# Testing for Differential Expression with DESeq2
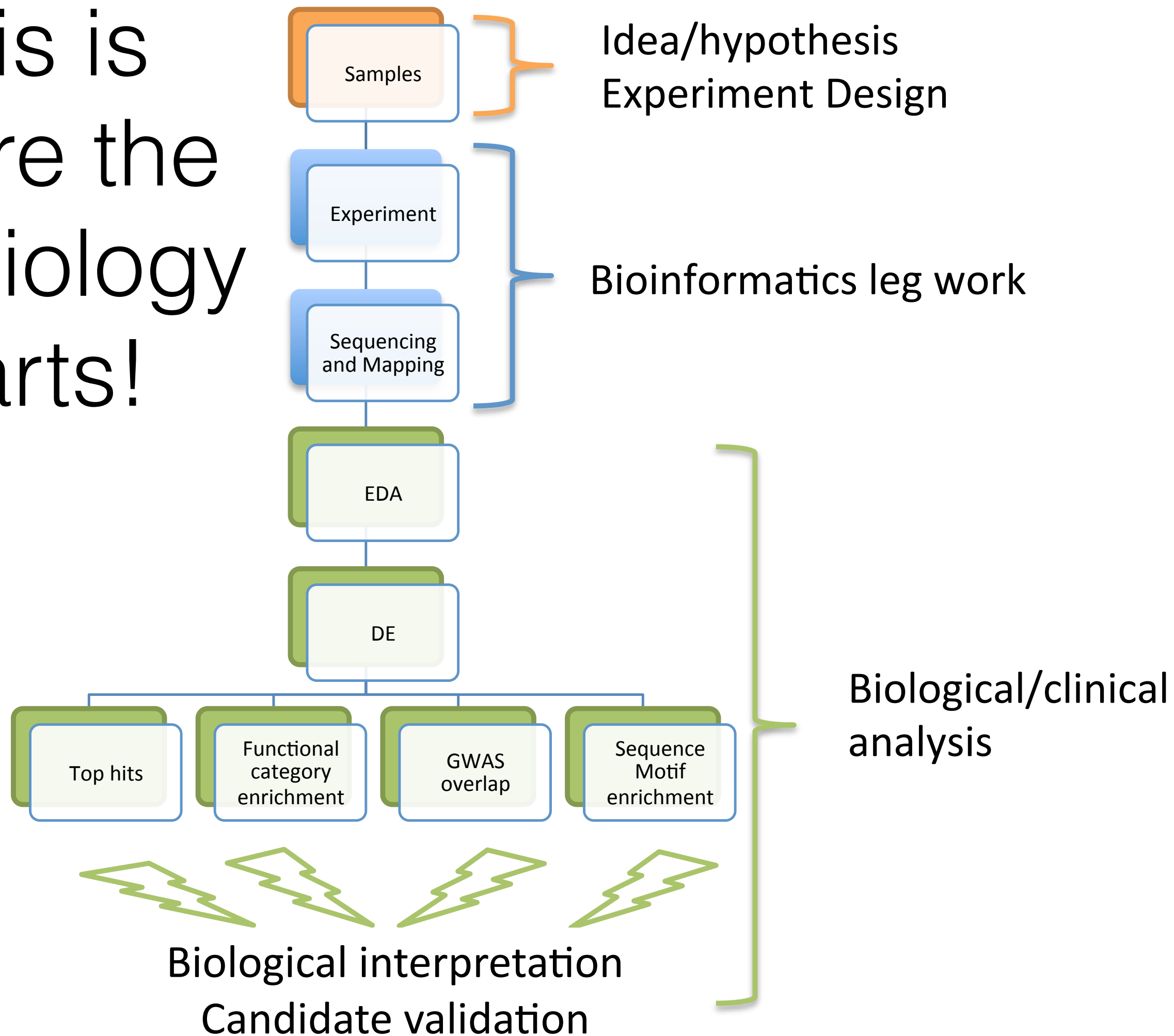
Laura Pikkupeura

# Teacher

- Laura Pikkupeura

- PhD student in Sandelin and Jensen Labs

- Mail: laura.pikkupeura@bric.ku.dk

- Background:

  - MSc Molecular Biomedicine

  - PhD Stem cell biology, genomics and transcriptomics

# Lecture outline

- **Theory:**

  - Introduction to testing in transcriptomics

  - Overview of DESeq2

- **Practical:**

  - Gene expression in the fissing yeast stress response

  - Gene expression in Inflammatory bowel disease.

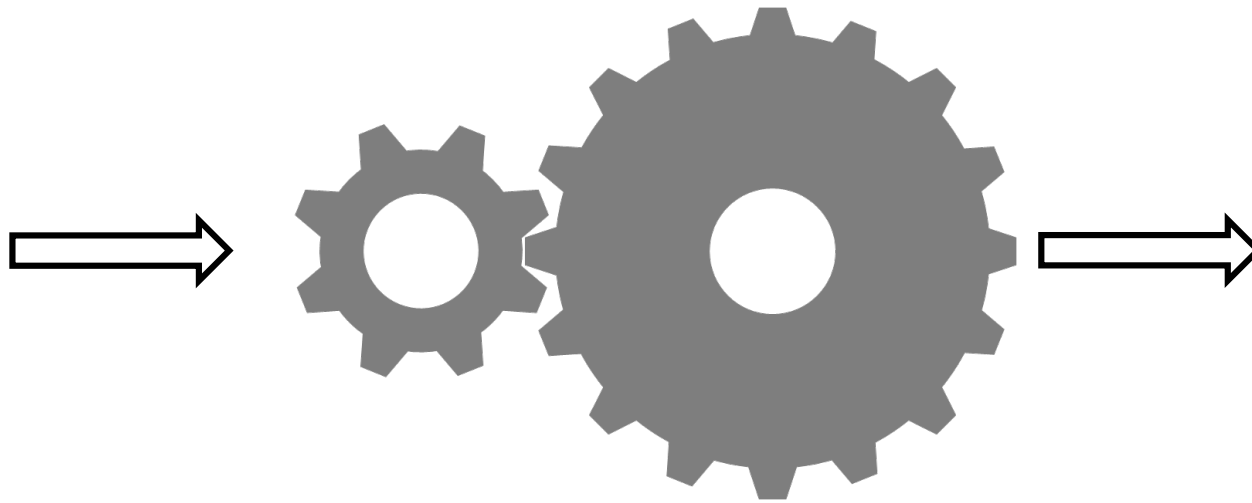  - DESeq2 applied to pseudo-aligned reads.

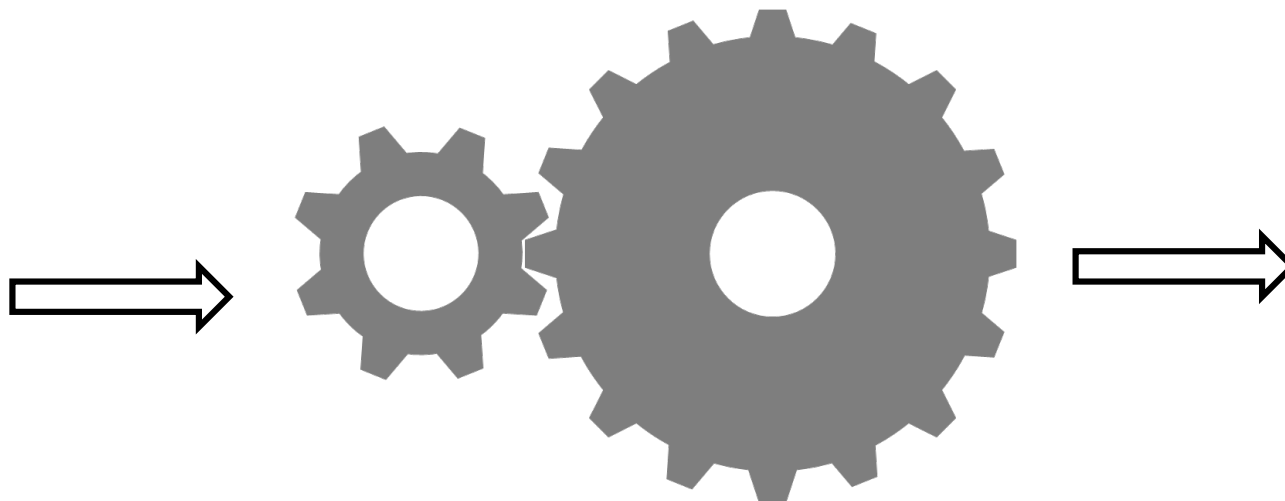# Importance of
# QC and filtering



Cleaned,
Trimmed and
prepared data
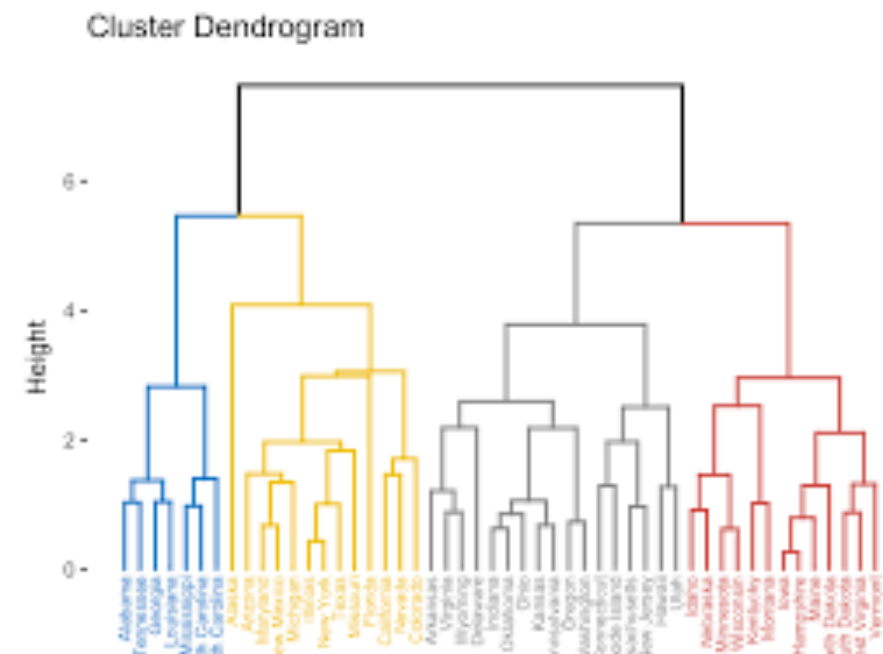
Nature publication

Shitty data

Even more
shitty data

# Testing in transcriptomics

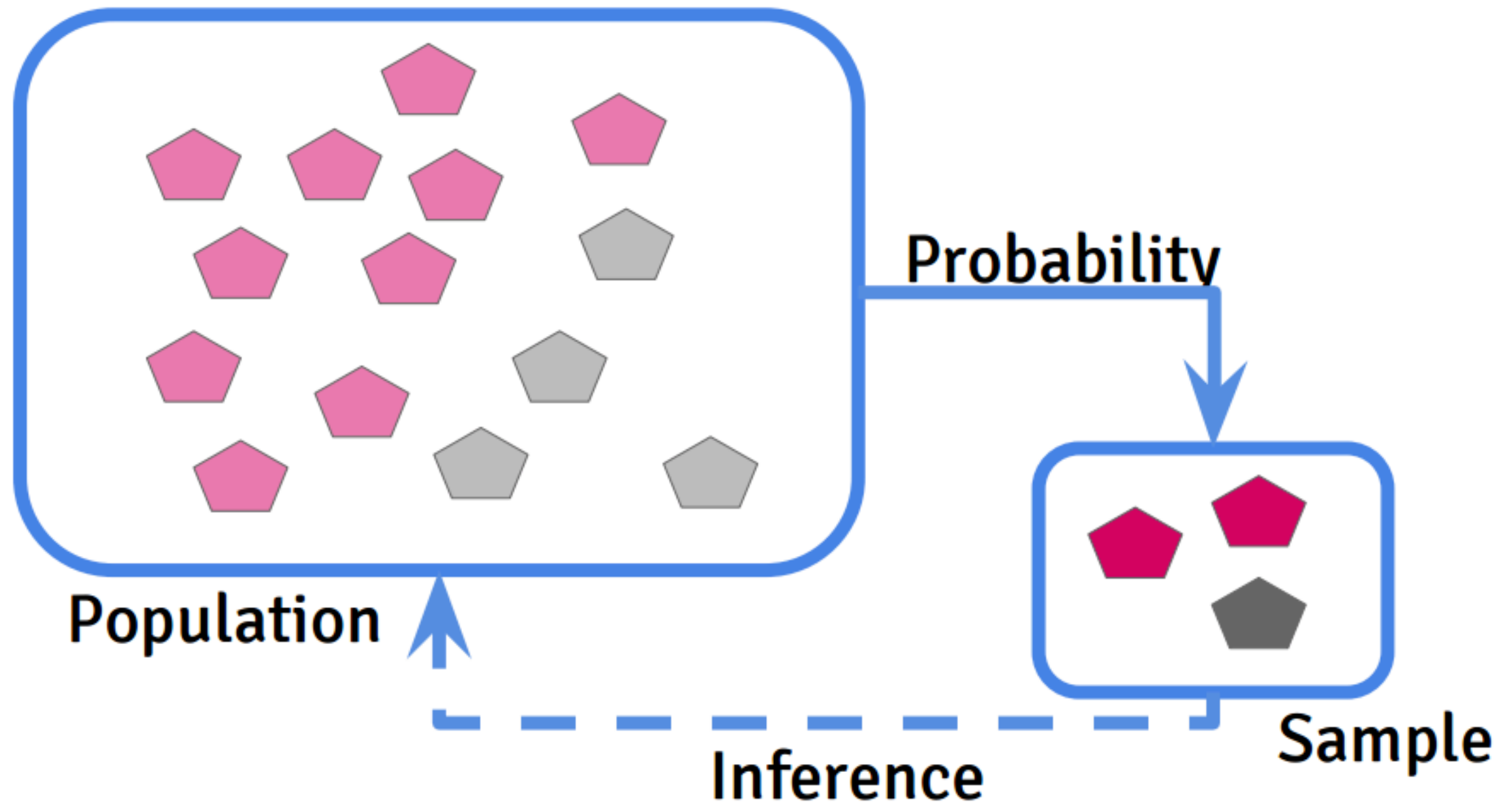# EDA

**Exploratory data analysis deals with *global differences***



**But what if we want to say something about whether *a single gene* is changing?**

# The basic idea of testing

# The basic idea of testing



~3 million men
~3 million women

1000s

# Effect size vs significance



**High variability**

**Medium variability**

**Low variability**

**Effect Size:**
Magnitude of the effect

**Significance (p-value):**
How much we believe the effect

*Large sample size + low variance  > even small effect sizes can be significant!*

# The basic idea of testing

# Different types of replication

# The use of technical replicates

- **Technical replicates:**

    - Useful for EDA.

    - Measure technical noise.

    - Determine whether the sample pipeline is stable - technical replicates should be most similar.

    - Samples are not independent: Should not be (directly) included in the DE analysis.

- **Biological replicates:**

    - Useful for DE analysis.

    - Determine the biological variation in gene expression.

# The inherent problem of testing in high-throughput data

- Very small sample size (usually 3 replicates per condition).
- Extremely high number of tests – problems with multiple testing.

*x 20.000*

# DESeq2

# Why DESeq2?



Bioconductor
OPEN SOURCE SOFTWARE FOR BIOINFORMATICS

**Top 75**

1 BiocInstaller (30871)
2 BiocGenerics (24259)
3 IRanges (21715)
4 S4Vectors (21507)
5 Biobase (20177)
6 AnnotationDbi (18365)
7 zlibbioc (16565)
8 GenomicRanges (15948)
9 limma (15149) ←
10 XVector (14705)
11 GenomeInfoDb (13745)
12 Biostrings (13522)
13 BiocParallel (13081)
14 SummarizedExperiment (12860)
15 annotate (10588)
16 GenomicAlignments (10213)
17 biomaRt (10090)
18 rtracklayer (9969)
19 Rsamtools (9770)
20 genefilter (9552)
21 DelayedArray (9180)
22 GenomicFeatures (8799)
23 graph (8539)
24 edgeR (7979) ←
25 preprocessCore (7447)

26 DESeq2 (7296) ←
27 geneplotter (7050)
28 affy (5888)
29 BSgenome (5812)
30 affyio (5525)
31 rhdf5 (5272)
32 RBGL (5122)
33 multtest (5047)
34 Rgraphviz (4985)
35 VariantAnnotation (4774)
36 impute (4656)
37 qvalue (4248)
38 AnnotationHub (4171)
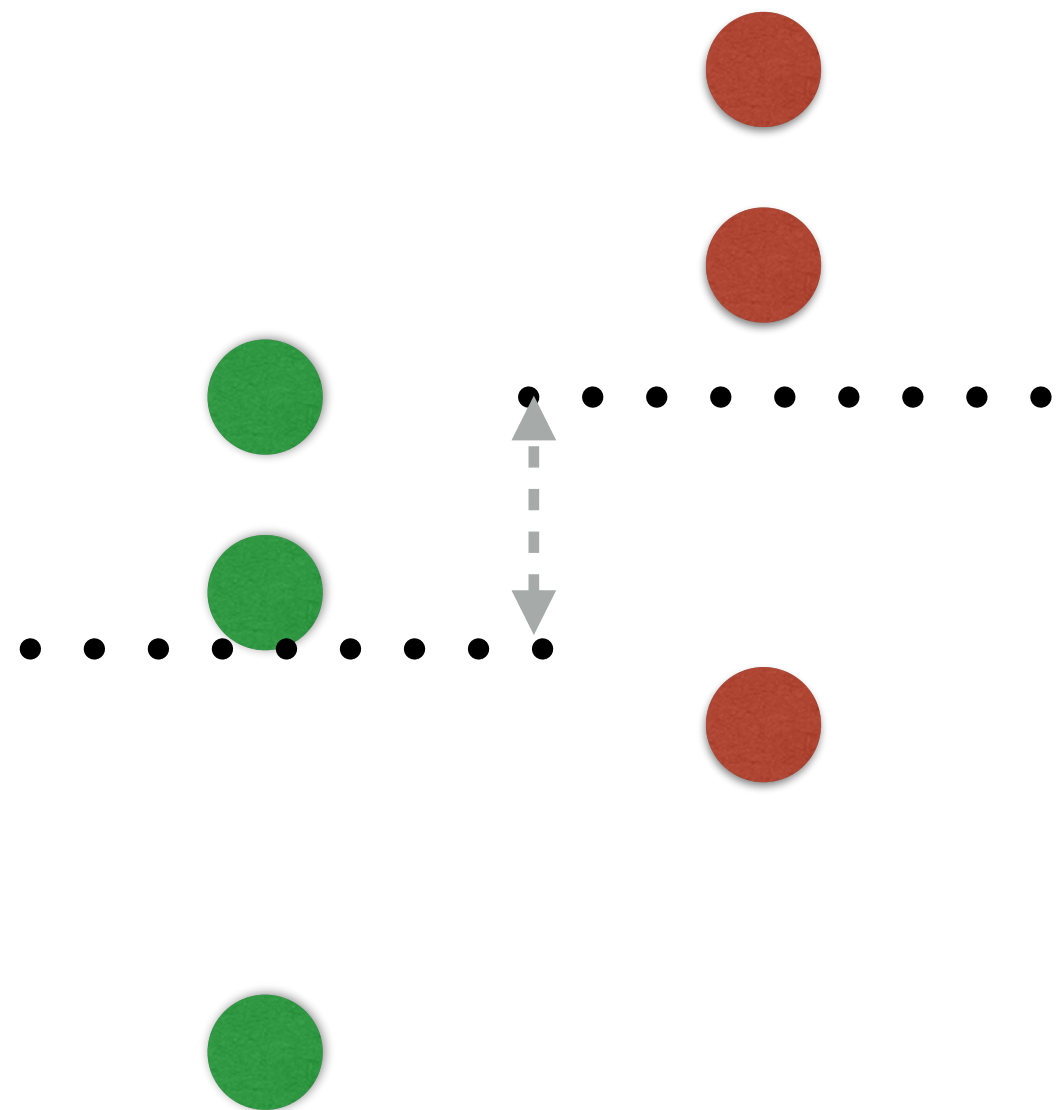39 GEOquery (4032)
40 ShortRead (3907)
41 ensembldb (3727)
42 interactiveDisplayBase (3569)
43 ProtGenerics (3395)
44 DNAcopy (3300)
45 GSEABase (3203)
46 DESeq (3116) ←
47 biovizBase (3050)
48 sva (2791)
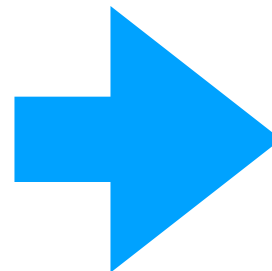49 Gviz (2657)
50 KEGGREST (2632)

- Third most downloaded R-package for DE

- Benchmarks show good performance in most cases

- Well made vignettes and guide available

- Defaults work well out of the box

- Same developers as Salmon, resulting in a highly optimised Salmon > tximport > DESeq2 pipeline.

# What DESeq2 does (tl;dr)

**Samples**

**Genes**

| | | | | | |
|---|---|---|---|---|---|
| 80822 | 90845 | 54240 | 52896 | 4729 | 54654 |
| 81205 | 181063 | 54218 | 104702 | 4849 | 54568 |
| 81394 | 91434 | 54175 | 52962 | 4877 | 55282 |
| 116837 | 205921 | 75741 | 100362 | 23484 | 122619 |
| 10047 | 82464 | 8674 | 11691 | 966 | 6285 |
| 796 | 4922 | 2404 | 3427 | 1003 | 2728 |
| 5185 | 6827 | 10631 | 4192 | 705 | 7515 |
| 35167 | 23788 | 31257 | 13726 | 1657 | 10852 |
| 35786 | 23955 | 31758 | 13857 | 1640 | 10666 |
| 2680 | 6431 | 2863 | 3709 | 634 | 3627 |
| 887 | 4837 | 4694 | 4562 | 852 | 3011 |
| 7 | 39 | 3 | 6 | 1 | 15 |
| 13 | 27 | 7 | 3 | 0 | 22 |

**+**

**~ Experimental groups**

**Effect Size**

**Significance**

| | log2FoldChange | Pvalue |
|---|---|---|
| ENSG00000274443 | 18.65756435 | 1.56E-08 |
| ENSG00000274766 | 18.65756435 | 1.56E-08 |
| ENSG00000275953 | 18.65756435 | 1.56E-08 |
| ENSG00000276375 | 18.65756435 | 1.56E-08 |
| ENSG00000278520 | 18.65756435 | 1.56E-08 |
| ENSG00000279173 | 18.65756435 | 1.56E-08 |
| ENSG00000279302 | 18.65756435 | 1.56E-08 |
| ENSG00000279756 | 18.65756435 | 1.56E-08 |
| ENSG00000202019 | 18.86227633 | 1.08E-08 |

# Components of DESeq2

1. Normalisation.

2. Estimating dispersion.

3. Sharing information across genes.

4. Testing (Using GLMs and Wald Tests).

5. Correction for multiple testing with independent filtering.

# Counts Per Million (CPM)

- As you already know, the individual samples in the EM cannot be compared before they have been **normalised**.
- The main reason for this is differences in **library size** or **sequencing depth** between samples in the same experiments.
- The basic way of dealing with this problem is to consider counts as fractions of total, usually scaled to a human scale unit, such as counts-per-million (cpm) or tags-per-million (tpm)

Samples →

Genes ↓

Counts

colSums/
Library sizes

# RNA composition

## Without DE

```
##    sample1 sample2 sample3
## 1      10      22      31
## 2      20      42      61
## 3      30      62      91
## 4      10      22      31
## 5      10      22      31
## 6      10      22      31
```

```
##          sample1   sample2   sample3
## [1,]  0.1111111 0.1145833 0.1123188
## [2,]  0.2222222 0.2187500 0.2210145
## [3,]  0.3333333 0.3229167 0.3297101
## [4,]  0.1111111 0.1145833 0.1123188
## [5,]  0.1111111 0.1145833 0.1123188
## [6,]  0.1111111 0.1145833 0.1123188
## attr(,"scaled:scale")
## sample1 sample2 sample3
##      90     192     276
```

## With DE

problems --> add a lot of count to many gene (eq to exp chnage in many genes), it will affect DE of other gene when we

use total count for normalization

we need assumption below

```
##    sample1 sample2 sample3
## 1      10      22      31
## 2      20      42      61
## 3      30      62      91
## 4      10     110     124
## 5      10     110     124
## 6      10     110     124
```

```
##          sample1   sample2    sample3
## [1,]  0.1111111 0.04824561 0.05585586
## [2,]  0.2222222 0.09210526 0.10990991
## [3,]  0.3333333 0.13596491 0.16396396
## [4,]  0.1111111 0.24122807 0.22342342
## [5,]  0.1111111 0.24122807 0.22342342
## [6,]  0.1111111 0.24122807 0.22342342
## attr(,"scaled:scale")
## sample1 sample2 sample3
##      90     456     555
```

DESeq2 solves this by estimating normalisation factors.
Importantly, this **assumes that most genes are not DE!**

# 2. Dispersion

$Y_{gj}$ is the count for gene $g$ in sample $j$

$$Y_{gj} \sim NB\left(\mu_{gj}, \mu_{gj} + \mu_{gj}^2 \phi_g\right)$$

negative binomial dist

Where $u_{gj}$ is the mean expression:

$$\mu_{gj} = M_j p_{pj}$$

Where $M$ is the library size and $p_{pj}$ the fraction of the gene within M.

$$BCV = \sqrt{\phi_g}$$

**Intuitively:**

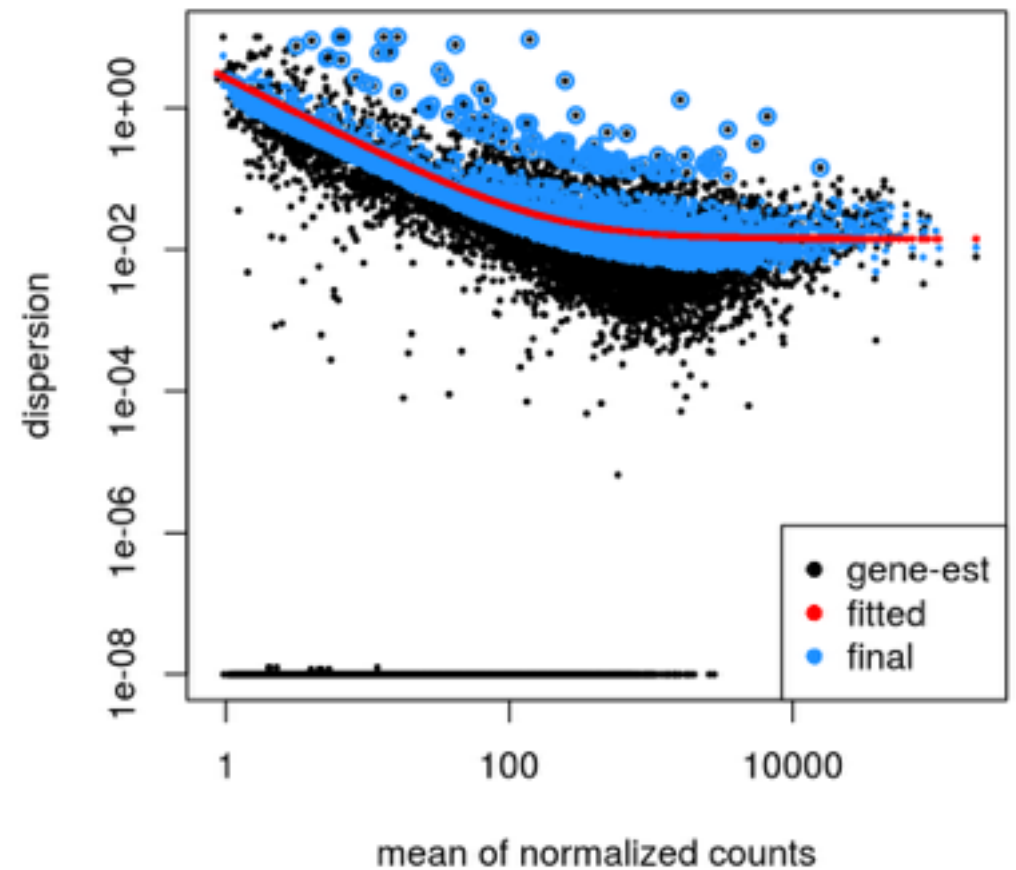- Gene expression varies between samples due to sequencing noise (Modelled with a poisson distribution)

- Gene expression varies between samples due to biological variation

- DESeq2 models gene expression as a negative binomial distribution, where the biological coefficient of variation (BCV) measures excess variance compared to poisson

$$\text{Total CV}^2 = \text{Technical CV}^2 + \text{Biological CV}^2$$

# 3. Information sharing

1.  It is difficult to estimate the variance of the expression of a single gene with very few replicates

2.  DESeq2 solves this issue by assuming genes with the same expression tend to have the same variance (BCV)

3.  By "shrinking" the individual estimates of variance towards the global trend DESeq2 obtains better downstream results.



*\* Additionally, DESeq2 also filters out outliers using Cook's distance*

# 4. Testing

**Black box warning!**

DESeq2 uses
generalised linear models (GLMs)
to test for differential expression,
allowing for highly complex
experimental setups.

Importantly, this allows for
correction of *batch effects* as we
discussed in the PCA lecture.

Here, we will only deal with
simple group-wise comparisons
in abstract terms

*If you want to know more, take the Advanced Bioinformatic Course in block 1!*

# 4. Testing

- Once shrunken dispersions estimates have been obtained, DESeq2 uses the Wald test to test for differences between groups

- Hypotheses:

  - $H_0$: $\log_2$FC between two groups is zero

  - $H_1$: $\log_2$FC between two groups is not zero

- Alternatively, the threshold can also be set to something other than zero:

  - $H_0$: $\log_2$FC between two groups is less than 1 (in either direction)

  - $H_1$: $\log_2$FC between two groups is more than 1 (in either direction)

- The results is an estimate of the **logFC** and a **p-value** for each gene!

# Summary on Dispersion and Testing

Var(Expression) = Across group variability + Biological Variability + Measurement Error

Effects of conditions, this is what we want to find.

Variation in expression between individuals.

Variation in expression due to sequencing.
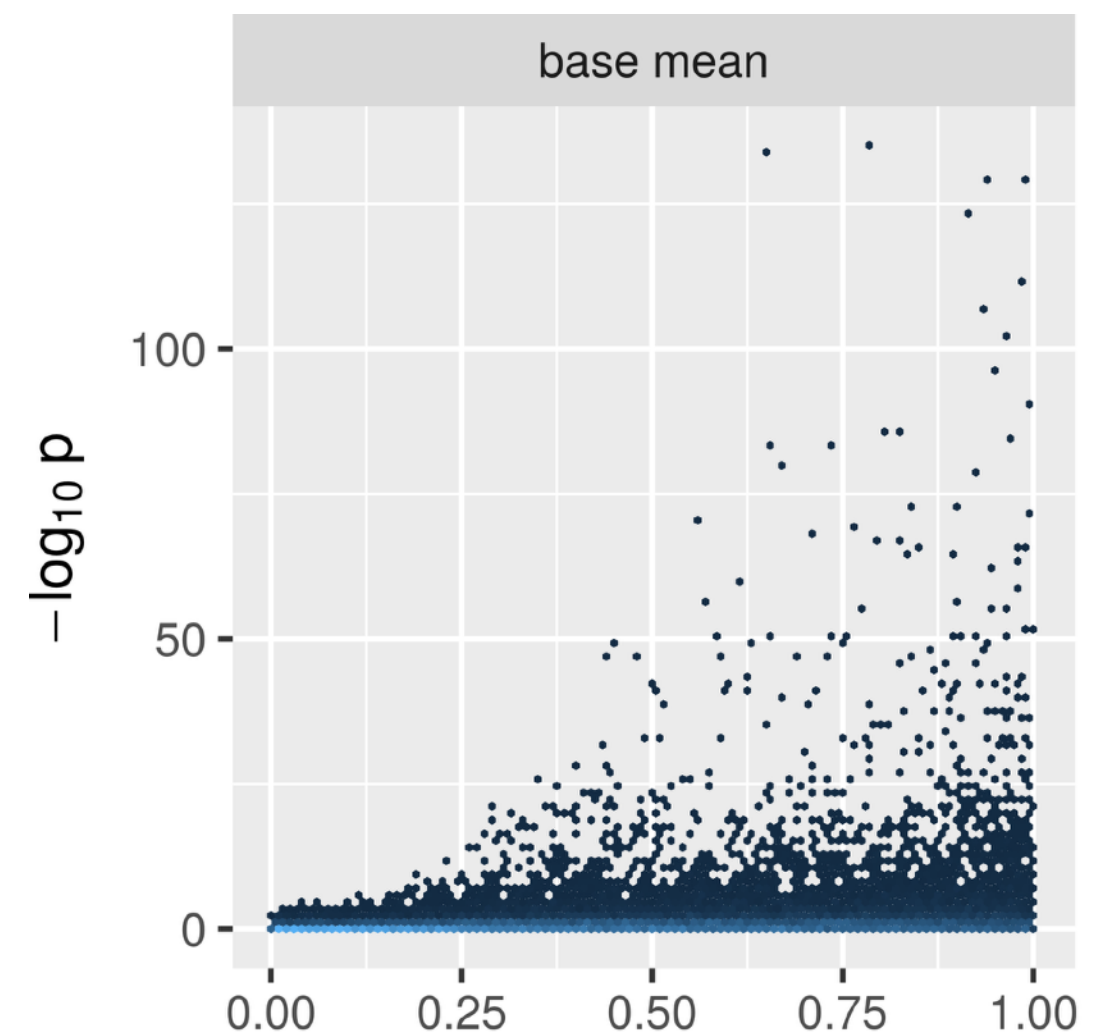
This is what we test for!

Shared across genes.

Sample independent (hopefully!).

Dominant for highly expressed genes

Dominant for lowly expressed genes.
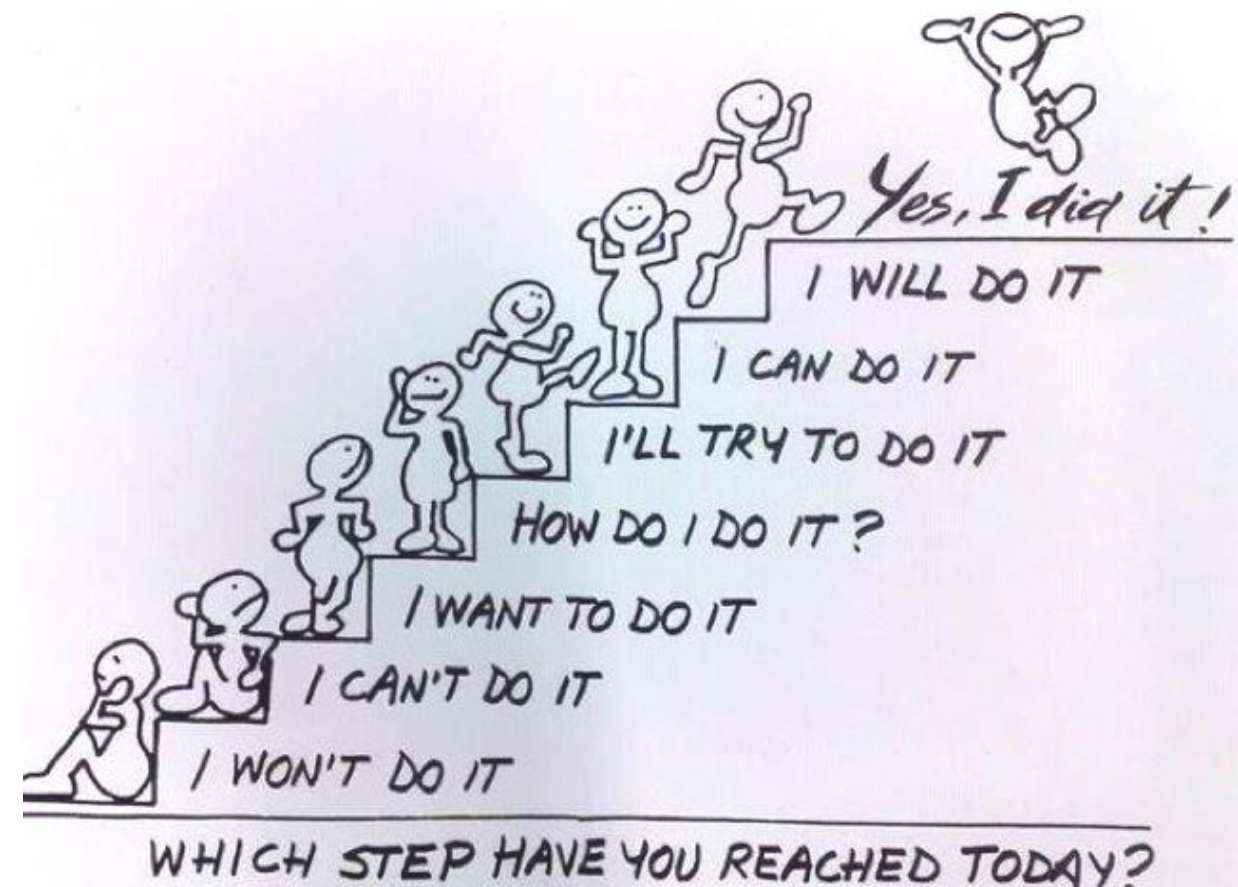
# Correction for multiple testing

- As we are testing thousands of genes, we must correct for multiple testing.

- DESeq2 be defaults uses an extension of the FDR-correction procedure based on *independent filtering:*

  - Due to sequencing noise, lowly expressed features will often vary too much to detect DE even very high logFCs.

  - Including many extra lowly expressed genes therefore decreases the number of DE genes due to more strict correction for multiple testing.

  - DESeq2 removes lowly expressed genes before FDR-correction to maximise the number of genes detected as DE.

- **Statistical note:** This is only valid if the filtering statistic (mean expression in this case) is *independent under the null hypothesis* (there is no relation between mean expression and tendency to be DE)



**Genes ranked by mean expression**

# Components of DESeq2

1. Normalisation.

2. Estimating dispersion.

3. Sharing information across genes.

4. Testing (Using GLMs and Wald Tests).

5. Correction for multiple testing with independent filtering.

# Other functionalities

- Variance-stabilizing transformations: An alternative to the standard log2-transformation (vst and rlog).

- Shrink *both* dispersion and logFCs

- Make PCA-plots

- Use a more advanced independent filtering approach called independent hypothesis weighting (IHW)

- Test complex multifactorial designs using Likehood ratio tests.

- New: Tissue decomposition (unmix)

# DESeq2 practical

- Example of complete analysis starting from a count matrix:

  - Setting up the data.

  - Running the DESeq2 functions.

  - Producing essential diagnostic plots.

  - Inspecting the results.

- You will redo the analysis on another dataset

- Demonstration of how to use tximport to import Salmon quantifications into DESeq2.

# Warning!

- DESeq2 is based on Bioconductor

- Bioconductor uses a quite rare R-programming system called S4.

- **Bioconductor IS NOT the tidyverse!**

- **DO NOT** expect everything to be a tibble!

- You might to read help files for the functions to know what to do!

# Next lecture

- Now you have used DESeq2 to find genes that show statistically significant changes between groups.

- Can you say something about these sets of genes?

- You can look at whether the genes as a group are associated to some known biological pathway or process.

- Monday you will look at **GO-term enrichment.**