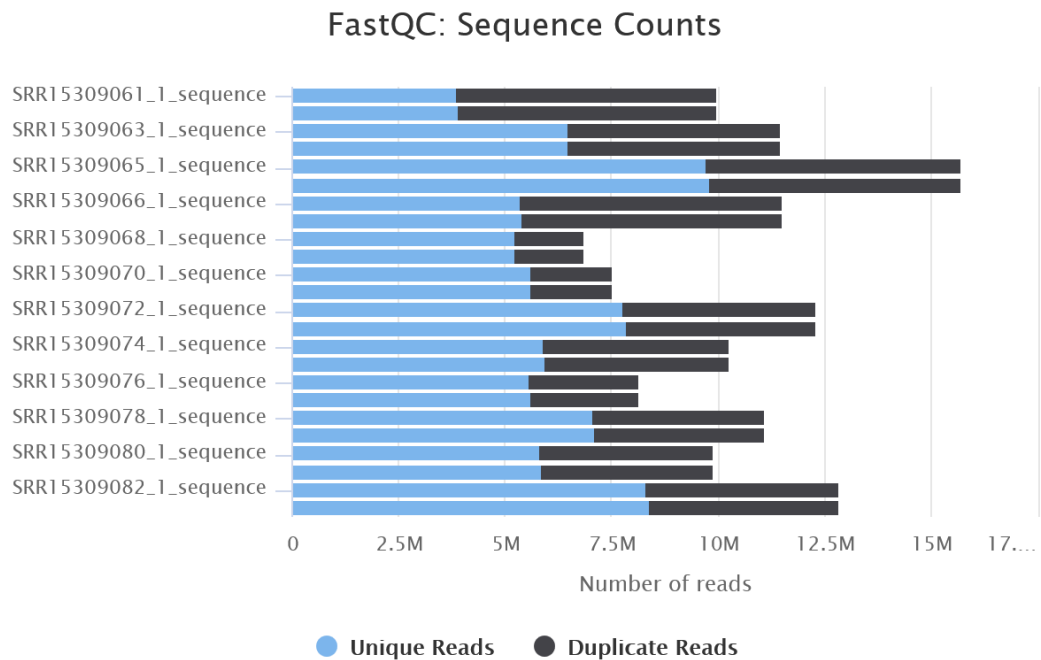


Figure 1: Barplot of the number of read pairs per sample.

Answer: As can be seen in Figure 1 the minimum number of reads per sample is around 6.8 million and the maximum number of reads per sample is around 15.7 million.

- b) What is the most overrepresented sequence (string of nucleotides) that was found by FastQC?

Answer: According to the MultiQC report the most overrepresented sequence was:

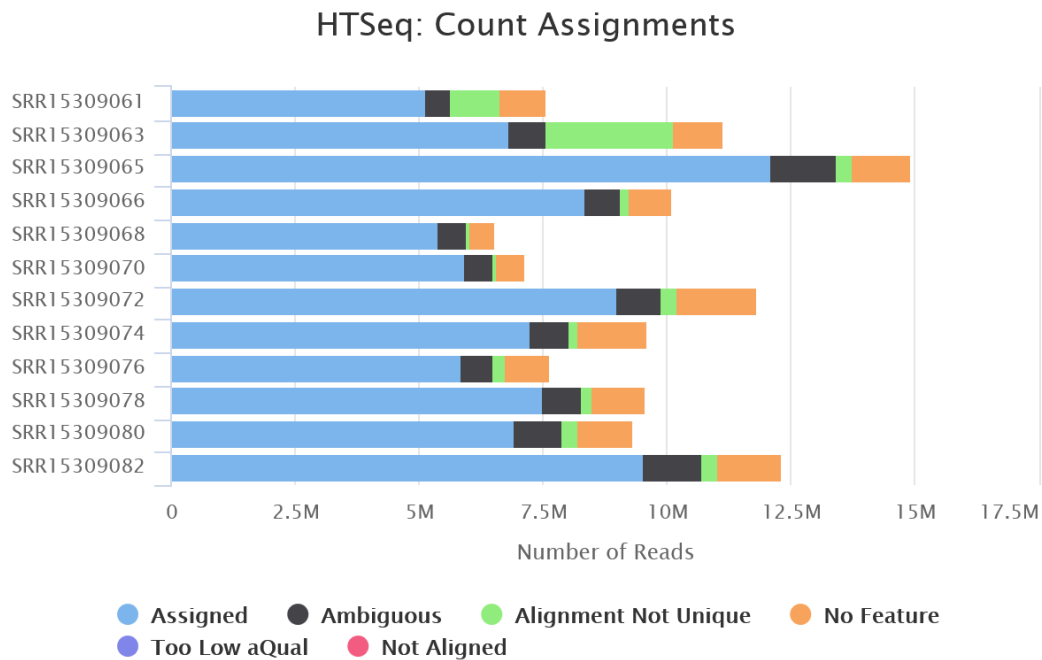
"AA"

- c) What might be the reason there is so much of this specific sequence/homopolymer ?

Answer: This could be from A tailed adapter dimers and PCR slippage products that outcompete genuine RNA fragments in the library.

2.2 HTSeq-count:

- d) According to HTSeq Count: What was the minimum and the maximum number of read pairs reported per sample?



Created with MultiQC

Figure 2: Barplot of the number of read pairs, per sample HTSeq.

Answer: As can be seen in Figure 2 the minimum number of reads per sample is around 5.9 million and the maximum number of reads per sample is around 14.8 million.

- f) Which was the minimum and maximum percentage of reads uniquely assigned to a gene, as reported by HTSeq-count?

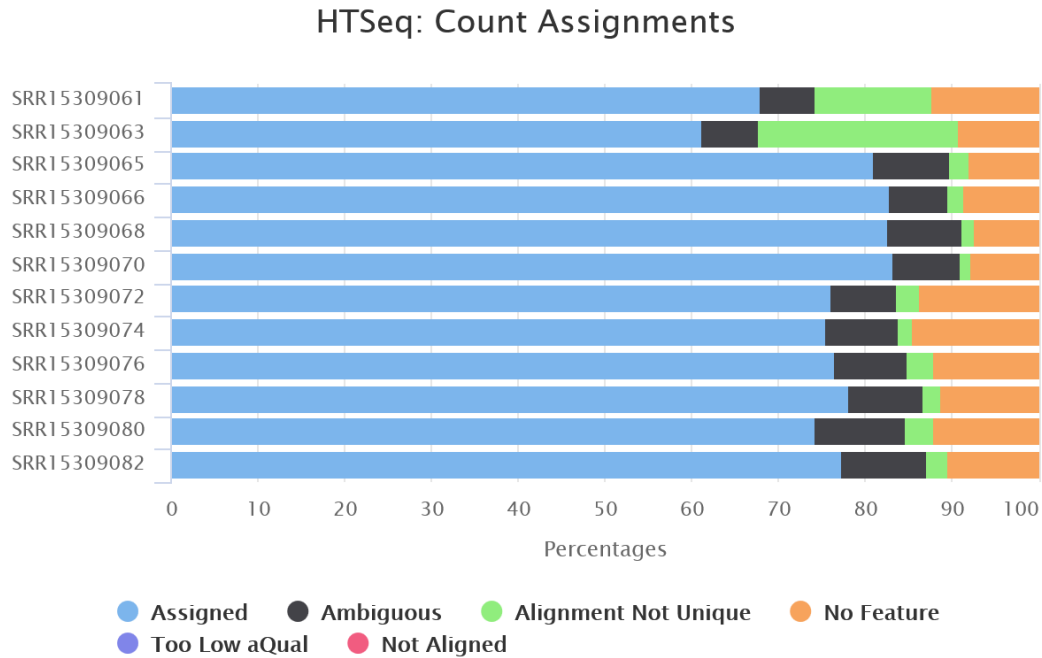


Figure 3: HTSeq assignment plot in percentages.

Answer: As can be seen in Figure 3 the minimum percentage of reads uniquely assigned to a gene is 61.3% and the maximum percentage of reads uniquely assigned to a gene is around 83.3%.

3 Dataset

3.1 Raw read counts

3.2 Sample information

- a) Which columns define the used base strain and substrain (WT or TGT mutant), respectively? Can you spot an error in one of those columns?

Answer: The genotype column shows if the wild type or tgt mutant was used. The strain column shows which *E.coli* strain was used. The error is in the strain column suggesting that for every experiment a different strain was used. But comparing this with the strains and method subsection of the methods section of the paper shows that only “*Escherichia coli* K-12 MG1655 was used as the WT strain.”[1]

- b) Which column defines if nickel was added to the media?

Answer: The treatment column defines if nickel was added to the media.

- c) Find the following information on the SRA Study page:

Which type of Illumina machine was used for sequencing?

Answer: Illumina HiSeq 2500

What was the library layout?

Answer: PAIRED

When was the data released?

Answer: 2022-05-16

4 Preprocessing of the data

4.1 Filtering of the data

- a) For how many genes did we originally retrieve count data?

```
dim(rawCounts)
#[1] 4295  12
```

Answer: Originally count data were retrieved for 4295 genes

- b) How many will be left after applying the filter?

```
dim(rawCounts[rowSums(rawCounts) > 10, ])
#[2] 4221  12
```

Answer: After applying the filter 4221 genes will be left.

5 DGE Analysis

5.1 Differential Expression Analysis

5.2 Extracting results

Interpreting the summary:

- a) How many genes are significantly up-regulated and how many are significantly down-regulated in the nickel treated WT strain as compared to the untreated WT strain, using the default cutoff for the adjusted p-value?

```

# Extract results with default alpha (0.1)
DESeq2Results_WT_nickel <- results(DESeq2Data,
                                   contrast = c("group", "WT.Nickel", "WT.none"))

# View summary
summary(DESeq2Results_WT_nickel)

# out of 4221 with nonzero total read count
# adjusted p-value < 0.1
# LFC > 0 (up)      : 1069, 25%
# LFC < 0 (down)    : 1063, 25%
# outliers [1]      : 6, 0.14%
# low counts [2]     : 0, 0%
# (mean count < 1)
# [1] see 'cooksCutoff' argument of ?results
# [2] see 'independentFiltering' argument of ?results

```

Answer: 1069 genes are significantly up-regulated and 1063 genes are significantly down-regulated in the nickel treated WT strain as compared to the untreated WT strain, using the default cutoff for the adjusted p-value of 10%.

- b) What is the standard cutoff used for the significance level (adjusted p-value), if we don't change it?

Answer: The standard cutoff used for the significance level (adjusted p-value), if we don't change it is 10% (p-value < 0.1).

- c) How many significantly differentially expressed genes does that make in total?

```

# Extract TGT.Nickel vs TGT.none with alpha = 0.1
DESeq2Results_TGT_nickel <- results(DESeq2Data,
                                   contrast = c("group", "TGT.Nickel", "TGT.none"),
                                   alpha = 0.1)

# View summary
summary(DESeq2Results_TGT_nickel)

# out of 4221 with nonzero total read count
# adjusted p-value < 0.1
# LFC > 0 (up)      : 986, 23%
# LFC < 0 (down)    : 877, 21%
# outliers [1]      : 6, 0.14%
# low counts [2]     : 164, 3.9%
# (mean count < 5)

```

```
# [1] see 'cooksCutoff' argument of ?results
# [2] see 'independentFiltering' argument of ?results
```

Answer: If we add up the up and down regulated genes we get the total amount of significantly differentially expressed genes. Looking at the summary this would be 1863 genes.

Changing the alpha factor:

- d) For the comparison of the genotypes under standard conditions. How many significantly differentially expressed genes in total are reported for a significance level of 0.05? (Go to the RStudio Help and search for “results” function, to identify the attribute you have to change.)

```
# Extract genotype comparison (TGT vs WT, no nickel) with alpha = 0.05
DESeq2Results_genotype <- results(DESeq2Data,
                                   contrast = c("group", "TGT.none", "WT.none"),
                                   alpha = 0.05)

# View summary (this shows the numbers you need)
summary(DESeq2Results_genotype)

# out of 4221 with nonzero total read count
# adjusted p-value < 0.05
# LFC > 0 (up)      : 187, 4.4%
# LFC < 0 (down)    : 275, 6.5%
# outliers [1]      : 6, 0.14%
# low counts [2]     : 0, 0%
# (mean count < 1)
# [1] see 'cooksCutoff' argument of ?results
# [2] see 'independentFiltering' argument of ?results
```

Answer: If we add up the up and down regulated genes we get the total amount of significantly differentially expressed genes. Looking at the summary this would be 462 genes at a p-value < 0.05.

Comparing the nickel treatment to no treatment in the TGT-mutant:

- e) Repeat the steps above for the comparison of the TGT-mutant strain treated with nickel to the TGT-mutant strain not treated with nickel. How many significantly differentially expressed genes in total are reported for a significance level of 0.05?

```
# Extract TGT mutant nickel effect with alpha = 0.05
DESeq2Results_TGT_nickel <- results(DESeq2Data,
                                     contrast = c("group", "TGT.Nickel", "TGT.none"),
                                     alpha = 0.05)

# View summary (this shows the numbers you need)
```

```
summary(DESeq2Results_TGT_nickel)

# out of 4221 with nonzero total read count
# adjusted p-value < 0.05
# LFC > 0 (up)      : 826, 20%
# LFC < 0 (down)    : 752, 18%
# outliers [1]      : 6, 0.14%
# low counts [2]    : 82, 1.9%
# (mean count < 3)
# [1] see 'cooksCutoff' argument of ?results
# [2] see 'independentFiltering' argument of ?results
```

Answer: If we add up the up and down regulated genes we get the total amount of significantly differentially expressed genes. Looking at the summary this would be 1578 genes at a p-value < 0.05.

6 Vizualising data

6.1 Experimental QC - Clustering of samples (PCA)

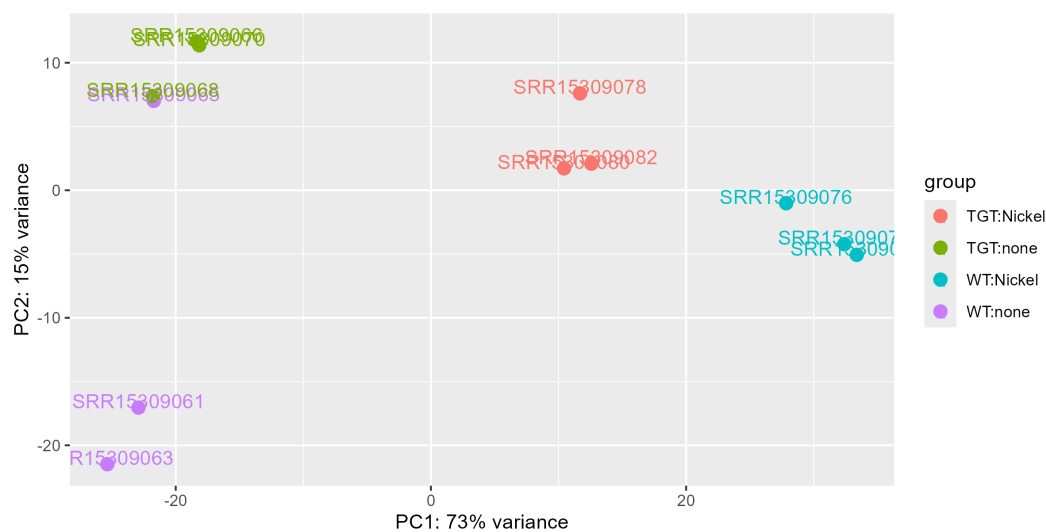


Figure 4: Clustering of samples PCA dot plot.

a) Do the groups of replicates behave as expected?

Answer: Looking at the plot in figure 4 we can observe that the individual groups are clearly separated from each other. It also makes sense that there is a big distance between untreated and Nickel treated strains (WT as well as TGT strains).

b) Which sample would you identify as an outlier?

Answer: Looking at figure 4 we would identify the WT strain untreated as an outlier as it does not completely cluster with its biological replicates but one point also clusters with the TGT untreated cluster.

6.2 Viewing counts for a single geneID

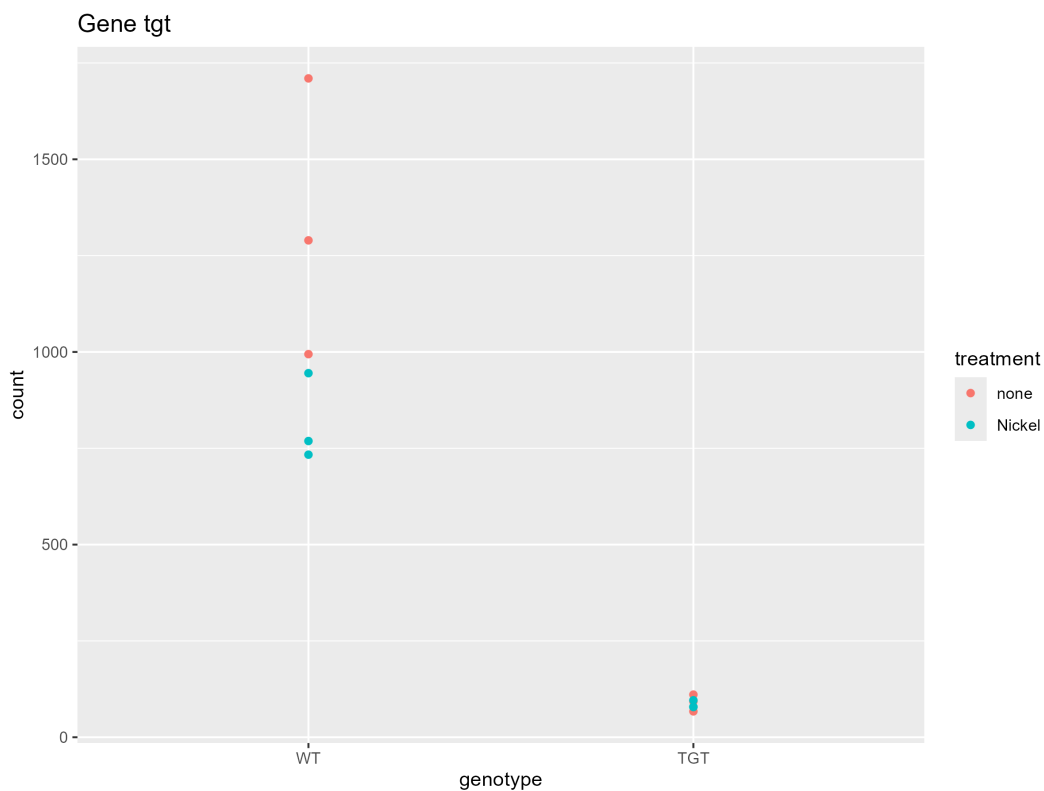


Figure 5: Viewing counts for a single geneID in a dot plot.

c) Think about the mutations in the E. coli strains and how that influences the transcripts of a gene. Are the read counts for the *tgt* gene in the wild type and the knockout strain what you expected? Explain why.

Answer:

Wild type + none: tgt encodes tRNA guanine transglycosylase, an essential enzyme for queuosine synthesis. It is expressed under normal growth conditions to modify tRNAs. Expected Counts are going to be high.

Wild type + Nickel: tgt is transcriptionally repressed by nickel stress. Expected Counts are going to be lower than that of Wt + none but still higher than the expected counts of the knock out mutants.

tgt + none: The tgt gene is deleted so No functional tgt should be transcribed. Expected counts are going to be very low to near zero.

tgt + Nickel: The tgt knock out stays the same but now Nickel is added. This should not have a significant affect on the already very low to near zero counts.

References

- [1] L. Pollo-Oliveira *et al.*, “The absence of the queuosine tRNA modification leads to pleiotropic phenotypes revealing perturbations of metal and oxidative stress homeostasis in escherichia coli K12,” *Metallomics*, vol. 14, no. 9, p. mfac065, Sep. 2022, doi: 10.1093/mtomcs/mfac065.
- [2] L. Pollo-Oliveira *et al.*, “The absence of the queuosine tRNA modification leads to pleiotropic phenotypes revealing perturbations of metal and oxidative stress homeostasis in escherichia coli K12,” *Metallomics*, vol. 14, no. 9, p. mfac065, Sep. 2022, doi: 10.1093/mtomcs/mfac065.
- [3] M. Love, S. Anders, and W. Huber, *DESeq2: Differential gene expression analysis based on the negative binomial distribution*. 2025. doi: 10.18129/B9.bioc.DESeq2.
- [4] H. Wickham *et al.*, *ggplot2: Create elegant data visualisations using the grammar of graphics*. 2025. Available: <https://ggplot2.tidyverse.org>
- [5] Y. Xie, *Knitr: A general-purpose package for dynamic report generation in r*. 2025. Available: <https://yihui.org/knitr/>
- [6] J. Allaire *et al.*, *Rmarkdown: Dynamic documents for r*. 2025. Available: <https://github.com/rstudio/rmarkdown>
- [7] H. Wickham, *Tidyverse: Easily install and load the tidyverse*. 2023. Available: <https://tidyverse.tidyverse.org>
- [8] M. I. Love, W. Huber, and S. Anders, “Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2,” *Genome Biology*, vol. 15, p. 550, 2014, doi: 10.1186/s13059-014-0550-8.
- [9] H. Wickham, *ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York, 2016. Available: <https://ggplot2.tidyverse.org>
- [10] Y. Xie, *Dynamic documents with R and knitr*, 2nd ed. Boca Raton, Florida: Chapman; Hall/CRC, 2015. Available: <https://yihui.org/knitr/>
- [11] Y. Xie, “Knitr: A comprehensive tool for reproducible research in R,” in *Implementing reproducible computational research*, V. Stodden, F. Leisch, and R. D. Peng, Eds., Chapman; Hall/CRC, 2014.
- [12] Y. Xie, J. J. Allaire, and G. Golemund, *R markdown: The definitive guide*. Boca Raton, Florida: Chapman; Hall/CRC, 2018. Available: <https://bookdown.org/yihui/rmarkdown>
- [13] Y. Xie, C. Dervieux, and E. Riederer, *R markdown cookbook*. Boca Raton, Florida: Chapman; Hall/CRC, 2020. Available: <https://bookdown.org/yihui/rmarkdown-cookbook>
- [14] H. Wickham *et al.*, “Welcome to the tidyverse,” *Journal of Open Source Software*, vol. 4, no. 43, p. 1686, 2019, doi: 10.21105/joss.01686.