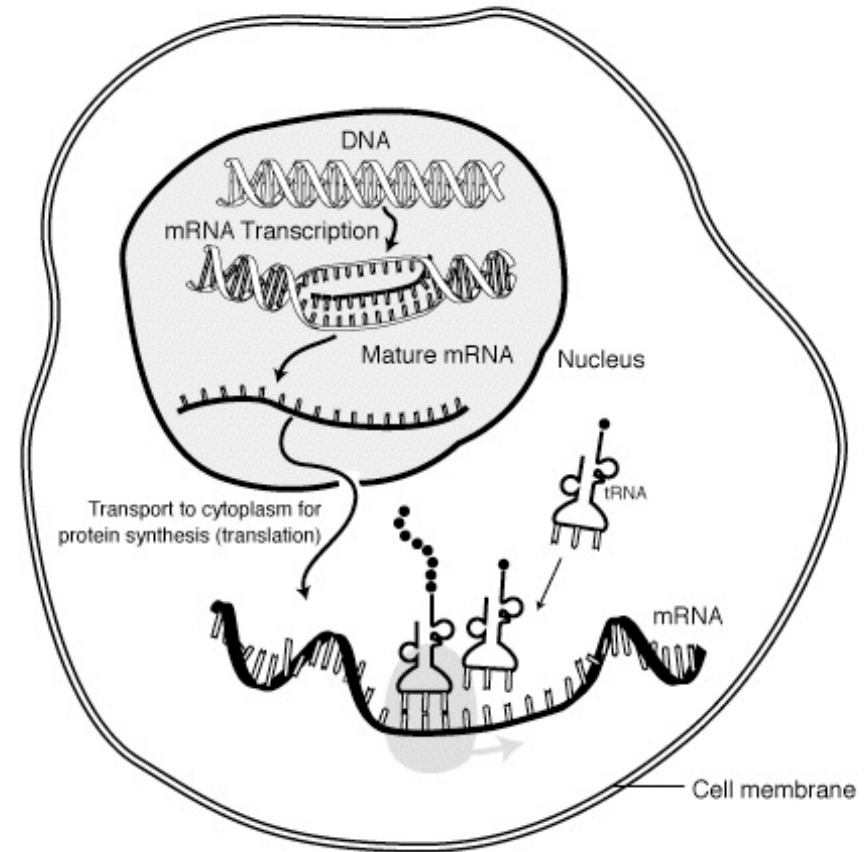


RNA-seq

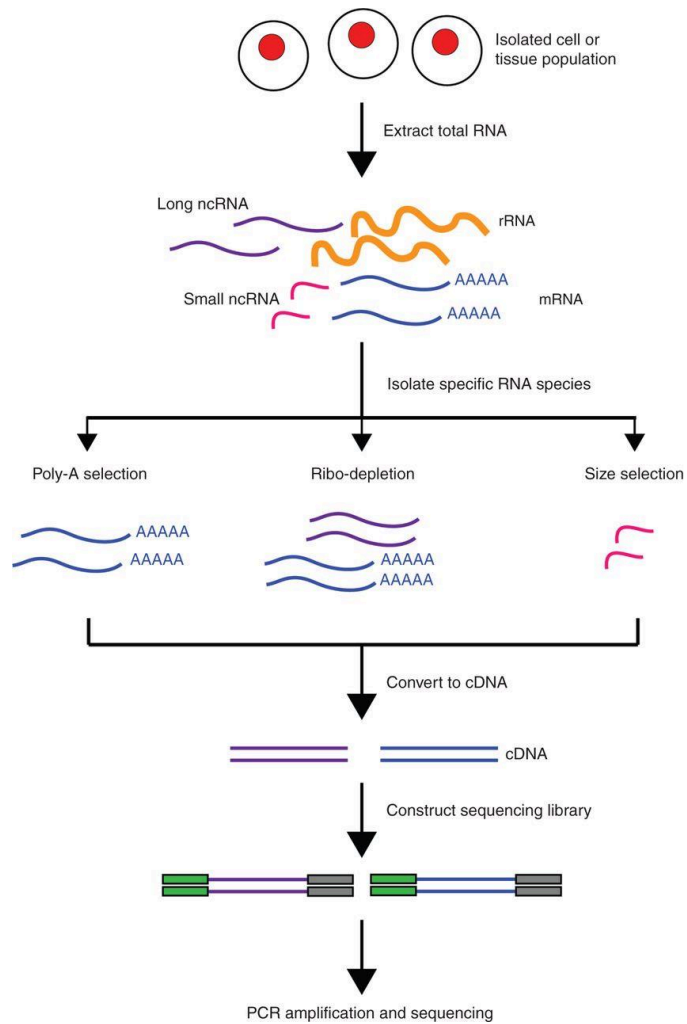
What is RNA-seq and how does it work?

What is RNA-seq?

- RNA-seq = RNA sequencing
- Use next generation sequencing to quantify RNA in a cell
- Mainly when people do RNA-seq they're concerned with mature messenger RNA, because they want to know what genes are expressed

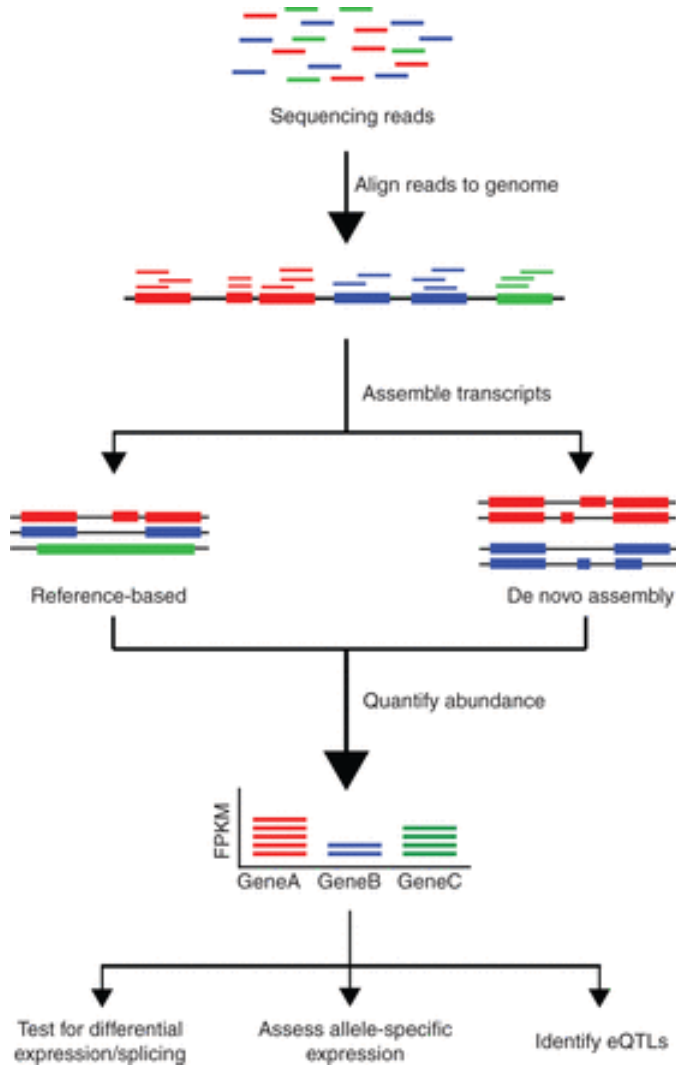


RNA-seq workflow



- RNA is extracted from cell
- Only mRNA (sometimes lincRNA) enriched
- Converted to cDNA
- Construct library
- Sent to sequencing

RNA-seq from sequencing reads to expression



unaligned reads



alignment to genome
(while aware of exon structure)



If we care about alternative transcripts,
we can assemble transcriptome



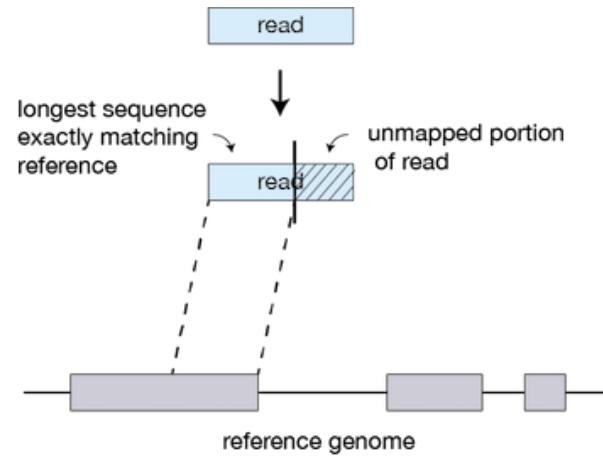
count reads per gene



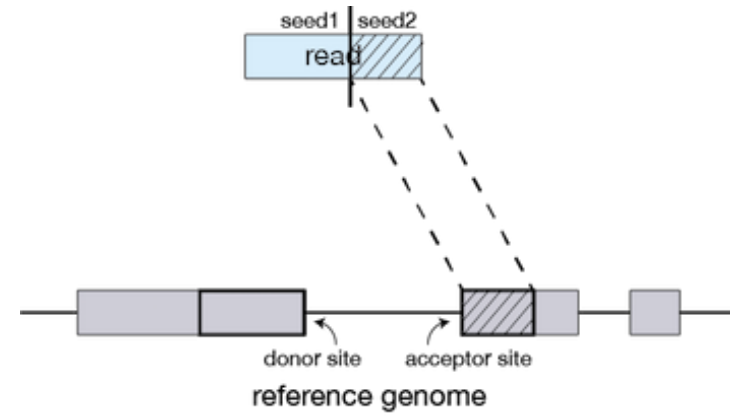
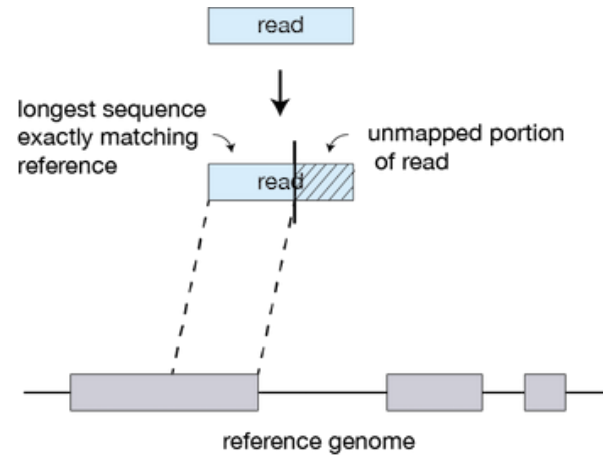
Actual analysis

differential gene expression / allele specific expression / expression QTL

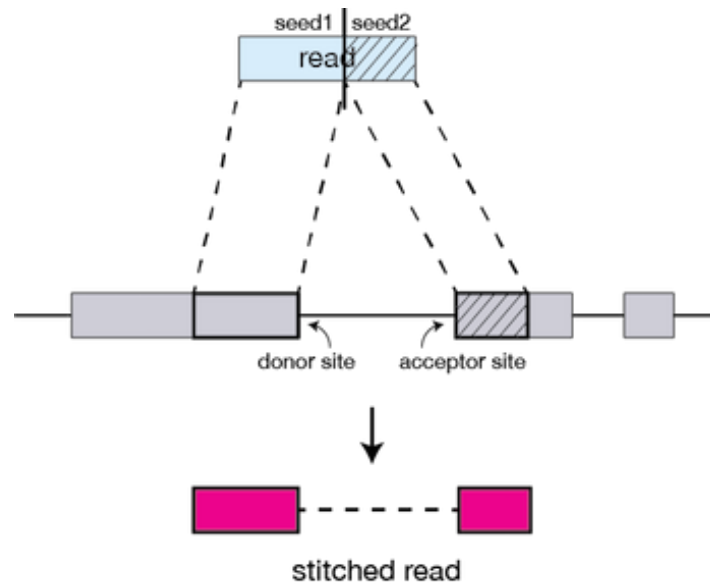
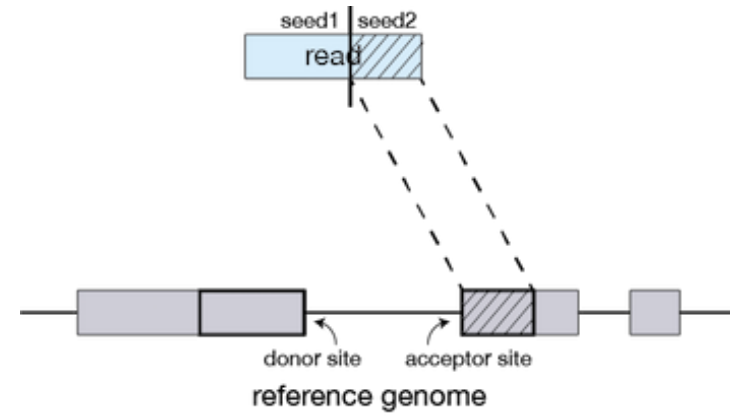
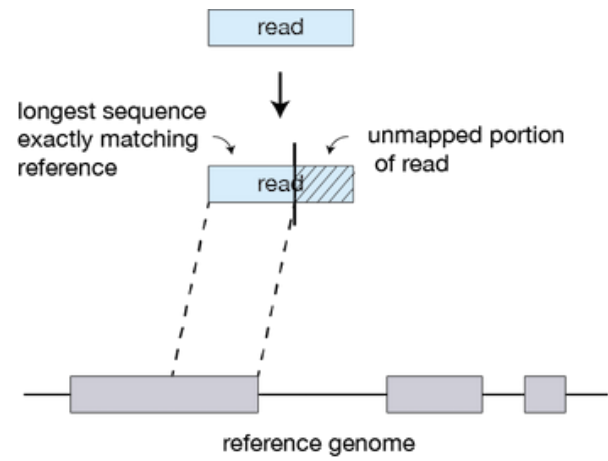
RNA-seq split reads mapping



RNA-seq split reads mapping



RNA-seq split reads mapping



FASTA format

We start with a reference genome to map to

The reference sequence
(chromosome)

Sequence description

```
>20 dna:chromosome chromosome:GRCh37:20:1:63025520:1  
NNNNNNNNNTACTTCGATTGCGTATTTACGGACGTAGCGAGTCTTTAGAGTCTTTTAGTCTGTATC  
GTCGTAGTGTCAGTTCGTAGTCTATGTCGTATTCGTAGGCGTACGTAGTCGTGTAGTCAGTCGTGTT
```

The diagram illustrates the FASTA format with three labels and arrows pointing to their respective parts in a FASTA entry. 'The reference sequence (chromosome)' points to the sequence identifier '>20'. 'Sequence description' points to the text 'dna:chromosome chromosome:GRCh37:20:1:63025520:1'. 'DNA sequence' points to the two lines of nucleotide bases: 'NNNNNNNNNTACTTCGATTGCGTATTTACGGACGTAGCGAGTCTTTAGAGTCTTTTAGTCTGTATC' and 'GTCGTAGTGTCAGTTCGTAGTCTATGTCGTATTCGTAGGCGTACGTAGTCGTGTAGTCAGTCGTGTT'.

DNA sequence

FASTQ format

and sequence reads to map

DNA sequence

```
@EAS139:136:FC706VJ:2:2104:15343:197393 1:Y:18:ATCACG
GTCGTAGTGTCAGTTCGTAGTCTATGTCGTATTCGTAGGCGTACGTAGTCGTGTAGTCAGTCGTGTT
+
!#!#!!*!#^^!#(#*&^*$^^#%&$*$&(^&$^^^$)$*&$**$))$*&^$*##$))$)))))
```

| | |
|----------------|--|
| EAS139 | the unique instrument name |
| 136 | the run id |
| FC706VJ | the flowcell id |
| 2 | flowcell lane |
| 2104 | tile number within the flowcell lane |
| 15343 | x'-coordinate of the cluster within the tile |
| 197393 | y'-coordinate of the cluster within the tile |
| 1 | the member of a pair, 1 or 2 (<i>paired-end or mate-pair reads only</i>) |
| Y | Y if the read is filtered, N otherwise |
| 18 | 0 when none of the control bits are on, otherwise it is an even number |
| ATCACG | index sequence |

Quality scores

$$Q = -10 \log_{10}(P_{\text{err}})$$

(Some) of the Many Gene IDs

- Ensembl ids
 - ENS[species prefix][feature type prefix][a unique eleven digit number].[version number]
 - Human: **ENSG00000168769**, Mouse: **ENSMUST00000101509.9**
- HUGO Gene Naming Consortium sets gene names like BRCA1
- Other ids:
 - Entrez
 - UCSC
 - NCBI
 - RefSeq
 - OMIM
 - **Uniprot**
- If you don't know what type of id it is, google the id

Gene name errors are widespread in the scientific literature

[Mark Ziemann](#), [Yotam Eren](#) & [Assam El-Osta](#) ✉

Genome Biology 17, Article number: 177 (2016) | [Cite this article](#)

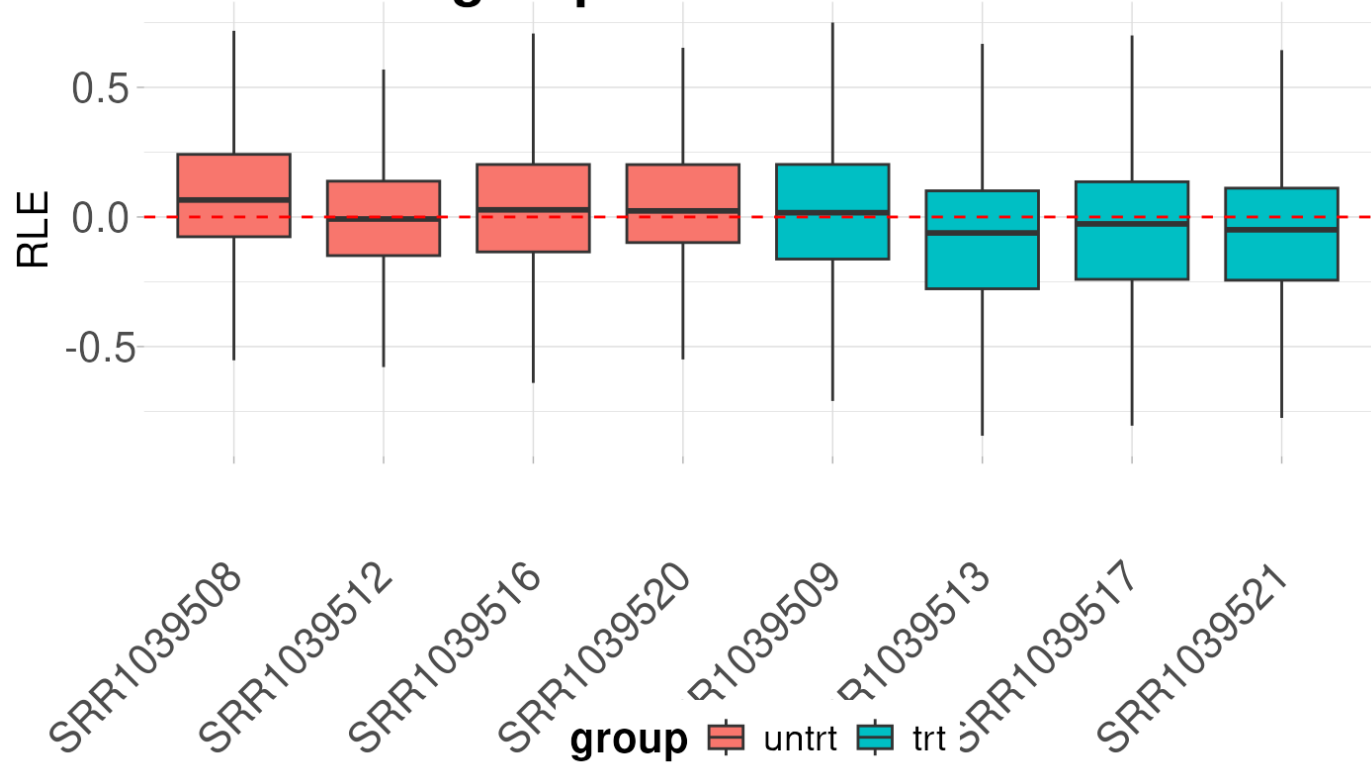
134k Accesses | 53 Citations | 3209 Altmetric | [Metrics](#)

Abstract

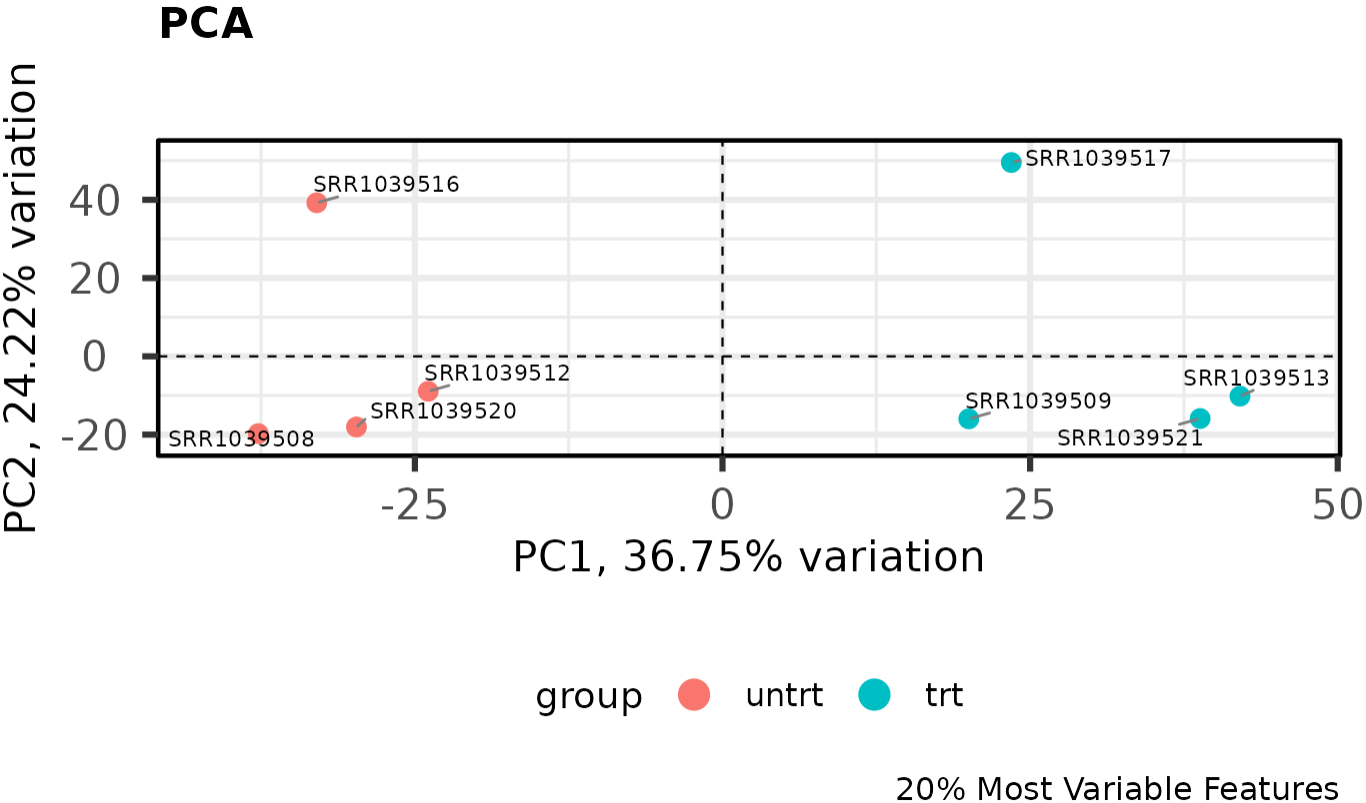
The spreadsheet software Microsoft Excel, when used with default settings, is known to convert gene names to dates and floating-point numbers. A programmatic scan of leading genomics journals reveals that approximately one-fifth of papers with supplementary Excel gene lists contain erroneous gene name conversions.



Relative Log Expression

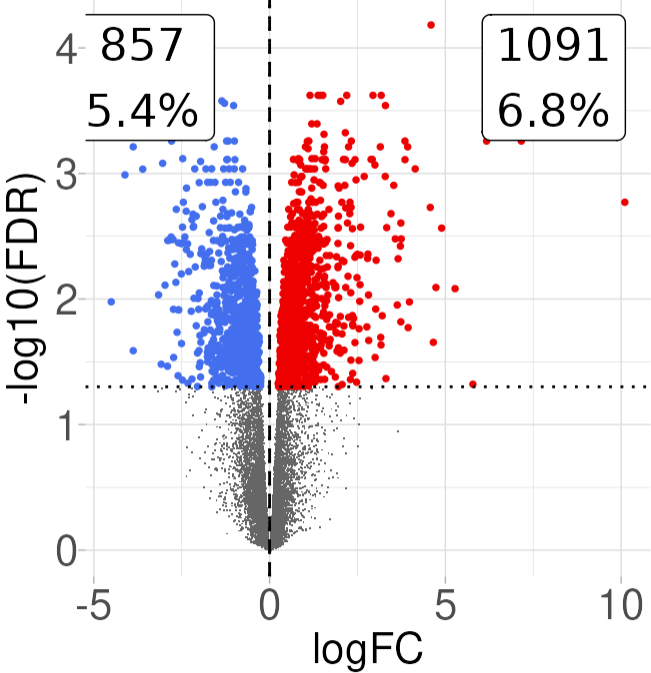


Plotting DE results: PCA

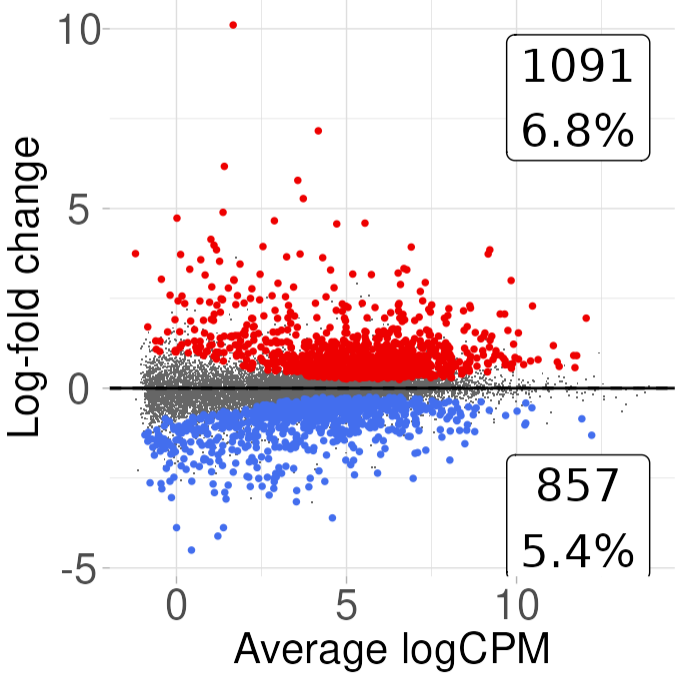


Plotting DE results: Volcano and MA plot

Treatment vs. Control



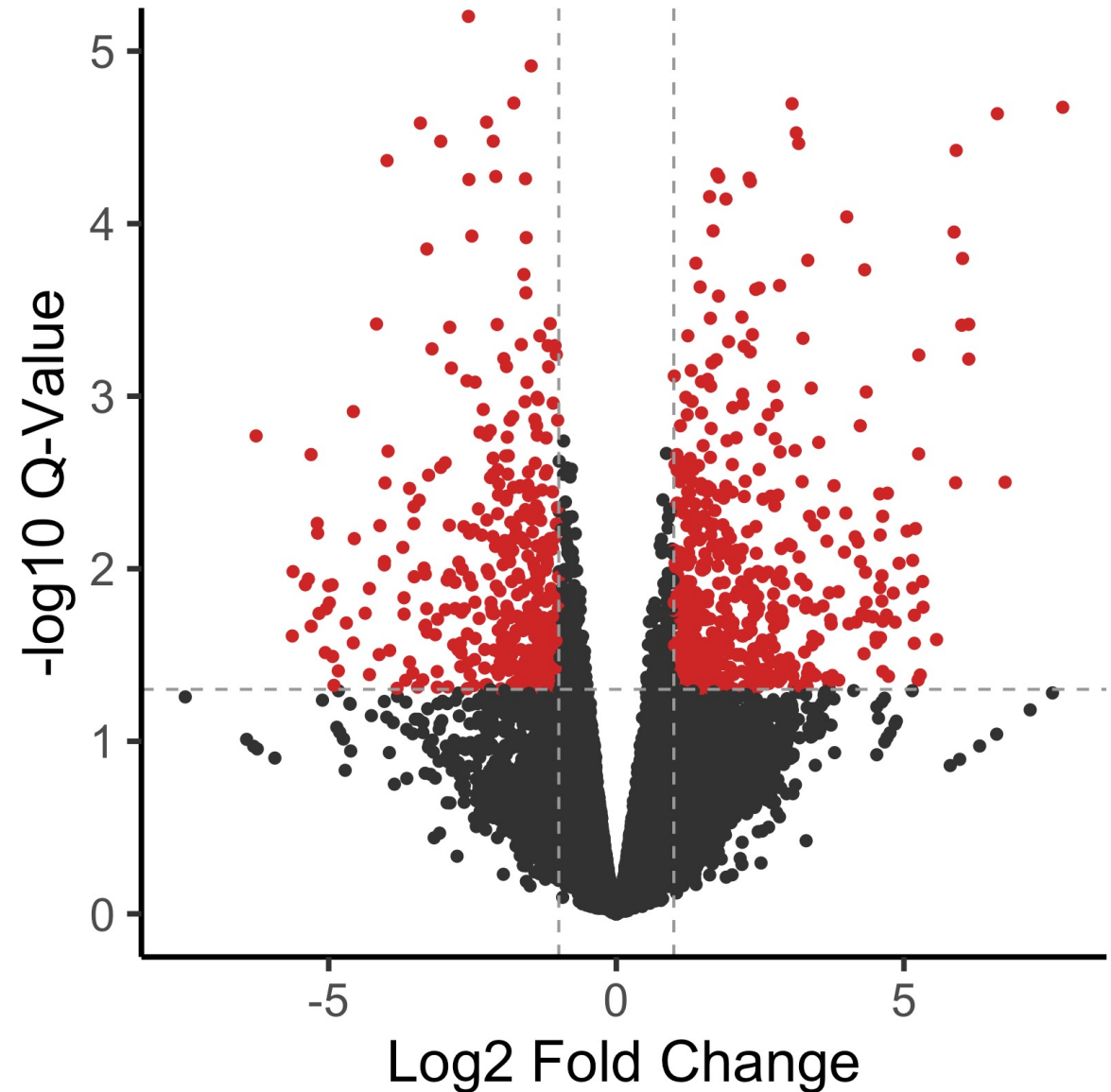
FDR = 0.05
lfc cutoff = 0



FDR = 0.05
lfc cutoff = 0

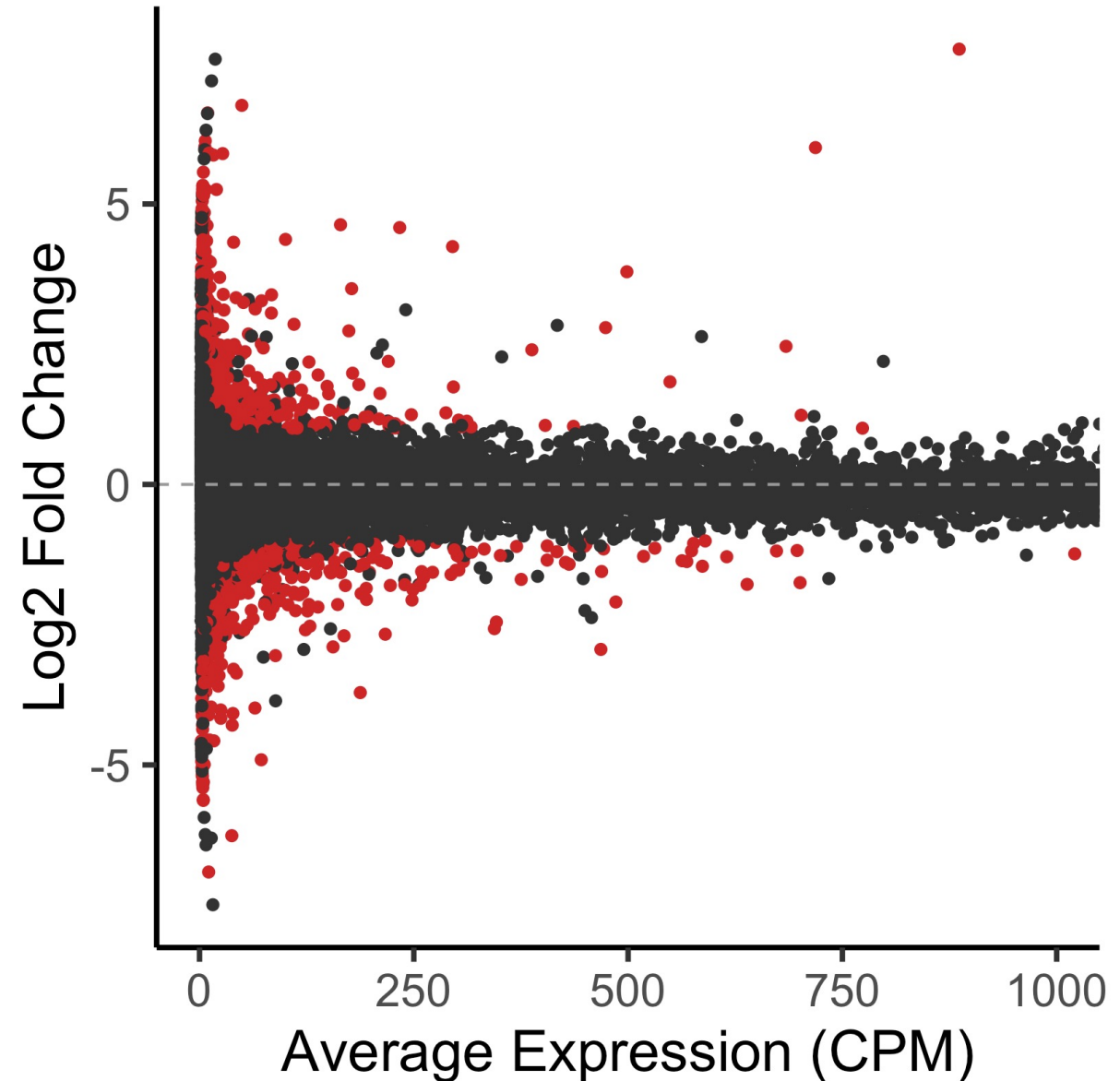
Volcano Plot

- X-axis shows the average log₂ fold change for the variable of interest
- Y-axis negative log₁₀ corrected p-value
- Overall shows how the expression in the treatment group changes and how statistically significant the changes are

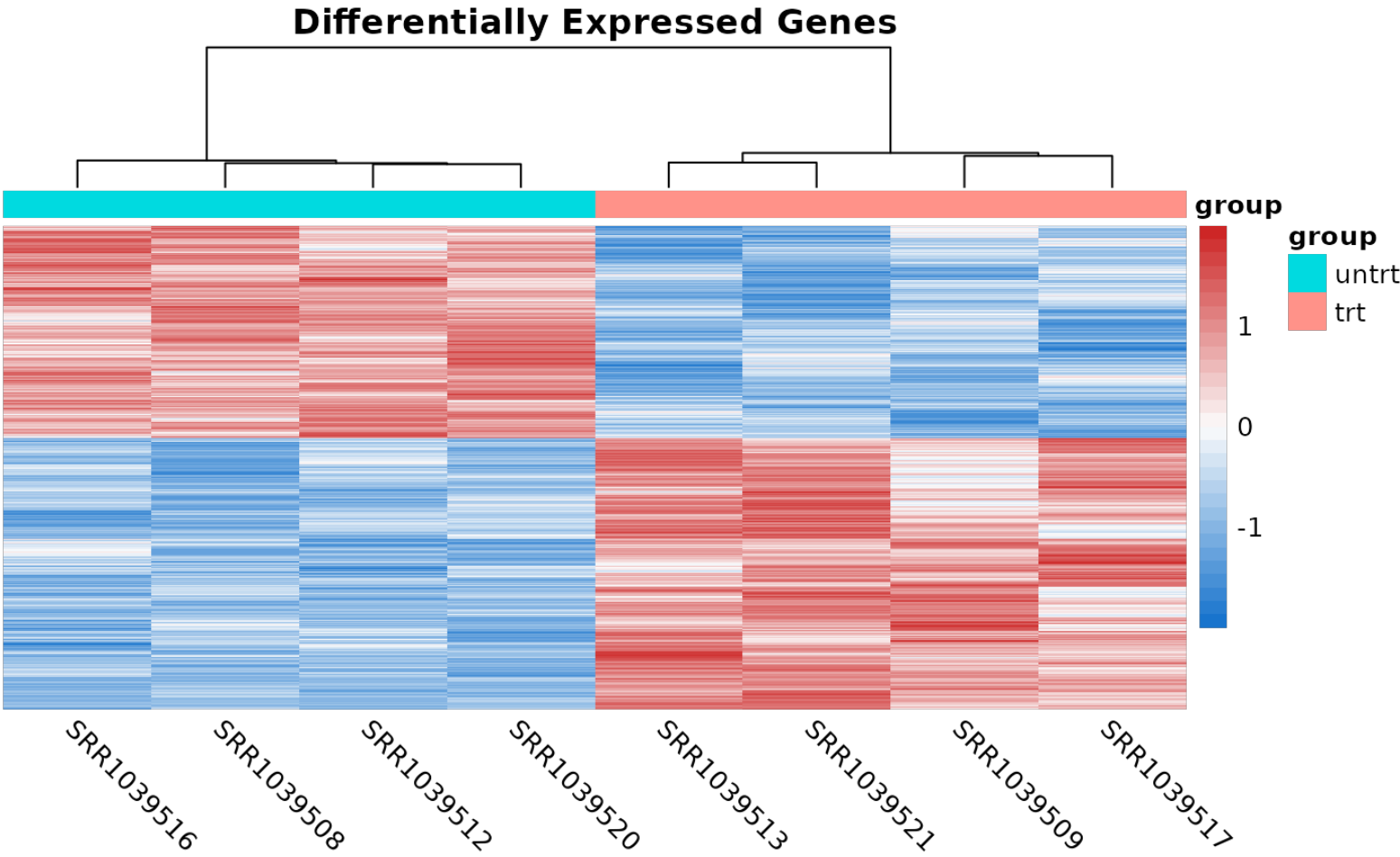


MA Plot

- X-axis shows the average expression for all samples
- Y-axis shows the log₂ fold change for the variable of interest
- Overall shows how the expression in the variable group changes relative to all other samples

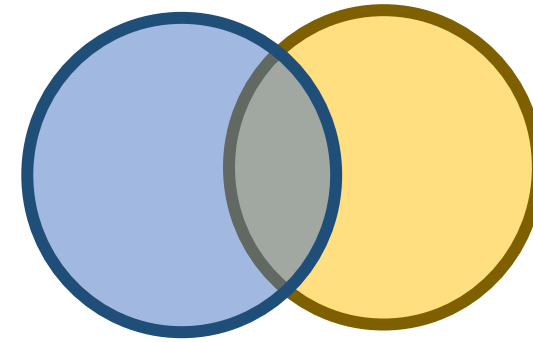


Plotting DE results: Heatmap



UpSet Plot

- New variation on a venn diagram because venn diagrams are not great at visualizing more than 2 or 2 sets
- https://en.wikipedia.org/wiki/Venn_diagram
- Points and line on the bottom show the overlap between the conditions on the left
- Bar height and counts show the number of features that overlap



B

