Course overview: Introduction to single cell

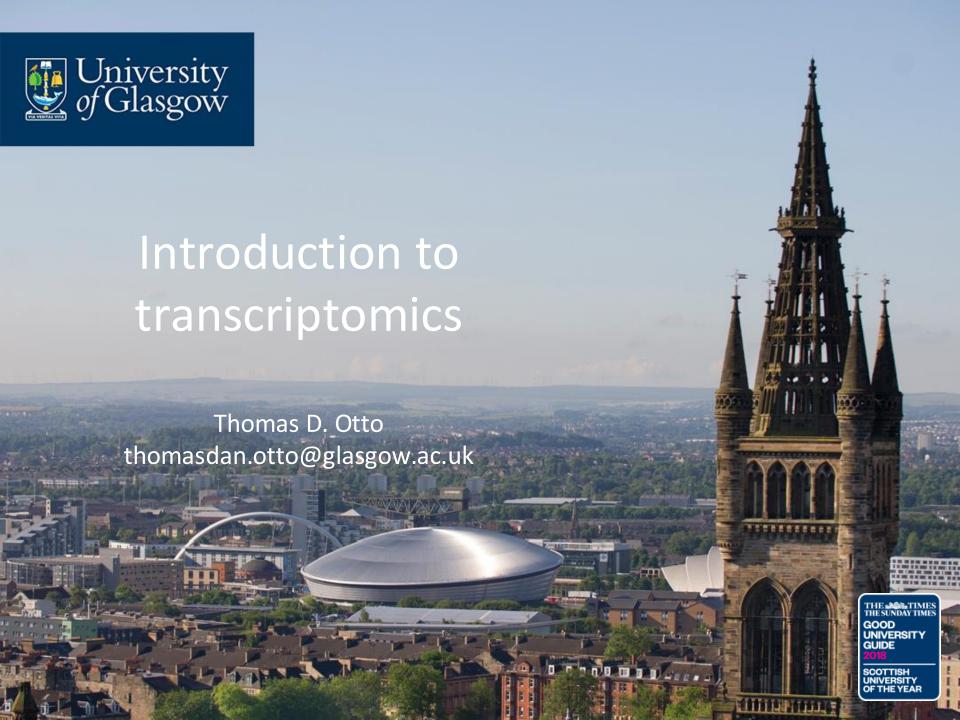
- A. Introduction to the need to transcriptomics, from Qpcr to scRNA-Seq
- B. Introduction to scRNA-Seq
- C. Computational background
- D. Pseudo time
- E. More scRNA methods
- F. Spatial
- G. LLM & co in exploration

Time

- Monday, Wednesday & Thursday 16-19h
- Break when?
- Thursday 14h talk

Informal

 Small group, ask question when something is not clear



Learning outcomes

- Describe the benefit of RNA-Seq when is a gene differentially expressed
- Develop a critical understanding of the need for replicates and a controlled experimental setup
- Perform (and understand) a t-test
- Reflect of the best use to represent any data
- Reflect on potential use of RNA-Seq for your future

Overview

- Why visualizing data?
- How and why to plot data
- Transcriptomics: From Q-PCR to Microarray / RNA-Seq – differential expression analysis
- Compare two conditions what is different?

Overarching aim

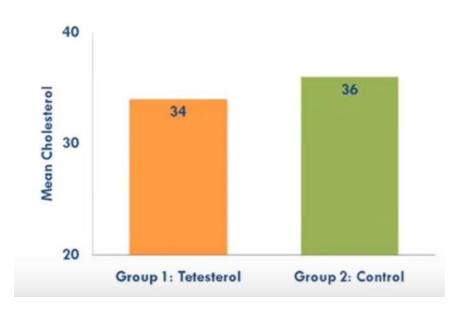
 We want to compare two different samples groups, and find genes that are differently expressed – statistically significant

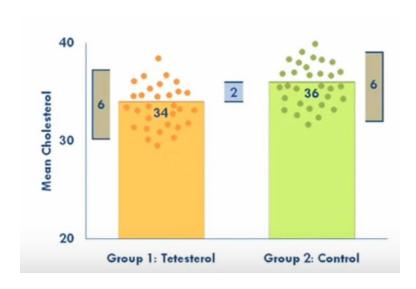
What is significance?

A t-test...

- ... "checks" if something the average of two "groups" are "reliably" (or significantly) different
- That means, that the difference is not obtained by chance

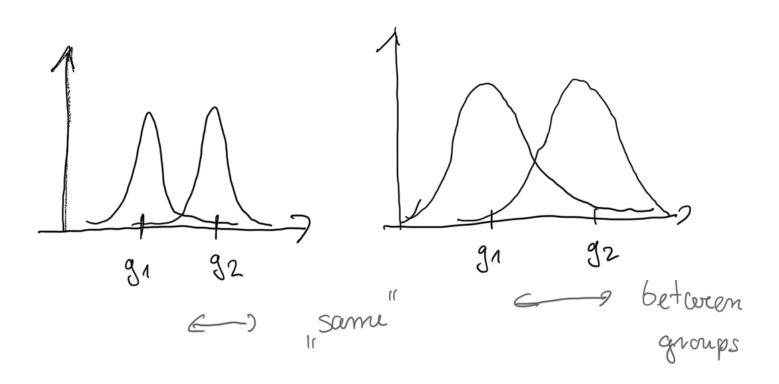
t = variance between groups variance within groups





Difference between groups?

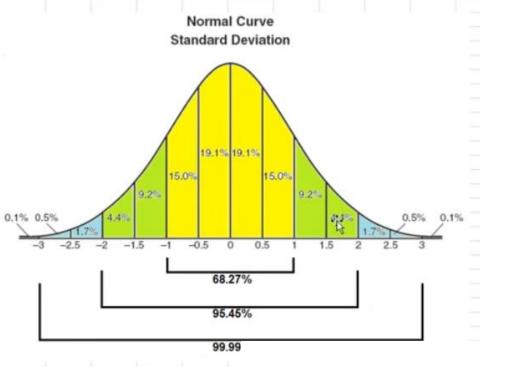
Intuition



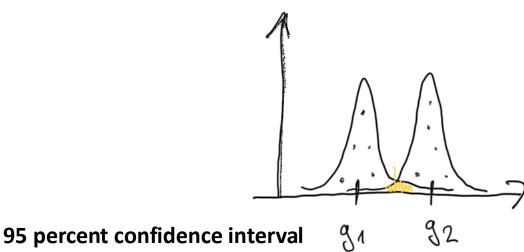
with in group?

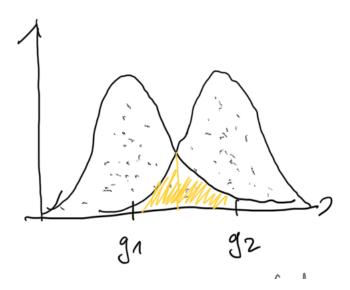
p-value

- each test has a p-value
- p-value tells us the likelihood that there is a real difference
- Specifically, the p-value is the probability that the pattern of data in the sample could be produced by random data
- p=0.001 there is a 0.1% chance that this result was obtained with random data



Gut feeling for P-value





https://www.youtube.com/watch?v=qvRWQre03tQ&t=64s

Sample size – biological replicates

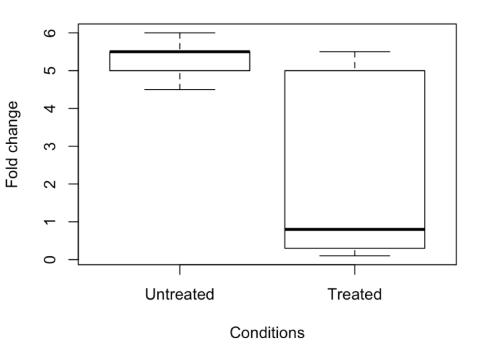
Why is the title in bold?

 The more the better – too small, less likely to find any difference

 "degrees of Freedome" (df) is equal the sample size minus one.

Example of T-test

Boxplot to compare fold expression of gene X



In R: t.test(data[,1],data[,2])

Welch Two Sample t-test

data: data[, 1] and data[, 2]
t = 2.4204, df = 4.3627, p-value = 0.06745
alternative hypothesis: true difference in
means is not equal to 0
95 percent confidence interval:
-0.3269542 6.2469542
sample estimates:
mean of x mean of y
5.30 2.34

Limitations

- 1. Findings limited to type of experiment
- 2. Data need to be normal distributed!
- 3. Compare roughly that same number of datapoints for each group
- 4. Data should be independent to each other
- 5. Units to compare (not ranks)

How to write a t-test results?

"An independent-samples t-test was used to check the effectiveness of a rheumatoid arthristis drug t(99)=0.33, p=0.37, but no significant difference was found (Drug M=34; Control M=36)"

In a more formal way: t-test

- Null Hypothesis: What we are aim to disprove, that taking the drug does not change the expression of our genes of interest.
- Alternative Hypothesis: Our effect of interest, or the antithesis of the null hypothesis: The drug is working
- Significance level (α): Probability of incorrectly rejecting the null hypothesis. This is always established at the beginning of every hypothesis test. **p<0.05** (5% change that the result is due to chance alone)

Transcriptomics

- From Q-PCR
- to
- transcriptomics sequencing

Rheumatoid Arthritis

Rheumatoid arthritis (RA) is a long-term autoimmune disorder that primarily affects joints.

While the cause of rheumatoid arthritis is not clear, it is believed to involve a combination of genetic and environmental factors. The underlying mechanism involves the body's immune system attacking the joints.

The goals of treatment are to reduce pain, decrease inflammation, and improve a person's overall functioning.

Cascades of drugs: Steroids, disease-modifying antirheumatic drugs (DMARDs), biologic DMARDs (using biotechnology)

Imaging:

- A lab where you do your master project has a new drug that might "cure" RA!
- You have been asked to analyse Q-PCR data that where performed on one gene from a clinical trial. "Blood treated with the drug"
- The values are 5.5 and 1.8, for untreated and treated patient samples, respectively.
- What do you do?

What do you do?

- How could you plot the data?
- Does this experiment make any sense?

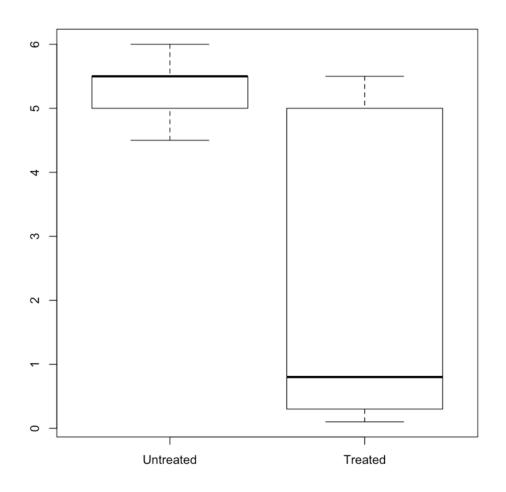
- If yes -> why?
- If no -> why not?

What about more data?

Expression fold change against 3 core genes from 10 patients				
Untreated	Treated			
5.5	0.8			
5	5 5			
4.5	0.3			
ϵ	5.5			
5.5	0.1			

How can we re-present these? Can we do in excel?

The boxplot



Is that a well annotated plot, yes or no?

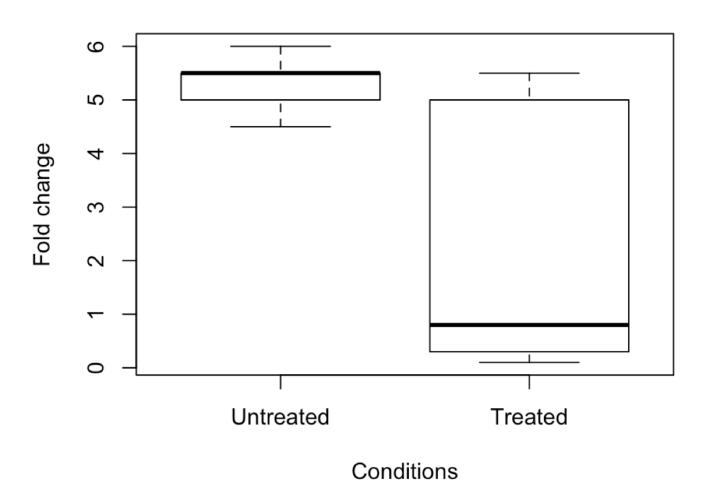
Intro Transcriptomics

• 1112100

- What is the average?
- What is the median?

Is the annotation better?

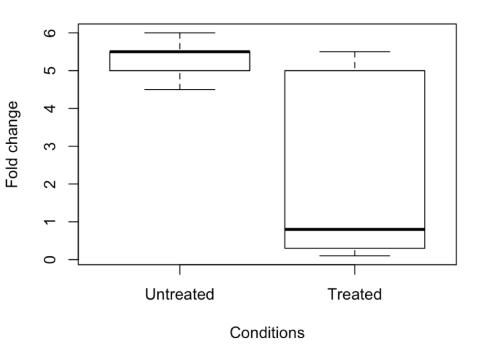
Boxplot to compare fold expression of gene X



boxplot(d[,1],d\$Treated, main="Boxplot to compare fold expression of gene X",xlab="Conditions",ylab="Fold change",names=6("Untreated","Treated"))

The boxplot

Boxplot to compare fold expression of gene X



In R: t.test(data[,1],data[,2])

Welch Two Sample t-test

data: data[, 1] and data[, 2]
t = 2.4204, df = 4.3627, p-value = 0.06745
alternative hypothesis: true difference in
means is not equal to 0
95 percent confidence interval:
-0.3269542 6.2469542
sample estimates:
mean of x mean of y
5.30 2.34

Let's take another gene

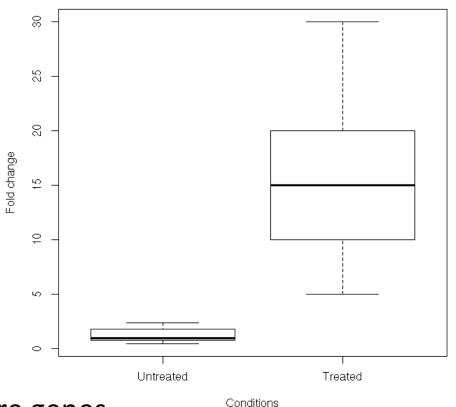
Untreated	Treated
1.8	10
0.8	15
2.4	20
1	5
0.5	30

p-value = 0.02657



- So we don't have to do more genes...
- Is there something wrong here?
 - Correction for Multiple testing





Intro Transcriptomics

Let's split this group

And prove that there is a real difference

Multiple testing, let's make an example

	Null hypothesis is True (H ₀)	Alternative hypothesis is True (H ₁)	Total
Declared significant	V	S	R
Declared non-significant	U	T	m-R
Total	m_0	$m - m_0$	m
• m is the total number by • m_0 is the number of true • $m-m_0$ is the number	•	ttoday	

- m is the total number hypotheses tested
- m_0 is the number of true null hypotheses
- ullet $m-m_0$ is the number of true alternative hypotheses
- ullet V is the number of false positives (Type I error)
- S is the number of true positives (also called "true discoveries")
- T is the number of false negatives (Type II error)
- II is the number of true negatives
- R is the number of rejected null hypotheses (also called "discoveries")
- In m hypothesis tests of which m_0 are true null hypotheses, R is an observable random variable, and S, T, U, and V are unobservable random variables.

Methods for multiple correction

Bonferroni Correction

The most conservative of corrections, the Bonferroni correction is also perhaps the most straightforward in its approach. Simply divide α by the number of tests (m).

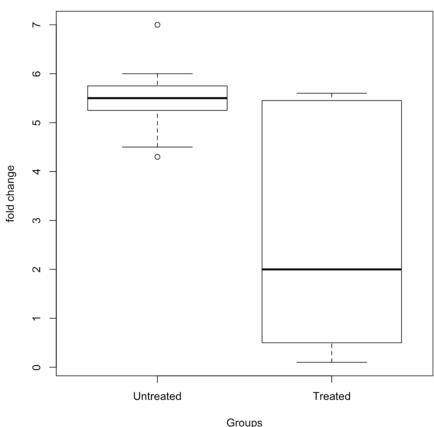
– New α is 0.05/ 2 -> 0.025, so our test before is not significant \odot

What can we do next?

OK, more measuring points

Untreated	Treated	
5.5	0	.8
5		5
4.5	0	.3
6	5	.5
5.5	0	.1
4.3	5	.6
5.8	0	.2
5.9	5	.4
5.5	0	.5
7	5	.5
5.3	0	.5
5.6	5	.6
5.4	0	.6
5.7	5	.2
5.2		2

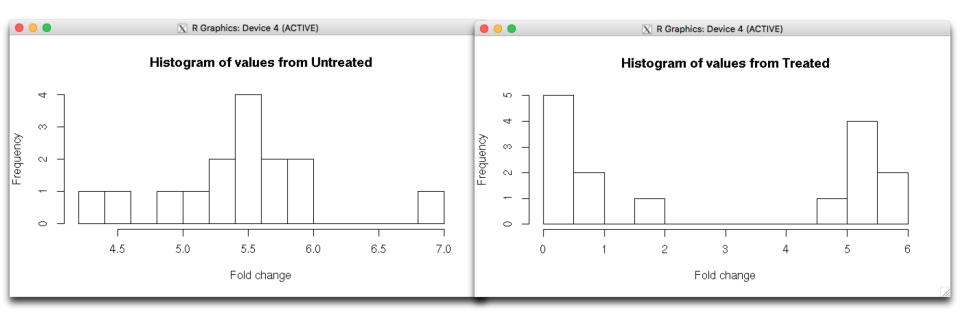
Q-PCR for RA drug assay



p-value = 0.001211

boxplot(data,xlab="Groups", ylab="fold change",main ="Q-PCR for RA drug assay")

Are we happy?



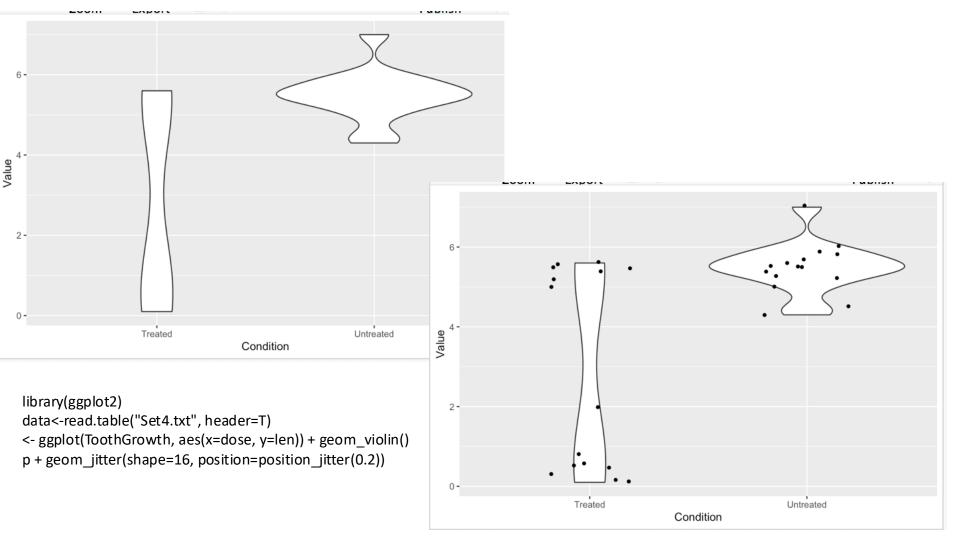
What could that immunological mean? (to see two peaks?)

Some of the individual are non-responders!

hist(data[,2],breaks=10,main="Histogram of values from Treated", xlab="Fold change")

Intro Transcriptomics

Better: Violin plots - ggplot2



Are you happy with the annotation of the graphs? No title, no y-axis label!

Intro Transcriptomics

Conclusions (2/n)

- Think about your null hypothesis and set Significance level (α)
- Look for significance (eg t-test) have enough replicates
- Think about multiple test correction
- Visualise the data in different manners (histogram, boxplots, violin plots and many more)
- In a scientific project the **data** are the results, the heart of your experiment. The correct visualization is crucial.

Consideration

- From which part of the body did we take the sample? (Can we access all part of the body?)
- Is it a mixed or single cell type?
- Could that experiment be real?
- Include clinical data

 Are there other methods to access expression of genes?

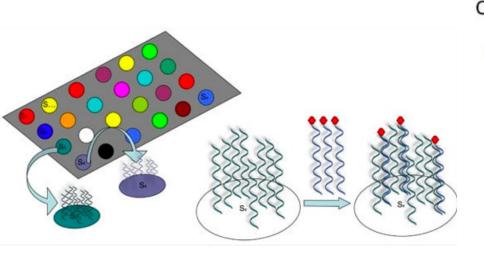
Part 2.2: Transcriptomics

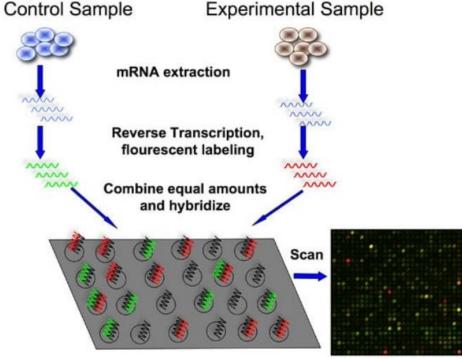
- Micro array not doing this anymore
- RNA-Seq
 - General idea
 - Visualisation
 - Differential expression analysis in R

Why do we want to access the complete transcriptome?

Q-PCR versus microarray

 A microarray is a laboratory tool used to detect the expression of thousands of genes at the same time.





Microarray analysis is challenging

- Access expression of all genes (the one you have probes - can buy)
- Capture of signal (light) is not perfect, especial low abundance.

 With the RNA-Seq (sequencing of RNA / transcripts) most groups moved away from arrays.

Example: Identify markers and mechanisms of resistance to adalimumab therapy

Research article

Open Access

Gene expression profiling in the synovium identifies a predictive signature of absence of response to adalimumab therapy in rheumatoid arthritis

Valérie Badot^{1,2}, Christine Galant³, Adrien Nzeusseu Toukap¹, Ivan Theate³, Anne-Lise Maudoux¹, Benoît J Van den Eynde⁴, Patrick Durez¹, Frédéric A Houssiau¹ and Bernard R Lauwerys¹

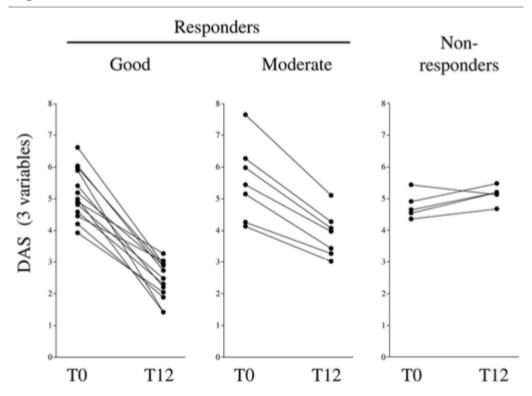
 Paired synovial biopsies were obtained from the affected knee of 25 DMARD (disease-modifying antirheumatic drug)-resistant RA patients at baseline (T0) and 12 weeks (T12) after initiation of adalimumab therapy.

Intro Transcriptomics

(PMID:19389237)

Group patients by responder

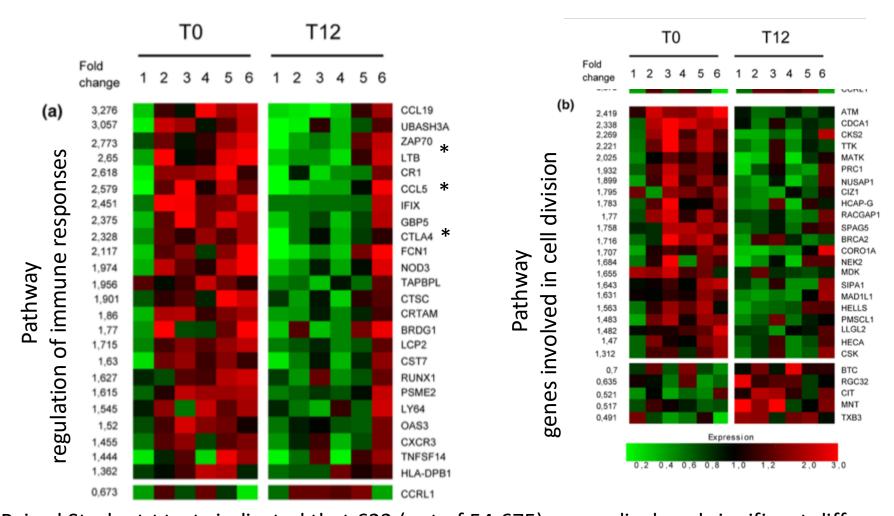
Figure 1



- New type of graph
- Baseline versus later time point

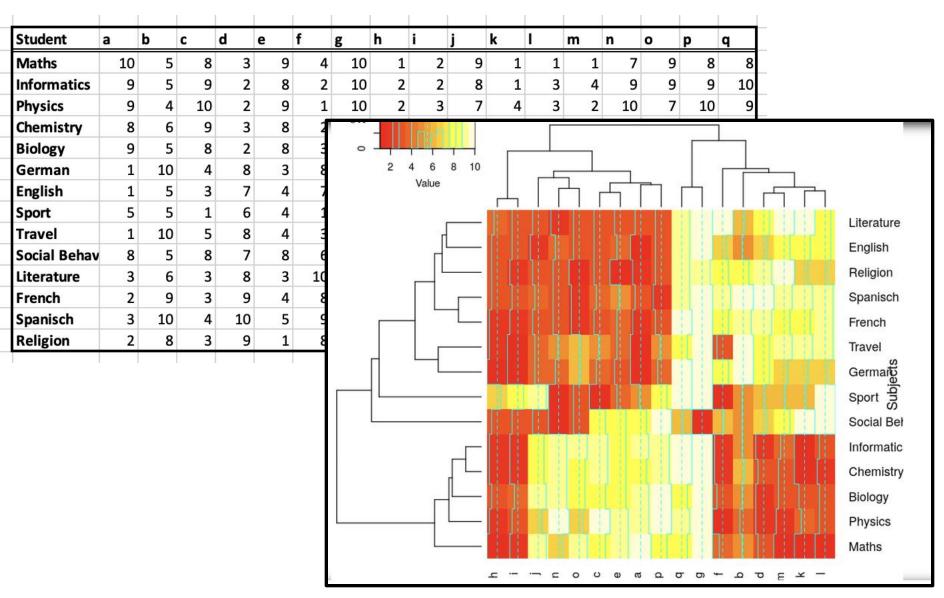
Evolution of disease activity score (DAS) (three variables) in 25 individual rheumatoid arthritis patients before (T0) and 12 weeks after (T12) initiation of adalimumab therapy. Patients are categorized into (good or moderate) responders or non-responders according to European League Against Rheumatism criteria.

Differentially expressed genes before (T0) and 12 weeks after (T12) start of adalimumab in synovial biopsy specimens of rheumatoid arthritis patients who responded to therapy.



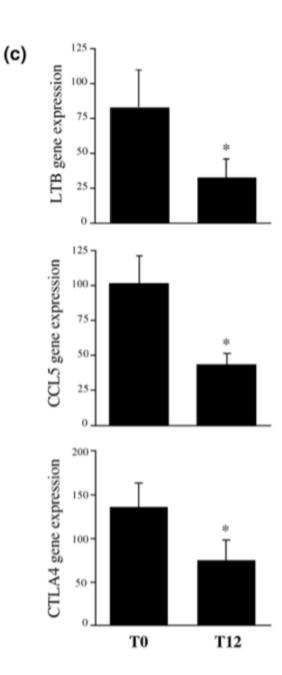
Paired Student *t* tests indicated that 632 (out of 54,675) genes displayed significant differences in expression between T0 and T12 in six synovial tissue samples obtained from RA patients who responded to adalimumab therapy

Do you know what a heatmap is?



Or explain better?

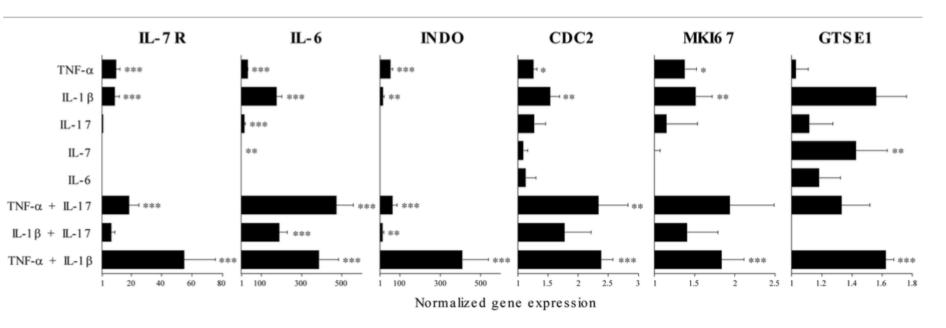
 /Users/thomasdan.otto/Library/CloudStorage /OneDrive-UniversityofGlasgow/Teaching/OtherWorksho ps/2025_Uruguay_SC/Data/



Q-PCR confirmation

- QPCR (T0) (n = 10) and 12 weeks after (T12) (n = 8) initiation of adalimumab therapy
- Samples were loaded in triplicate
- **P* < 0.05.
- CCL5, chemokine lig- and 5; CTLA4, cytotoxic T-lymphocyteassociated antigen 4; LTB, lymphotoxin beta.

fibroblast-like synovial cells cultured with of different cytokines



- evaluated in at least four different experiments
- mean fold change in gene expression and standard error of the mean, relative to the mean gene expression of the baseline condition normalized to
- *P < 0.05, **P < 0.005, ***P < 0.005 using Wilcoxon signed rank test

Summary of paper

- identified baseline markers of response to TNF blockade in a group of RA patients treated with adalimumab
- genes overexpressed in the poor responders are induced by TNF- α , but also by IL-1 β
- "initiate larger studies in order to confirm the prognostic value of our markers"

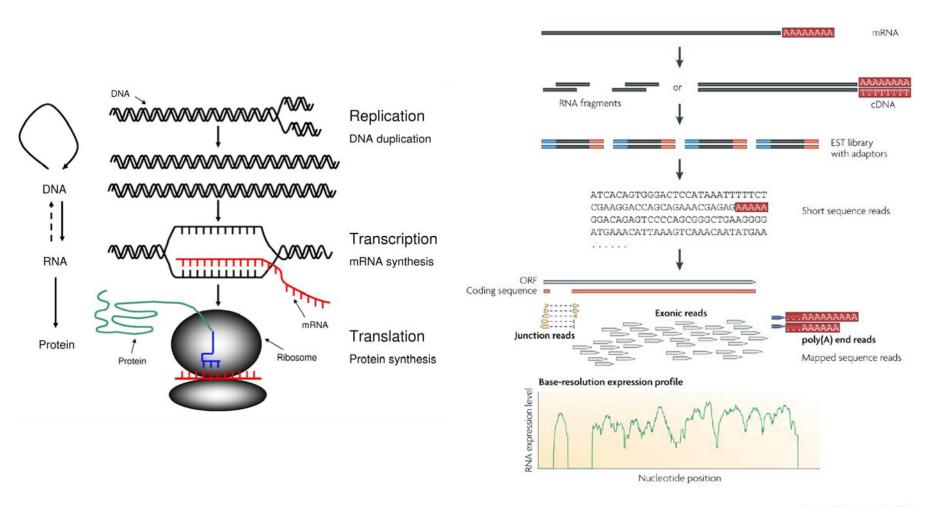
Summary

- Microarray analysis is challenging, need of many replicates
- Access expression of all genes (the one you have probes - can buy)
- Capture of signal (light) is not perfect, especial low abundance.
- With the RNA-Seq (sequencing of RNA / transcripts) most groups moved away from arrays.

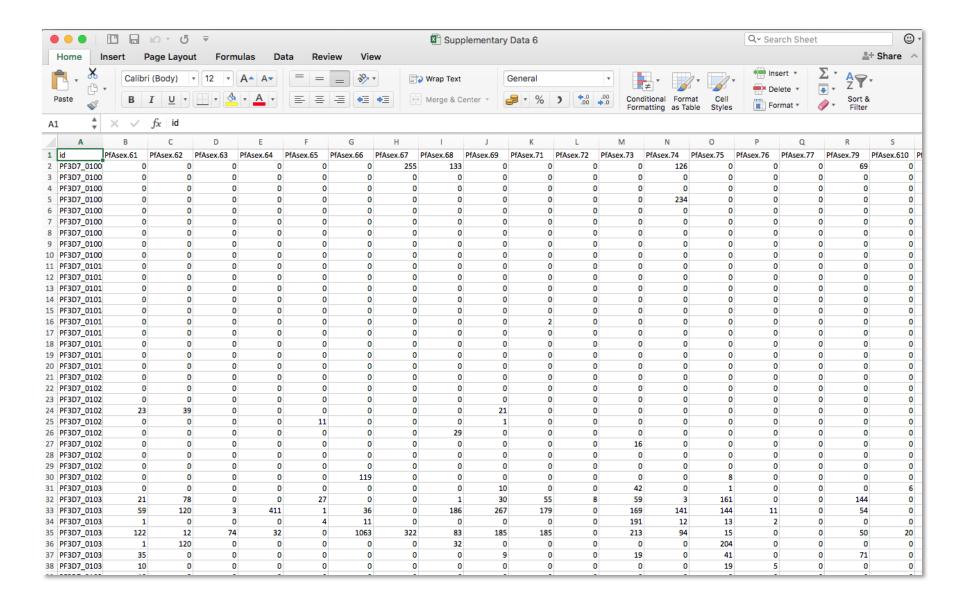
Probes / Antibodies

Concept will return

Transcriptome sequencing



How does the data looks like?



Microarray vs. RNA-seq





	Microarray	RNAseq	
Basis	Microarray platform required	Reference genome preferable	
Interspecific comparison	Tricky	Relatively straightforward	
Sample quantity	100 ng RNA per sample	1ug RNA per sample	
Coverage	Low abundance transcripts not detectable	Single copy transcripts detectable with enough sequence coverage	
Informatics	Relatively low requirements	Currently intense requirements	
Cost	Fairly cheap once platform setup	Relatively expensive but falling.	

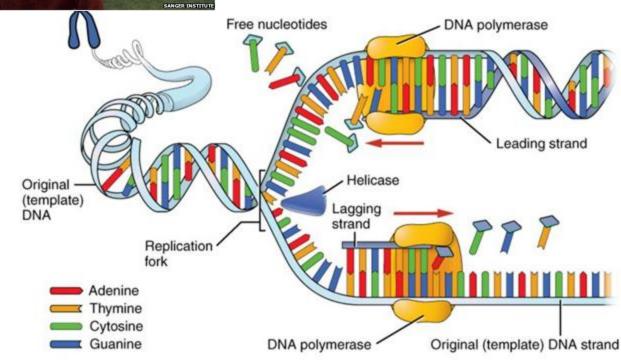
In fact RNA-seq is much more like EST sequencing, but truly global and quantitative

How Does Sequencing Work?

SANGER INSTITUTE

Sanger sequencing (1975)

- Carry out replication under controlled conditions
- Artificially slow down the reaction to see the order in which bases are incorporated



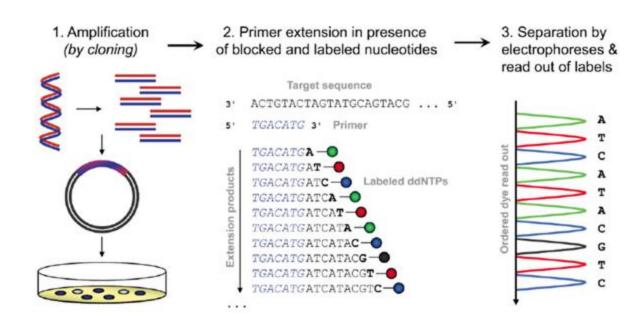
How Does Sequencing Work?

Sanger Sequencing (1975)

Uses modified nucleotides (didioxynucleotide - ddNTPs) that cannot be extended

Each ddNTP is labelled with a different dye so you can see the order in which they are incorporated

Long reads and low error rate, but low throughput



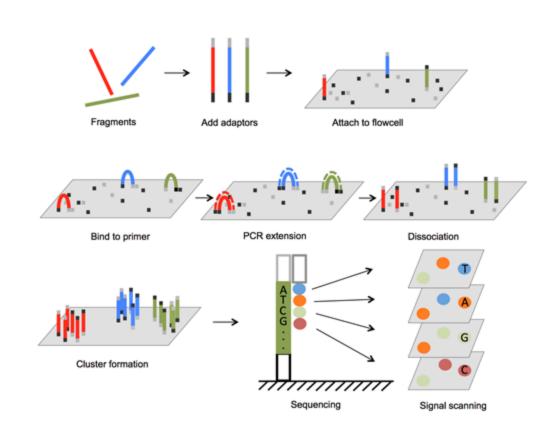
How Does Sequencing Work? Illumina Sequencing

Illumina Sequencing

DNA fragments are adaptor-ligated and attached to a flow cell

PCR is carried out in situ to form clusters

Sequencing can be carried out on millions of clusters simultaneously



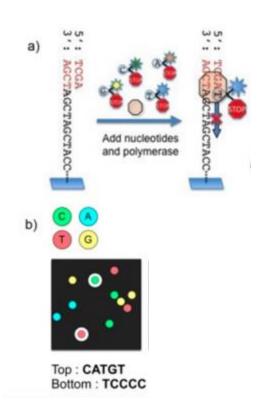
How Does Sequencing Work?

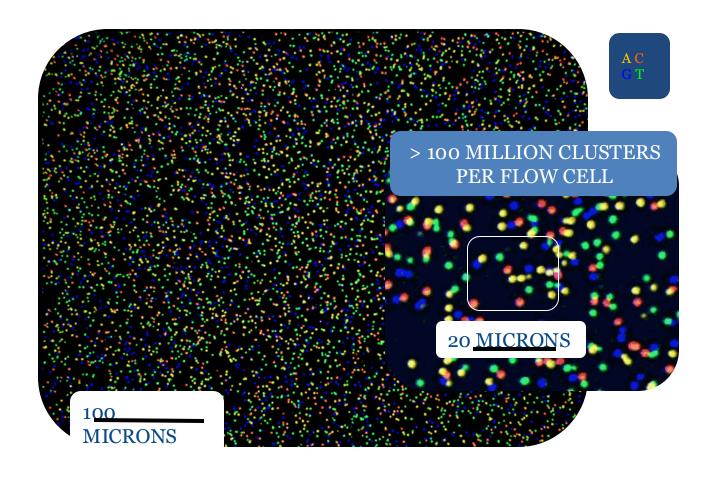
Illumina Sequencing

Blocked and labelled nucleotides are added

1 nucleotide is incorporated and an image is taken of the array

Label and block are removed and cycle repeats





Illumina reversible terminator

Figure S1. a. Structure of the reversible terminator 3'-O-azidomethyl 2'-deoxythymine triphosphate (T) labelled with a removable fluorophore. b. Structure of the incorporated nucleotide after removal of the fluorophore and terminator group. Each of the four nucleotides have an equivalent structure to the one shown here, except for the different base and a corresponding base-specific fluor.

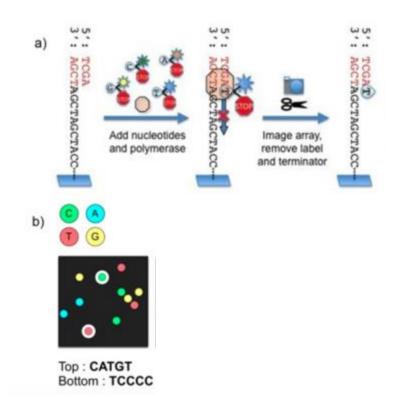
How Does Sequencing Work?

Illumina Sequencing

Blocked and labelled nucleotides are added

1 nucleotide is incorporated and an image is taken of the array

Label and block are removed and cycle repeats



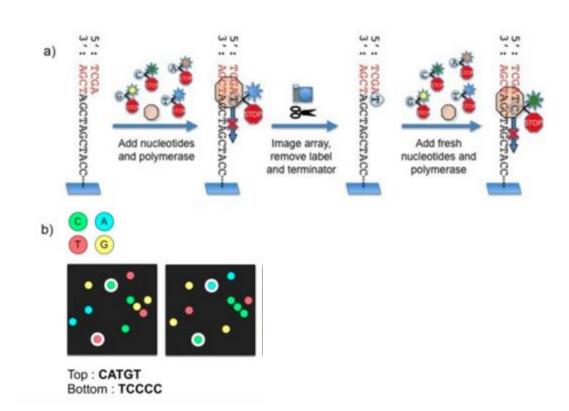
How Does Sequencing Work?

Illumina Sequencing

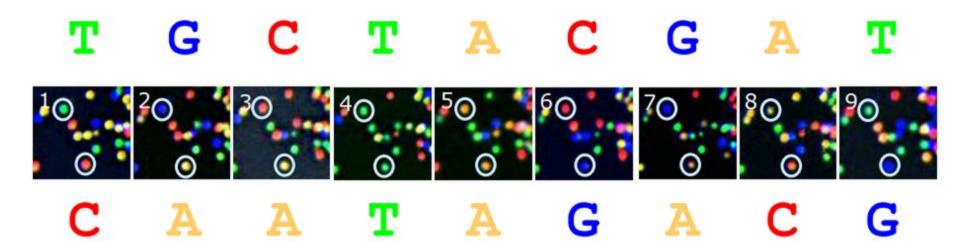
Blocked and labelled nucleotides are added

1 nucleotide is incorporated and an image is taken of the array

Label and block are removed and cycle repeats



Basecalling



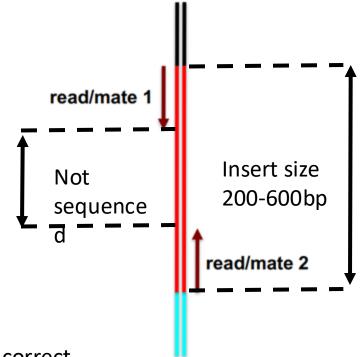
Paired End Sequencing

A short read is sequenced from each end of each fragment

Blue and black are adaptors

Insert or fragment size is the distance between adaptors

The region between the reads is not sequenced - it is covered by other fragments



Reads that map in the correct orientation and the expected distance apart are "concordant" or "proper pairs"

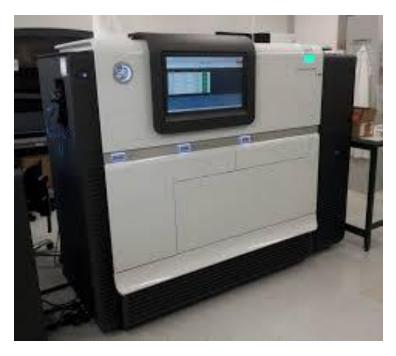
Concordant alignments are prioritised

Illumina HiSeq

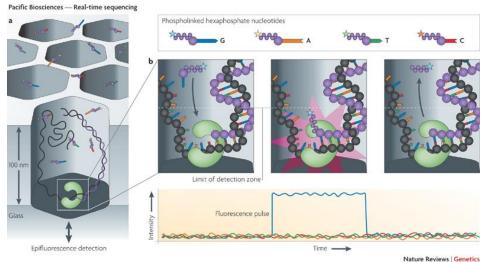




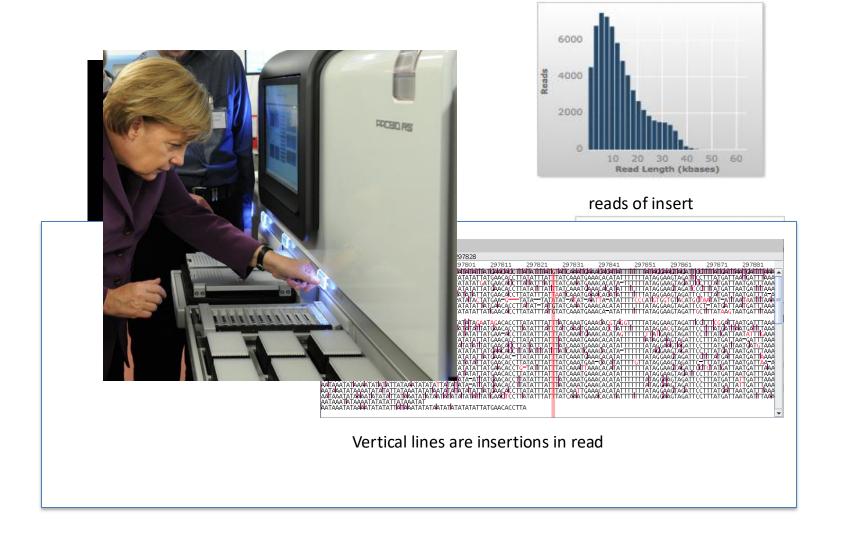
Pacific Biosciences and Single Molecule Real Time (SMRT) sequencing



- Sequencing via light detection
- Light passing through channels is attenuated
 - "Zero Mode Waveguides"
 - Tiny reaction chambers
 - A single polymerase is monitored



PacBio RS



raw reads

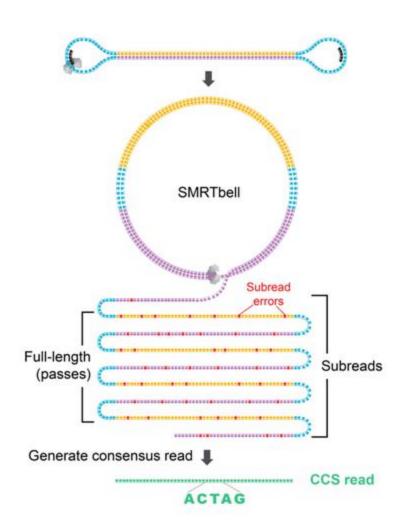
Hifi – High quality

Third Generation PacBio

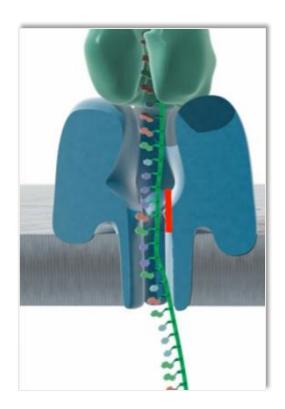
A polymerase trapped in a well on a plate synthesises DNA

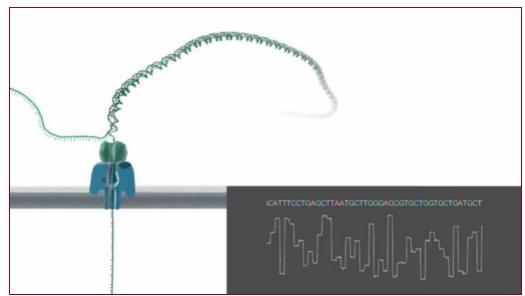
High stochastic error rate can be mitigated by circular consensus sequencing

This reduces the read length



Strand Sequencing by Oxford Nanopore







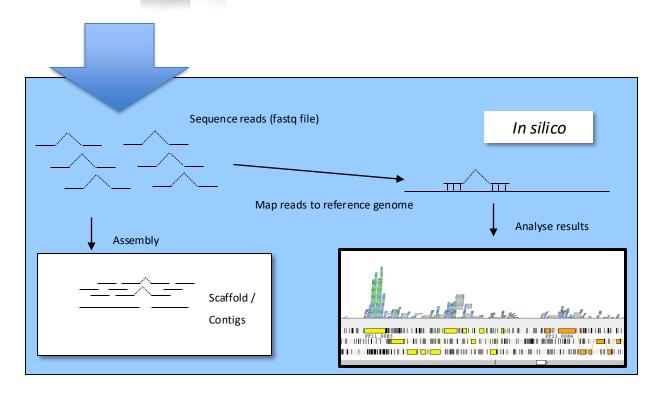


Which Technology to Use?

Sanger	Illumina	PacBio	Nanopore
Cheap	Mid	Expensive	Cheap
1000 bp	30 - 350 bp BEST 150bp	10,000 bp	100,000 bp (?)
>99.5% accuracy	>99% accuracy	85% accuracy Now 99%	70% accuracy Q20: 99%??
Low throughput			Portable
Sequence a mutant or construct	Genome wide SNV analysis, differential expression	Generating a new reference assembly	Reference assembly, field monitoring







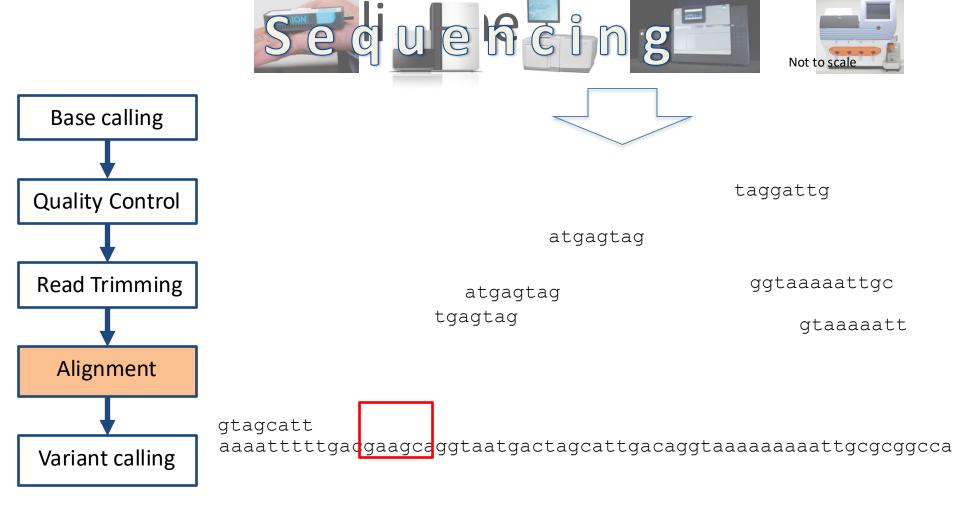
What is an alignment? • THOMAS

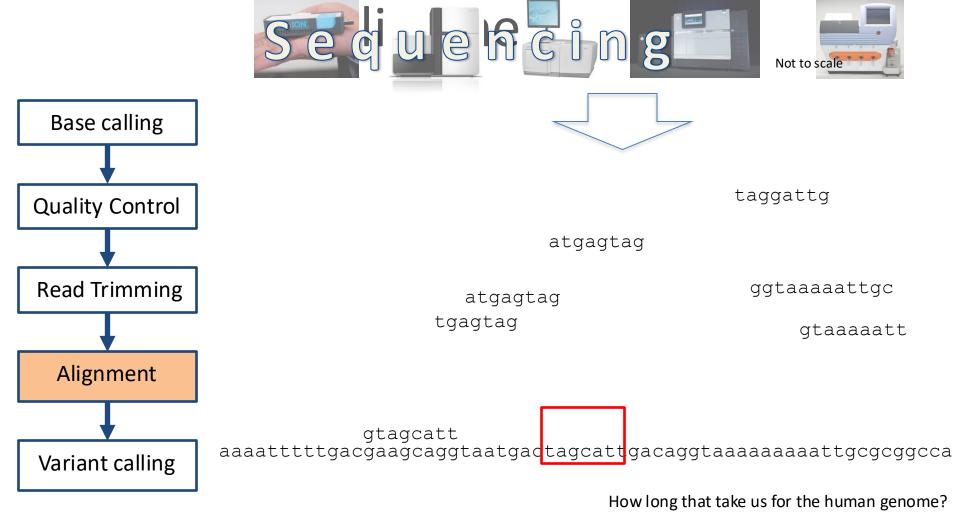
- T-OMAZ

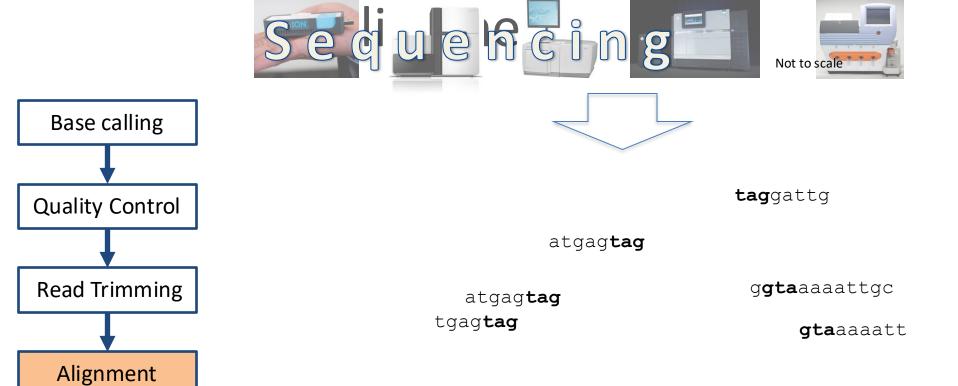
Align the following two sequences:

ATTGAAAGCTA GAAATGAAAAGG

THOMAS TOMAZ 1 "aligne" THOMAS T-011 + 2





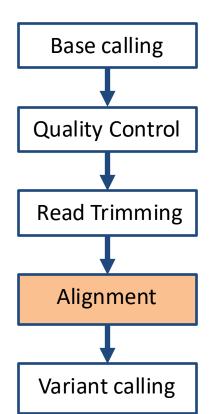


Variant calling

g**tag**catt aaaatttttgacgaagcaggtaatgac**tag**cattgacag**gta**aaaaaaaattgcgcggcca







```
atgagtag
tgagtag
tgagtag
tagcattg ggtaaaaa---ttgc
gtagcatt gtaaaaaaaattggggggg
```

aaaatttttgacgaagcaggtaatgac tag cattgacaggtaaaaaaaaattgcgcggcca

atgag**tag**

Reference





Base calling **Quality Control** Read Trimming Alignment Variant calling

```
atgagtag
tgagtag
tagcaatg ggtaaaaa----ttgc
gtagcatt gtaaaaa----tt
```

atgagtag

Reference – Could the SNP be in the COVID spike protein?

How to map reads

reference: GGGTTAGCGATGGAGA

read1: GA

read2: TAGCG

read3: GGCG

How to map reads

Read

Reference

TAGCG

| | | | | |

GGGTTAGCGATGGAGA

Read

Reference

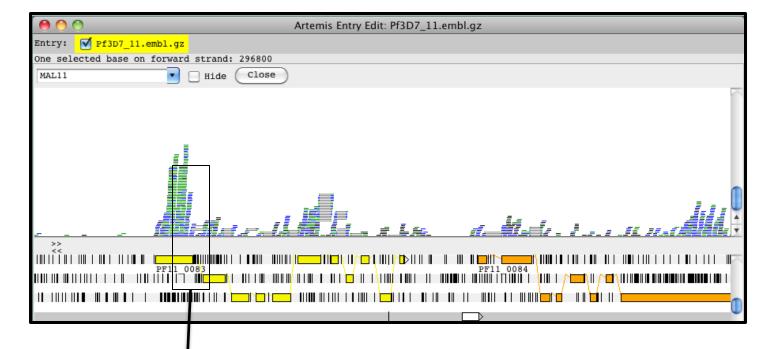
GG----CG

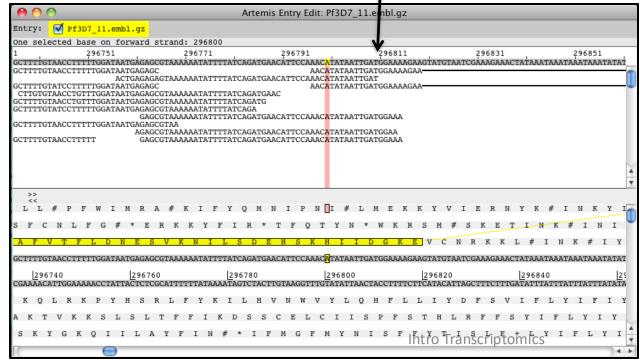
1 11

GGGTTAGCGATGGAGA

Splicing

Mapped reads in Artemis





Mapped reads to read count matrix

aagtaatga aggattga atgagtag taggattg taggattg atgagtag ctaggatt atgagtag ctaggatt tgactag tatggcg tatggcg aaaaatttttgacgtagcaggtaaatgactagcattgacaggtaaaaatttatggcgcgcgggccatatta geneB geneA

read count table

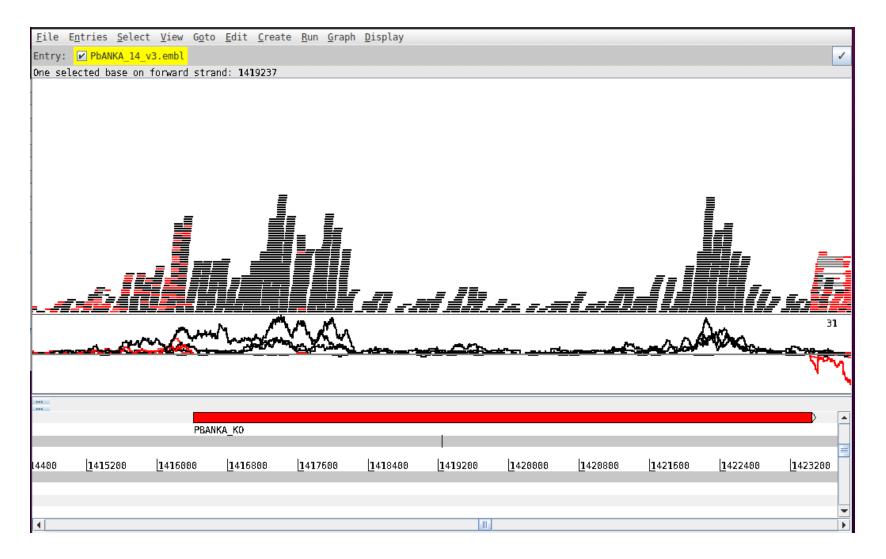
Genes	Condition1

read count table

Genes	Conditi on1	Cond 2

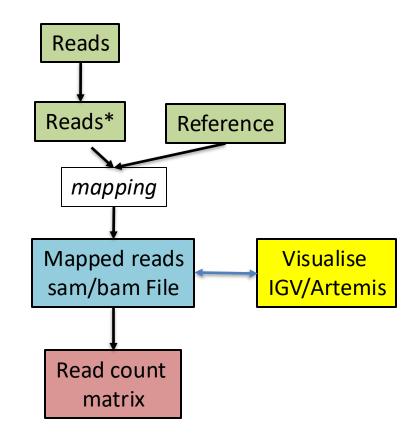
aagtaatga aggattga atgagtag taggattg atgagtag ctaggatt tgactag aggtaatgactagcattgaca	<u>ta</u>	tggcgc tggcg tggcgcgcgcggcca <mark>tat</mark>	ta
geneB	ger	neA	

Counting the reads...



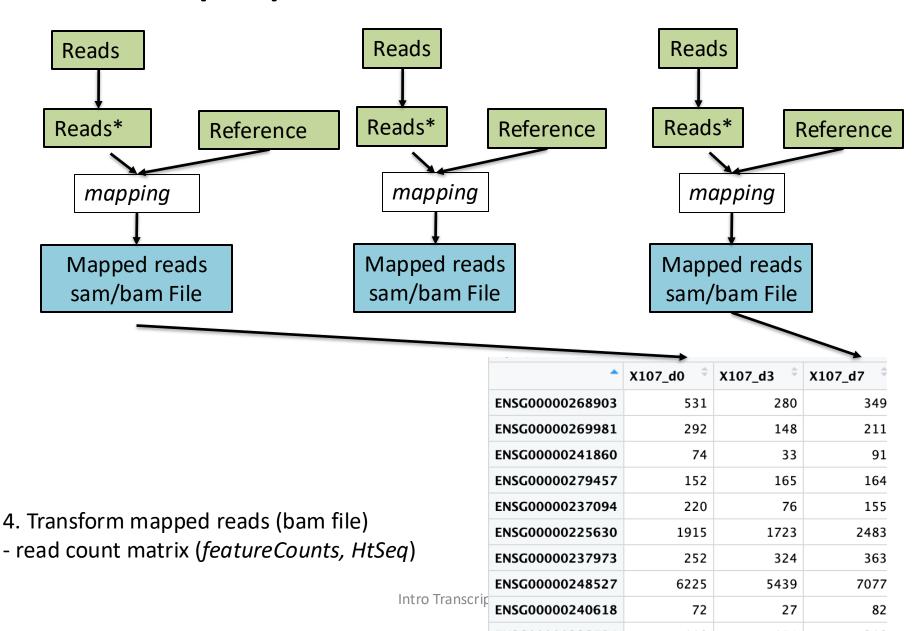
RNA-Seq pipeline

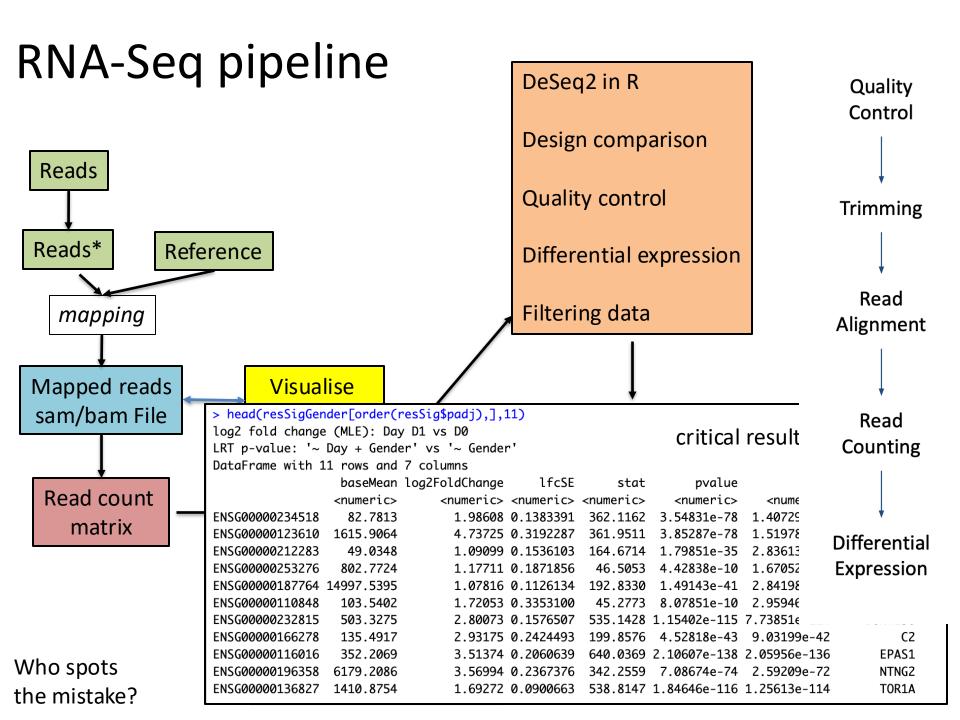
- 1. Quality control, Adapter trimming, low quality fastqQC, trimOmatic
- 2. Mapping (reads on reference) STAR, HiSat2
- 3. Transform mapping reads
- bam file (samtools)
- read count matrix (featureCounts, HtSeq)
- 4. Transform mapped reads (bam file)
- read count matrix (featureCounts, HtSeq)



^	X107_d0 [‡]
ENSG00000268903	531
ENSG00000269981	292
ENSG00000241860	74
ENSG00000279457	152
ENSG00000237094	220
ENSG00000225630	1915
ENSG00000237973	252
ENSG00000248527	6225

RNA-Seq Pipeline





Some definitions

- What does a RNA-Seq read represents?
- Why does that represent expression?
- What does the height of the coverage plot represents?

- 2 methods to count
 - read counts (how many reads map to each gene)
 - normalized : FPKM (fragments per kilobase of exon per million fragments mapped)

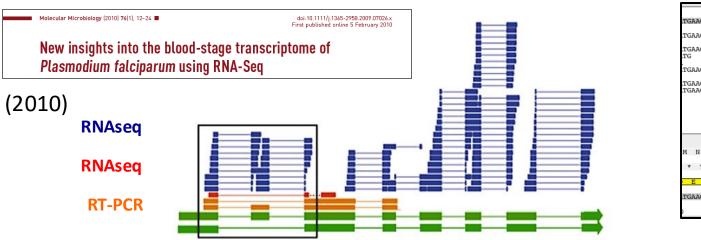
Overarching aim

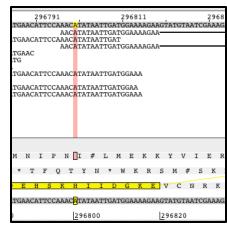
 We try to detect genes that are differentially expressed between two groups that can explain observed "differences"

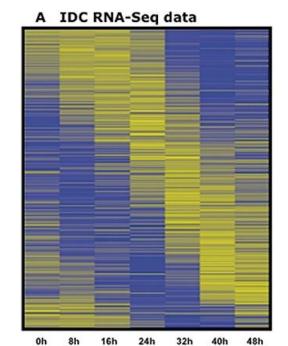
We need good quality data that we can trust

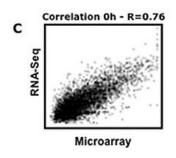
Two examples

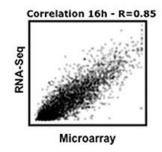
Plasmodium RNAseq

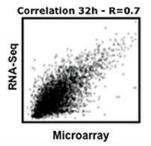








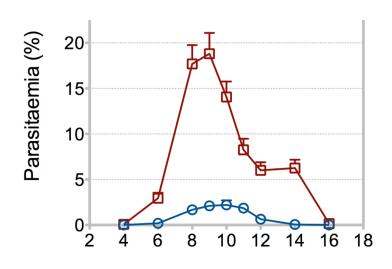




Berlin 2015

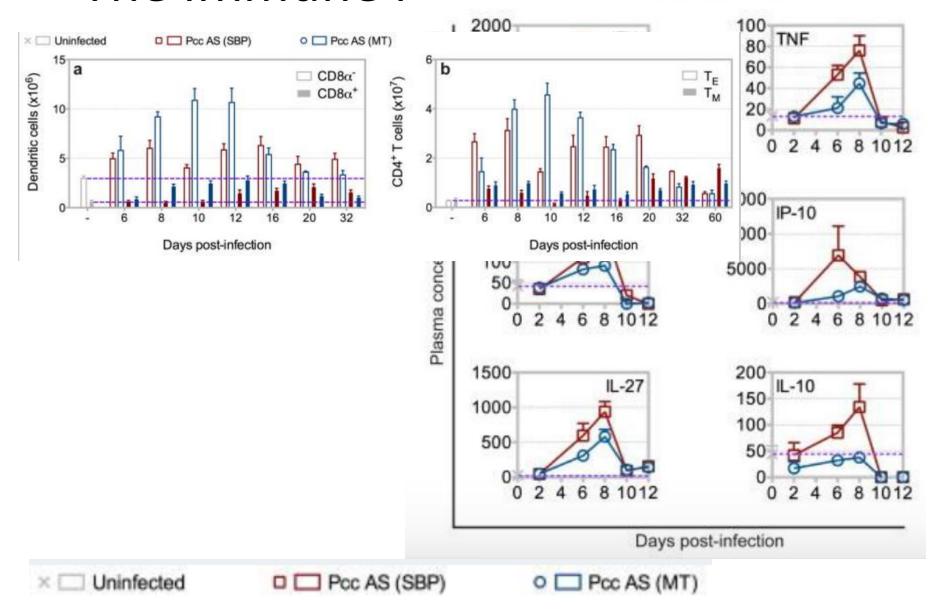
Differential expression: Vector transmission regulates immune control of virulence

- Serial blood passage of parasites leads to increased virulence
- Phil Spence & Jean Langhorne developed routine mosquito transmission of *P.c. chabaudi*.

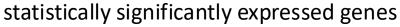


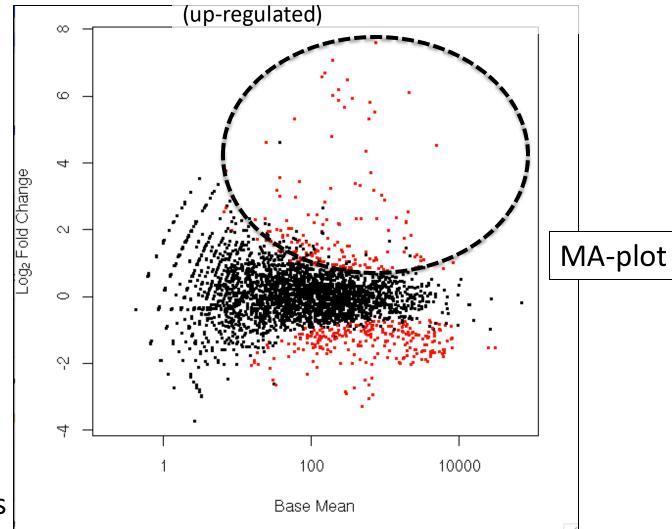
MT parasites show attenuated growth (not due to dose or pre-erythrocytic-stage of infection)

The immune resnance is different



Differential expression





3 biological replicates Use read counts

DESeq

assume negative binomial distribution of read counts

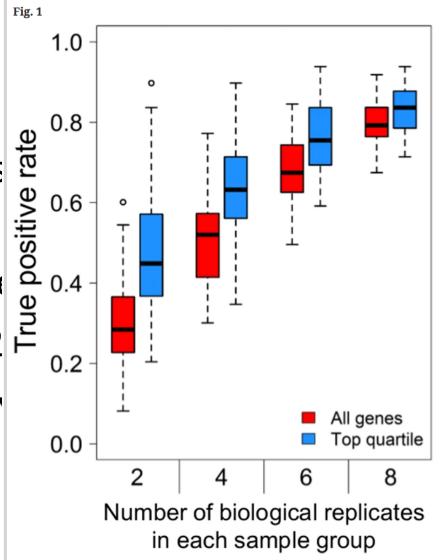
Intro Transcriptomics

Statistically differentially expressed?

- Need a test to understand which genes are "really" different
- Correction of multiple testing
- Tools: Sleuth, DEseq2, EDGER
- More robust results with more replicates

Statistically differ

- Need a test to unders "really" different
- Correction of multiple
- Tools: Sleuth, DEseq2
- More robust results v



Increasing numbers of sample replicates improves identification of schizont-stage genes varying in expression between different *P. falciparum* lines. Assessment of the proportion of genes captured as being differentially expressed between two different parasite clones (3D7 and D10), by taking 100 random samples of two, four, six and eight replicates of each (out of ten initially analysed replicates that identified 123 genes with \log_2 differences of > 2 in relative transcript levels between the two clones)

Intro Tra

Differential expression with DESeq2

- DESeq is an R package, part of Bioconductor
- Starts from read counts and does its own normalisation
- For each sample the size factor is calculated how many reads
- Dispersion estimated (how tied are the expression values)
- For DE uses a negative binomial distribution, and "performs an advanced t-test"
- We will focus to run the tools and have an idea of the functionality, rather than becoming statistician's!

Most significant

padj hgnc_symbol

GBP1

IL1RN

IFIH1

TDRD7

<numeric> <character>

Downregulated

How does the result looks like?

stat

pvalue

1421.24 7.24338e-308 1.02001e-303

1417.55 4.57195e-307 3.21911e-303

1280.36 2.69209e-277 1.26367e-273

1277.60 1.06827e-276 3.76085e-273

<numeric>

> head(resGender[order(resGender\$padj),],11)

LRT p-value: '~ Day + Gender' vs '~ Gender'

<numeric>

baseMean log2FoldChange

log2 fold change (MLE): Day D1 vs D0

DataFrame with 11 rows and 7 columns

ENSG00000117228 9186.124

ENSG00000136689 17782.754

ENSG00000115267 4708.550

ENSG00000196116 1635.390

```
ENSG00000170581 23333.275
                                4.07187 0.1401291
                                                    1270.88 3.06152e-275 8.62246e-272
                                                                                           STAT2
ENSG00000138035 1222.625
                                4.55961 0.1589723
                                                    1259.13 1.08754e-272 2.55246e-269
                                                                                            PNPT1
> head(resGender[order(resGender$log2FoldChange,decreasing = F),],11)
log2 fold change (MLE): Day D1 vs D0
LRT p-value: '~ Day + Gender' vs '~ Gender'
DataFrame with 11 rows and 7 columns
                 baseMean log2FoldChange
                                             1fcSE
                                                                   pvalue
                                                                                 padi hanc_symbol
                                                         stat
                                                                <numeric>
                                                                            <numeric> <character>
                <numeric>
                               <numeric> <numeric> <numeric>
ENSG00000131459
                   7.0171
                                -3.04647 0.582347
                                                     52.3370 2.53858e-11 1.08889e-10
                                                                                            GFPT2
ENSG00000182580
                  41.9070
                                -2.84842 0.260561 153.4997 4.63201e-33 6.75938e-32
                                                                                            EPHB3
                                                                                             TGM3
                 141.0757
                                -2.34984 0.419139
                                                     41.9035 4.20589e-09 1.42785e-08
ENSG00000125780
ENSG00000148600
                  46.3343
                                -2.32099
                                          0.264634
                                                    113.6641 1.78515e-24 1.85661e-23
                                                                                            CDHR1
                  70.7451
                                -2.26813
                                          0.221951
                                                    137.4909 1.31352e-29 1.69232e-28
                                                                                            WNT7A
ENSG00000154764
                                -2.20036
                                                                                         C16orf74
ENSG00000154102
                 142.0127
                                          0.181780
                                                    148.4089 5.80734e-32 8.16158e-31
ENSG00000279447
                  36.9222
                                -2.19621
                                          0.290009
                                                     70.4884 3.35492e-15 1.98839e-14
                  10.6137
                                          0.384055
                                                     44.4879 1.18876e-09 4.27478e-09
                                                                                            RAB7B
ENSG00000276600
                                -2.18503
                                          0.279827
ENSG00000261150
                  46.0225
                                -2.14882
                                                     70.2949 3.69071e-15 2.18006e-14
                                                                                            EPPK1
ENSG00000178947
                  22.8845
                                -2.08345
                                          0.239324
                                                     94.6484 2.19719e-20 1.80518e-19
                                                                                        SMTM10L2A
ENSG00000158270
                   9.9380
                                -2.05730
                                                     46.4807 4.48206e-10 1.68896e-09
                                                                                          COLEC12
                                          0.488719
```

lfcSE

<numeric> <numeric> <numeric>

4.85146 0.1558552

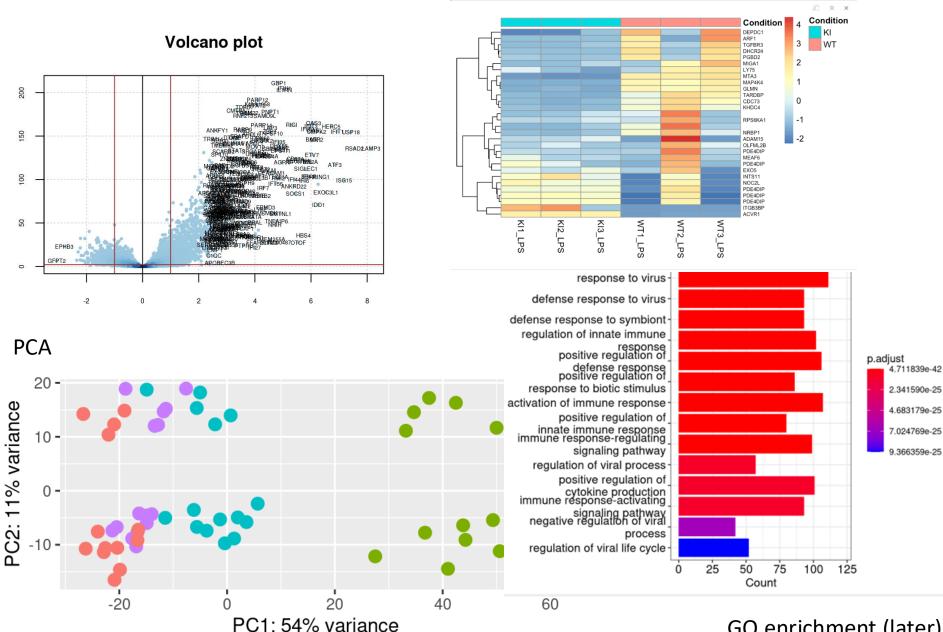
5.04210 0.1655905

5.05549 0.1640378

3.63350 0.1255556

Visualisation

heatmap most expressed genes



GO enrichment (later)

Differential expression

- Good biological question
- Good sample preparation
- Enough biological replicates (get information through published studies)
- Sequencing -> Quality control, fastqc and visualize the results
- Differential expression (DeSeq2, EdgeR, Sleuth)
- Check FDR, enrichment test (GO)
- Interpretation of the results
- Write the paper

DE genes, what to do now?

From Covid paper, CD14 cluster

```
5.387767e-159 5.0588481 0.998 0.233 7.571429e-155
  ISG15
## IFIT3
          1.945114e-154 6.1124940 0.965 0.052 2.733468e-150
          2.503565e-152 5.4933132 0.965 0.076 3.518260e-148
## TFT6
## TSG20
          6.492570e-150 3.0549593 1.000 0.668 9.124009e-146
## IFIT1
          1.951022e-139 6.2320388 0.907 0.029 2.741772e-135
## MX1
          6.897626e-123 3.9798482 0.905 0.115 9.693234e-119
## LY6E
          2.825649e-120 3.7907800 0.898 0.150 3.970885e-116
  TNFSF10 4.007285e-112 6.5802175 0.786 0.020 5.631437e-108
          2.672552e-108 5.5525558 0.786 0.037 3.755738e-104
## IFIT2
## B2M
          5.283684e-98 0.6104044 1.000 1.000 7.425161e-94
## PLSCR1 4.634658e-96 3.8010721 0.793 0.113 6.513085e-92
## IRF7
         2.411149e-94 3.1992949 0.835 0.187 3.388388e-90
## CXCL10
          3.708508e-86 8.0906108 0.651 0.010 5.211566e-82
## UBE2L6
          5.547472e-83
                        2.5167981 0.851 0.297 7.795863e-79
## PSMB9
          1.716262e-77 1.7715351 0.937 0.568 2.411863e-73
```

Study them one by one



Enrichment analysis: what do we need?

It is a computational biology method that identifies biological functions that are overrepresented in a group of genes more than would be expected by chance

DE genes from your dataset

Where I can find info about "group of genes"?

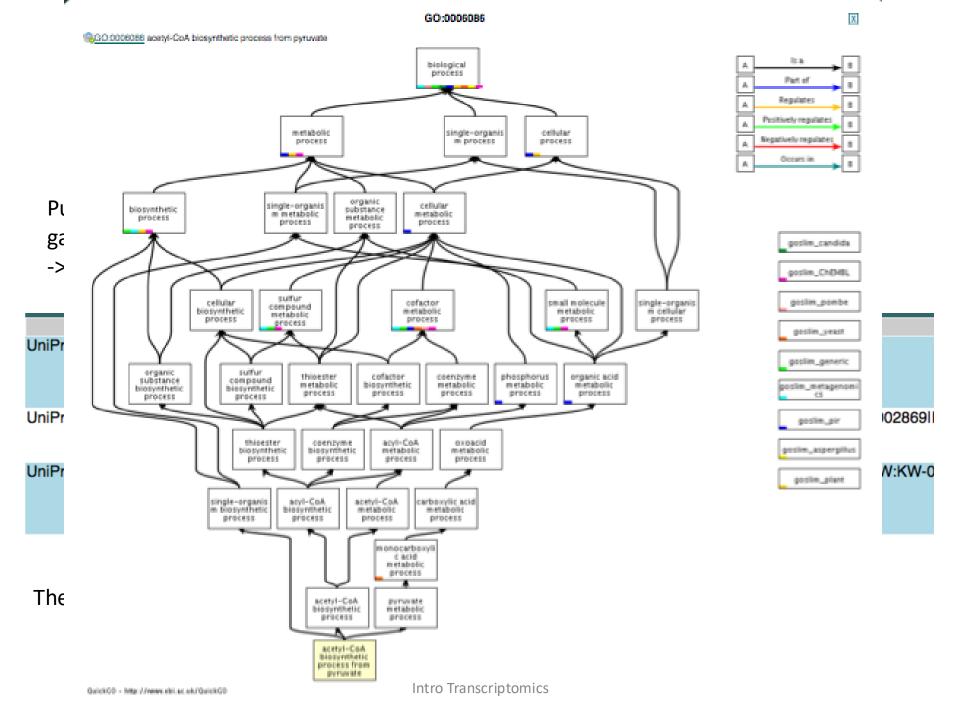
Databases:

- KEGG
- GO
- StringDB
- Reactome
- DAVID knowledge
- **—** ..

- From where are the info?
 - Literature
 - Automatic
 - Manual
 - Experimental
 - Modelling/Prediction

Stop here? what is GO?

- Gene Ontology
- Controlled vocabulary
- Saved as a graph
- The terms are grouped into three categories:
 - molecular function (describing the molecular activity of a gene)
 - biological process (describing the larger cellular or physiological role carried out by the gene, coordinated with other genes)
 - cellular process (describing the location in the cell where the gene product executes its function).



Or in R

Good =matrix(data=c(100,900,500,11500),nrow=2)
colnames(Good) = c("DE", "Rest")
rownames(Good) =c("in Pathway","somewhere else")
fisher.test(Good)

• p-value < 2.2e-14

_	DE [‡]	Rest [‡]
in Pathway	100	500
somewhere else	900	11500



NotG	
=matrix(data=c(100,900,1100,119000),nrow=2)
colnames(NotG)	= c("DE", "Rest")
rownames(Note	G) =c("in Pathway","somewhere
else")	
fisher.test(NotG))

?

•	DE ‡	Rest [‡]
in Pathway	100	1100
somewhere else	900	11900
		Rati



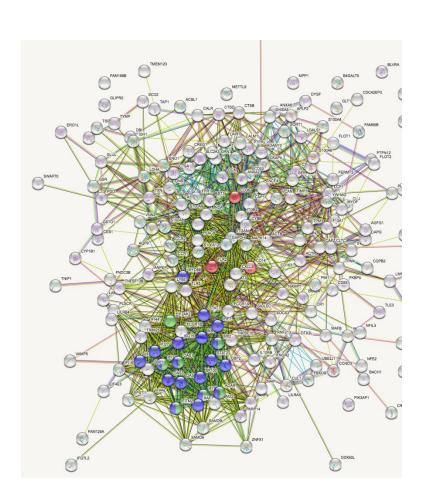
Ratio is The same

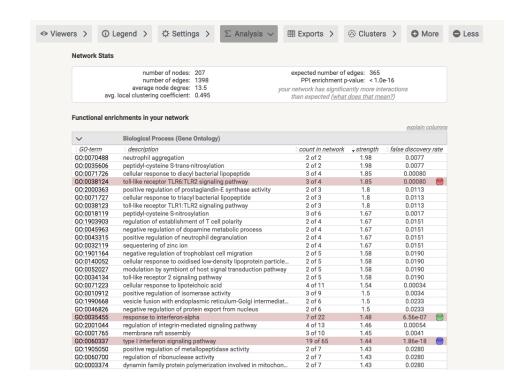


Multiple correction

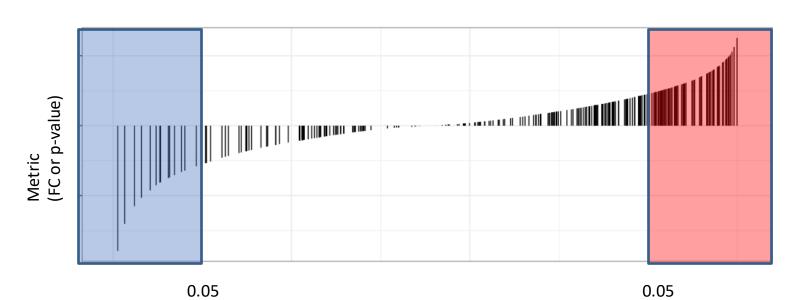
- If you do more test, you need to apply correct for multiple testing!!!
- If you have four conditions A B C D E and do
- A vs B; A vs C; A vs D; A vs E; B vs C; B vs D; B vs E; C vs D; C vs E & D vs E
 you have to correct ALL P-values for the 10 amount of test.
- Most conservative method is to multiple so an adjusted P-value of 0.005 becomes 0.05 and is on the boardline of the false discovery rate!
- Tools normally do multiple correction, but not between "experiments" of the user

Example 2 StringDB





DE genes/p-value and pathways relationship...



Downregulated

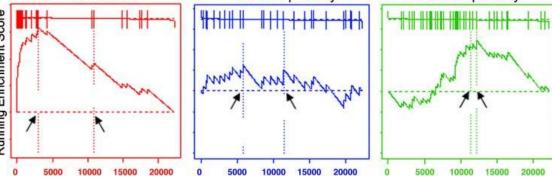
What if <u>small but coordinated</u> changes in sets of functionally related genes may also be important???

Upregulated

GSEA

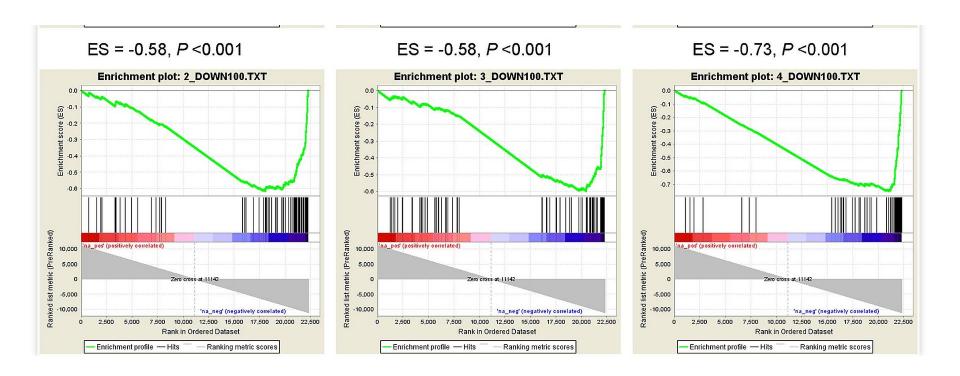
- Gene set enrichment analysis (GSEA)
 (also functional enrichment analysis)
- Uses the complete gene list (by P-value)
- Orders genes by P-values (down regulatec)

• Test the distribution of light terms of



https://www.gsea-msigdb.org/gsea/index.jsp

Examples



https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4470683/

Questions

- 1. What is the power of RNA-Seq and DE analysis?
- 2. Why do we need good quality data?
- 3. What is the importance of biological replicates?
- 4. Why is statistical analysis DE / correct of P-values so important?

Consider

- How powerful are enrichments? In some cases these make more sense than in others
- Databases are based on existing knowledge
- Annotation is not evenly distributed: There is more research in T-cell in cancer than in a rare novel subtype in rheumatoid arthritis — this might distracted
- Enrichments gives you a general overview, but you need to then dig deeper
- Check that your "most significant genes" for the gene list are also part of the enrichments
- In the end, it stills comes down to how much you and your collaborator know!

Conclusions (3/n)

- RNA-Seq is a powerful tool to access expression of all genes (and other elements miRNA)
- Differential expression analysis to find genes changing expression due to different condition
- Very population in Immunology
- In the exercise you are going to perform differential expression from a read count matrix with DESeq2
- You need good quality data and enough biological replicates to obtain statistically relevant results