

# RNA-seq data analysis: From counts to differentially expressed genes using edgeR-quasi

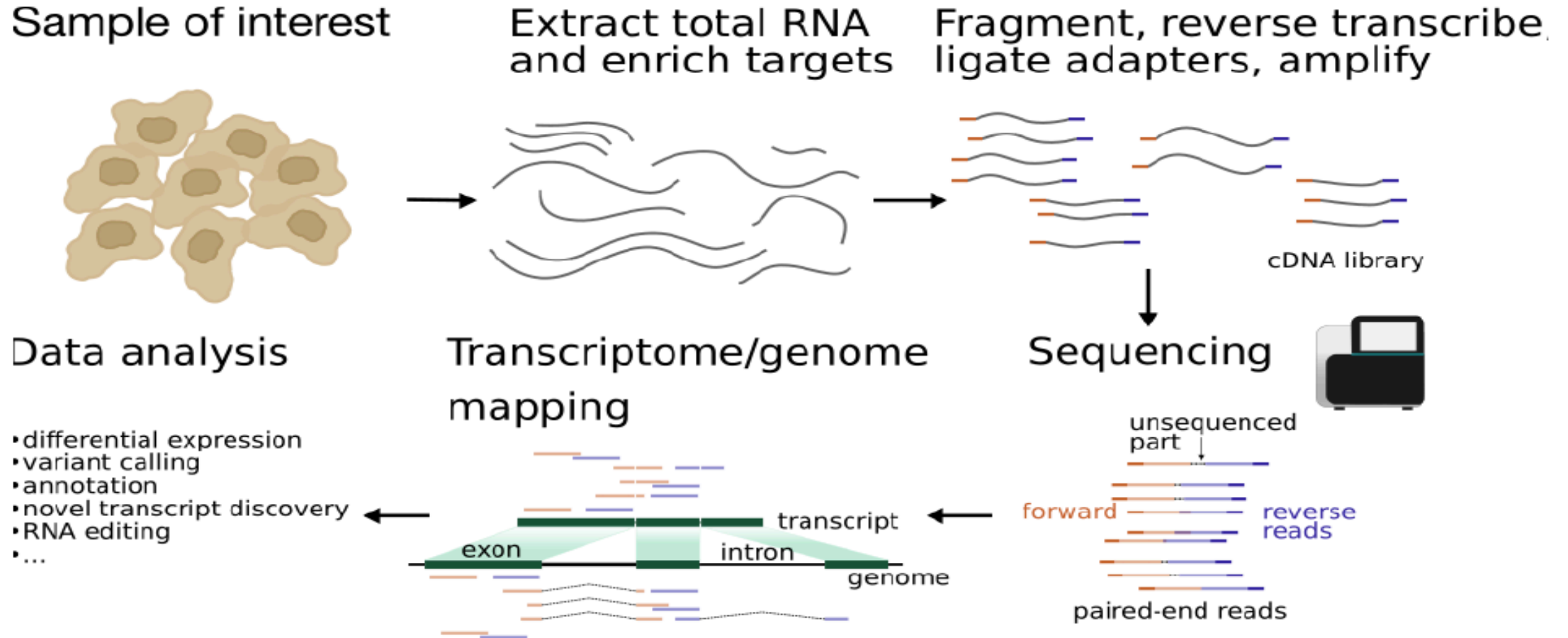
Gladstone Institutes

Krishna Choudhary  
Bioinformatics Core @ GIDB  
November 18, 2019

# Assumed background

- ◆ Familiarity with R and RStudio
- ◆ Familiarity with RNA-seq protocol
- ◆ Familiarity with basic concepts of hypothesis testing

# Typical protocol



# Experiment design influences data analysis. (should be planned to address relevant questions)

- ◆ What is the biological question that we seek to answer?
- ◆ How many tissue types and/or time points to compare?
- ◆ How deep should we sequence?
- ◆ Read length?
- ◆ Which sequencing platform?
- ◆ Single-end or paired-end?
- ◆ Pooling?
- ◆ Biological replicates?
- ◆ Technical replicates?
- ◆ Additional considerations?

Not the subject matter today!

- Workshop by Reuben Thomas:  
*Intro to statistics and experimental design.*
- Reading material in Dropbox:  
*RNA sequencing data : hitchhiker's guide to expression analysis* by Berge *et al.*, 2018

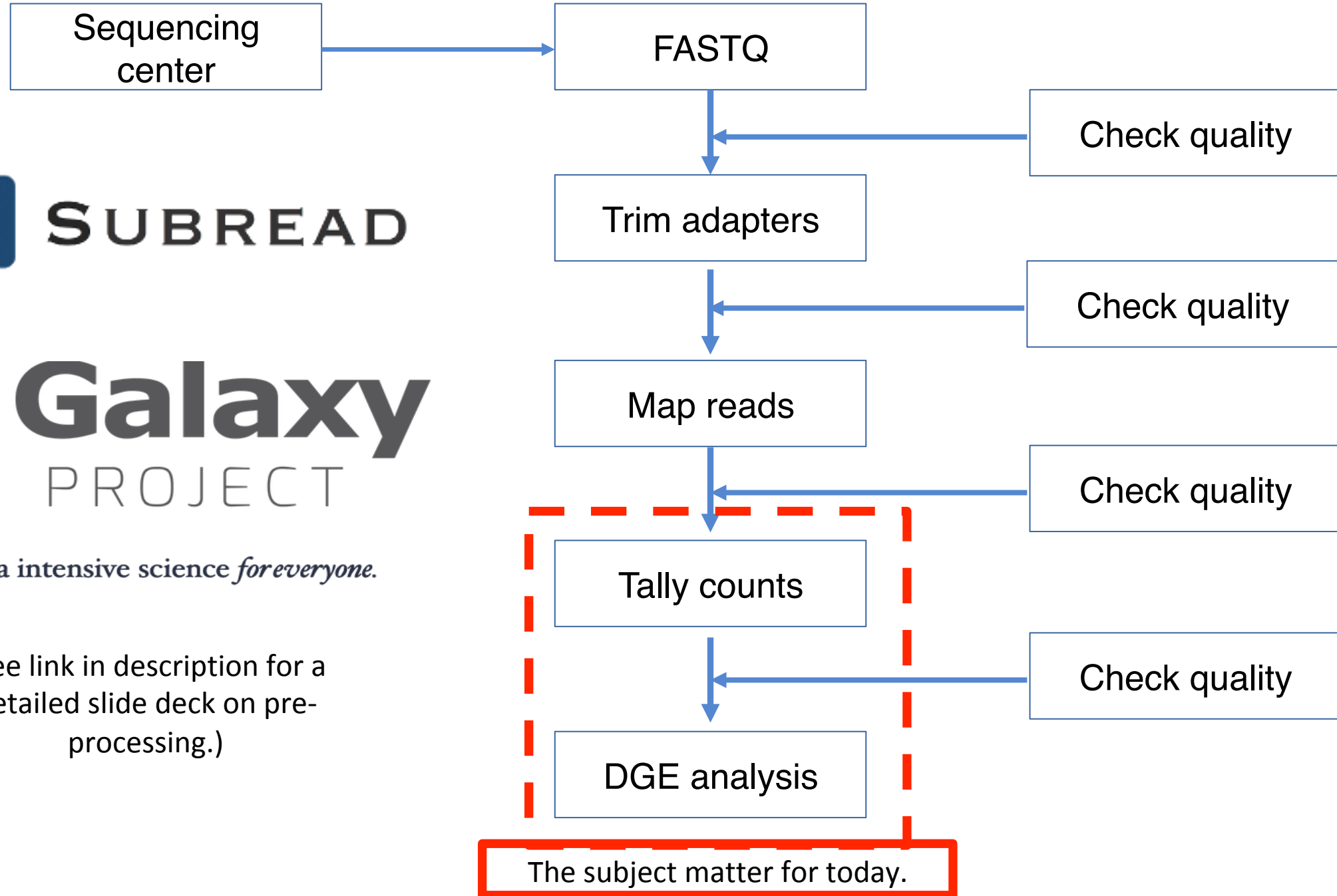


SUBREAD

 Galaxy  
PROJECT

Data intensive science *foreveryone.*

(See link in description for a  
detailed slide deck on pre-  
processing.)



# Outline



- ◆ Load and reformat the data
- ◆ Exploratory visualization : MA plot
- ◆ Create DGElist object and retrieve gene symbols
- ◆ Filter genes with inadequate information
- ◆ Normalize counts : *What's under the hood?*
- ◆ Exploratory visualization : MDS and PCA plots
- ◆ Define and fit a model
- ◆ Hypothesis testing (four example hypotheses)
- ◆ Save results as a table and explore in Excel

# Reference for the workshop



## SOFTWARE TOOL ARTICLE

# **REVISED** From reads to genes to pathways: differential expression analysis of RNA-Seq experiments using Rsubread and the edgeR quasi-likelihood pipeline [version 2; referees: 5 approved]

Yunshun Chen<sup>1,2</sup>, Aaron T. L. Lun <sup>3</sup>, Gordon K. Smyth <sup>1,4</sup>

<sup>1</sup>The Walter and Eliza Hall Institute of Medical Research, Parkville, Victoria, 3052, Australia

<sup>2</sup>Department of Medical Biology, The University of Melbourne, Victoria, 3010, Australia

<sup>3</sup>Cancer Research UK Cambridge Institute, University of Cambridge, Li Ka Shing Centre, Cambridge, UK

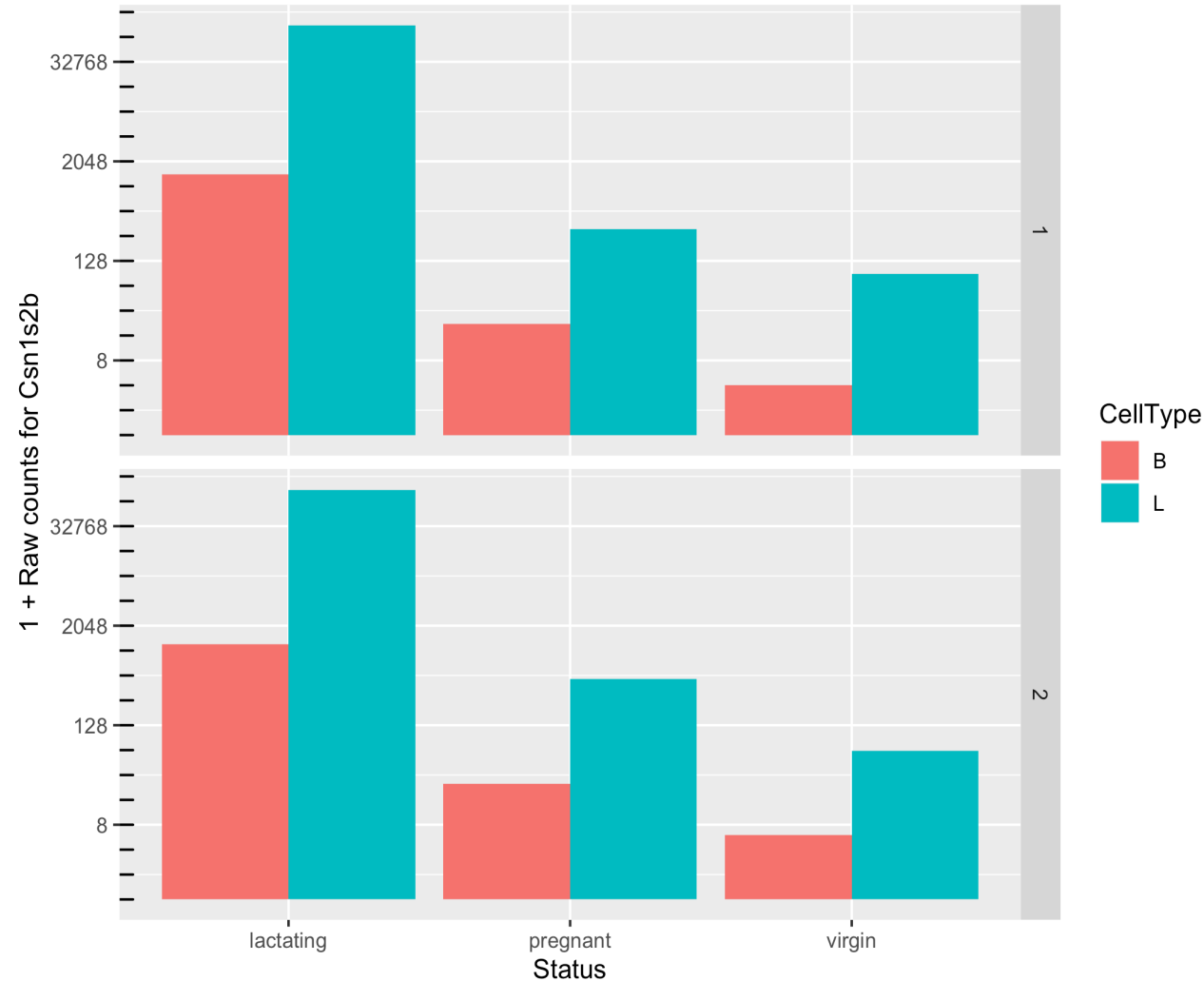
<sup>4</sup>Department of Mathematics and Statistics, The University of Melbourne, Victoria, 3010, Australia

# Dataset

- ◆ GEO accession: GSE60450
- ◆ Tissue of origin: Mammary glands of mouse
- ◆ Cell types: Basal stem-cell enriched cells (B) and committed luminal cells (L)
- ◆ Biological conditions: Virgin, Lactating and Pregnant
- ◆ # of groups: 2 cell types x 3 conditions = 6 groups
- ◆ # of replicates: 2 of each group



# Goal: To identify a set of genes that are differentially expressed

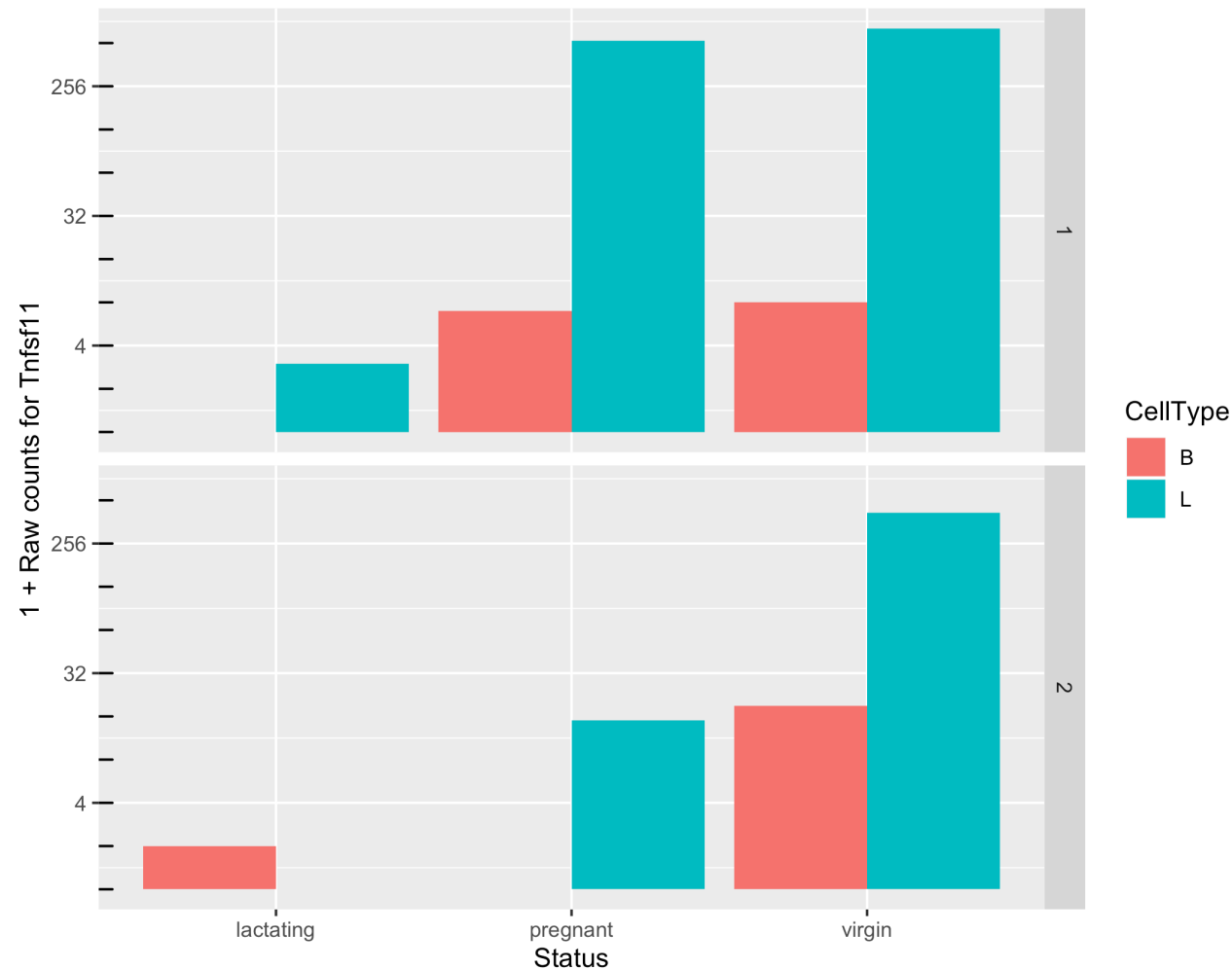


Which comparisons are we interested in?

Example:

1. B vs L,
2. B.lactating vs L.pregnant,
3. ...
4. All of them

# Goal: To identify a set of genes that are differentially expressed



How do we ensure high power for detection and high specificity when faced with noise?

Can we make reliable inferences for genes with very low counts? What should we consider “very low”?

# Factors other than differential gene expression that cause variation in counts

- ♦ Variation in sequencing depths => Need to normalize counts

Group	Total counts
B.virgin	23085177
B.virgin	21628857
B.pregnant	23919152
B.pregnant	22490570
B.lactating	21382233
B.lactating	19884434

Group	Total counts
L.virgin	20213223
L.virgin	21509988
L.pregnant	22073815
L.pregnant	21837341
L.lactating	24638939
L.lactating	24581591

# Observed counts depend on total reads sequenced and sample composition

- ◆ # reads for YFG =  $\text{Amount of nucleic acid from YFG} / \text{Total nucleic acid in sample} \times \text{Total reads}$
- ◆ Need to normalize for difference in total reads between samples.
  - ◆ Might be enough if total nucleic acid is the same in both samples.
  - ◆ Example: technical replicates
- ◆ Need to account for difference in sample composition.
  - ◆ Assume that the large majority of genes are not differential.
  - ◆ Adjust counts such that for most genes, counts are not differential.

# Normalization: Trimmed Mean of M-values

1. Choose a reference sample.
2. Compute the M and A values for all genes.
3. Filter genes that fall in the tails of M and A distributions.
4. Estimate variance of M values.
5. Estimate TMM --- the weighted average of trimmed M-values.
6. Size factor is  $2^{\uparrow TMM}$ .
7. Adjust such that these multiply to 1.

# Other approaches to normalization

- ◆ RLE approach by Anders and Huber (2010)
  - ◆ Reference: geometric mean of all samples
  - ◆ Normalization factor: median ratio of each sample to the reference
  - ◆ Identical to TMM approach
  - ◆ See link in description for complete reference
- ◆ Upper quartile normalization by Bullard et al (2010)
  - ◆ Normalization factor: 75% quantile of the counts for each sample
  - ◆ Not recommended in general
  - ◆ See link in description for complete reference

# Empirical Bayes estimates of dispersion parameters: Learning from the experience of others

- ✦ For an intuitive explanation by Bradley Efron, see link in description
- ✦ Original paper: Robbins, Herbert. *An Empirical Bayes Approach to Statistics*. (see description for complete reference)

# Your feedback is important to us!

- ◆ <https://bioinformatics-course-feedback.questionpro.com/>
- ◆ ~5 min.



Thank you.