

# Session IV. RNAseq

**Armando Reyes-Palomares**

Department Biochemistry and Molecular Biology

Complutense University of Madrid

[armandorp@ucm.es](mailto:armandorp@ucm.es)

# Practical Sessions

1. Introduction to R
2. Bioconductor: BiomaRt and GenomicRanges
3. Whole-Genome Analysis – SARS-CoV2
4. Transcriptome analysis (e.g. RNAseq)

# Table of contents

1. Survey (educational innovation project).
2. Nextclade
3. Gene expression resources
4. Align and quantification of raw data from RNA-seq assays
5. Differential Expression Analysis using DESeq2

# Parte 1<sup>a</sup>. Fase PREVIA a las sesiones Análisis Genómico (GAG)

Responde considerando tus conocimientos y estatus previo a estas sesiones de Análisis Genómico en GAG en el Máster BBMBiomed. Encuesta anónima relacionada con diversos métodos asociados al proyecto de innovación docente relacionado con el área de bioinformática de procesado y el análisis de datos en biología molecular y bioquímica.



armreyes@ucm.es (not shared) [Switch account](#)



¿Tenías conocimientos en algún lenguaje de programación el curso académico anterior?

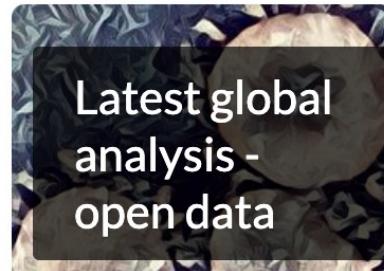
Si

No

# Nextstrain

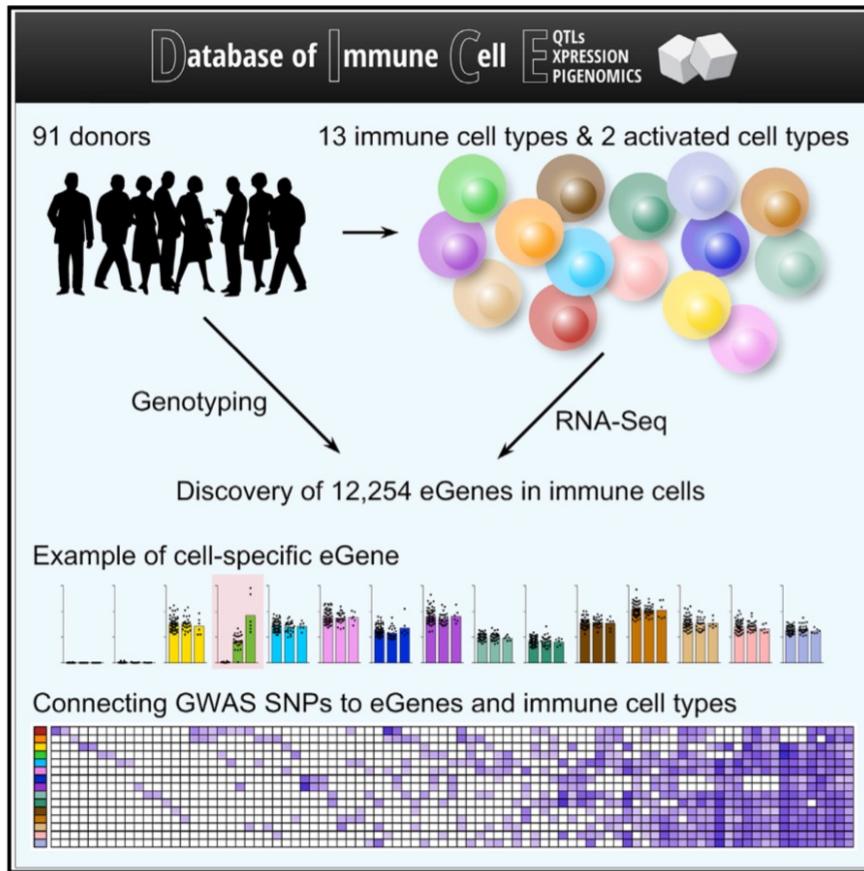
Real-time tracking of pathogen evolution  
SARS-CoV-2 (COVID-19)

We are incorporating SARS-CoV-2 genomes as soon as they are shared and providing analyses and situation reports. In addition we have developed a number of resources and tools, and are facilitating independent groups to run their own analyses.



# Impact of Genetic Polymorphisms on Human Immune Cell Gene Expression

## Graphical Abstract



## Authors

Benjamin J. Schmiedel, Divya Singh, Ariel Madrigal, ..., Mitchell Kronenberg, Bjoern Peters, Pandurangan Vijayanand

## Correspondence

vijay@lji.org

## In Brief

Surveying gene expression and SNP genotypes across immune cell types from healthy humans reveals cis-eQTLs affecting over half of all expressed genes and demonstrates that variant effects often manifest in cell types other than those with highest gene expression.



# GEO Datasets (NCBI)

GEO DataSets

GEO Profiles

This database stores original submitter-supplied study descriptions, as well as curated gene expression DataSets. DataSets form the basis of GEO's advanced data display and analysis tools, including gene expression profile charts and clusters.

## Search Examples:

Search by...	Search text
Free text	<b>smoking cancer</b>
Keywords and species	<b>(smok* OR diet) AND (mammals[organism] NOT human[organism])</b>
Studies in the <b>NIH Roadmap Epigenomics project</b>	<b>"roadmap epigenomics"[Project]</b>
Study type	<b>"expression profiling by high throughput sequencing"[DataSet Type]</b>
Studies with between 100 and 500 samples	<b>100:500[Number of Samples]</b>
Studies with CEL files	<b>"cel"[Supplementary Files]</b>
DataSets that have 'age' as an experimental variable	<b>"age"[Subset Variable Type]</b>
Author	<b>smith a[Author]</b>
Published between January and June 2007	<b>2007/01:2007/06[Publication Date]</b>
Platform accession	<b>GPL570</b>
Studies with PubMed identifiers	<b>"gds pubmed"[Filter]</b>

# ArrayExpress (EMBL-EBI)

EMBL-EBI

Services

Research

Training

About us

EMBL-EBI  Hinxton



## ArrayExpress

Search



Examples: [E-MEXP-31](#), [cancer](#), [p53](#), [Geuvadis](#)

 advanced search

Home

Browse

Submit

Help

About ArrayExpress

Contact Us

 Login

## ArrayExpress – functional genomics data

ArrayExpress Archive of Functional Genomics Data stores data from high-throughput functional genomics experiments, and provides these data for reuse to the research community.



Browse ArrayExpress

### Data Content

Updated today at 02:00

- 74979 experiments
- 2569020 assays
- 61.29 TB of archived data



UNIVERSIDAD  
COMPLUTENSE  
MADRID

armandorp@ucm.es

# GTEx portal (gtexportal.org)

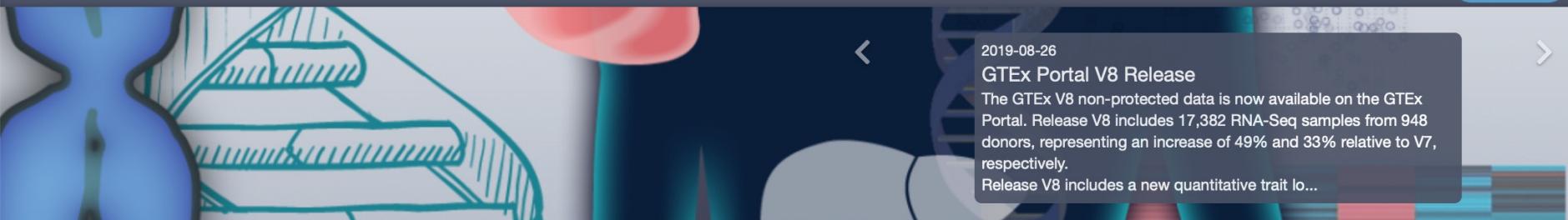


About GTEx Publications Access Biospecimens FAQs Contact

Home Datasets Expression QTLs & Browsers Sample Data Documentation

Search Gene or SNP ID...

Sign In



2019-08-26

## GTEx Portal V8 Release

The GTEx V8 non-protected data is now available on the GTEx Portal. Release V8 includes 17,382 RNA-Seq samples from 948 donors, representing an increase of 49% and 33% relative to V7, respectively.

Release V8 includes a new quantitative trait lo...

## Resource Overview

### Current Release (V8)

- Tissue & Sample Statistics
- Tissue Sampling Info (Anatomogram)
- Access & Download Data
- Release History
- How to cite GTEx?

The Genotype-Tissue Expression (GTEx) project is an ongoing effort to build a comprehensive public resource to study tissue-specific gene expression and regulation. Samples were collected from 54 non-diseased tissue sites across nearly 1000 individuals, primarily for molecular assays including WGS, WES, and RNA-Seq. Remaining samples are available from the GTEx Biobank. The GTEx Portal provides open access to data including gene expression, QTLs, and histology images.

## Explore GTEx

### Browse



By gene ID

Browse and search all data by gene



By variant or rs ID

Browse and search all data by variant



By Tissue

Browse and search all data by tissue



Histology Image Viewer

Browse and search GTEx histology images

### Expression



Multi-Gene Query

Browse and search expression by gene and tissue



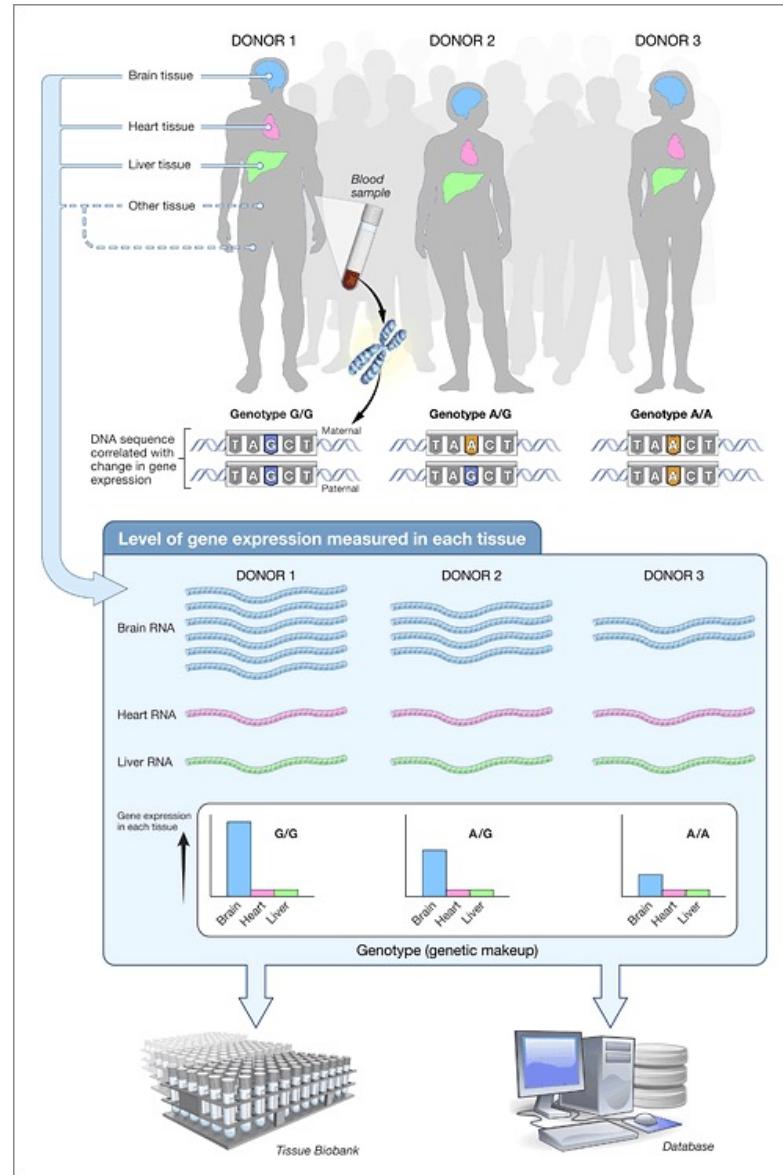
Top 50 Expressed Genes

Visualize the top 50 expressed genes in each tissue



Transcript Browser

Visualize transcript expression and isoform structures



# Expression Atlas (EMBL-EBI)

[EMBL-EBI](#)[Services](#)[Research](#)[Training](#)[About us](#)

EMBL-EBI



## Expression Atlas

Gene expression across species and biological conditions

Query single cell expression

[To Single Cell Expression Atlas](#)

[Home](#)[Browse experiments](#)[Download](#)[Release notes](#)[FAQ](#)[Help](#)[Licence](#)[Also in this section](#) ▾

Search across 65 species, 4,169 studies, 139,128 assays

Ensembl 99, Ensembl Genomes 46, WormBase ParaSite 14,

[Animals](#)[Plants](#)[Fungi](#)[Homo sapiens](#)

1518 experiments

Baseline: 79

Differential: 1439

[Mus musculus](#)

1185 experiments

Baseline: 49

Differential: 1136

[Rattus norvegicus](#)

152 experiments

Baseline: 3  
Differential: 149[Drosophila melanogaster](#)

142 experiments

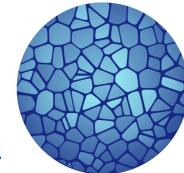
Baseline: 4  
Differential: 138[Gallus gallus](#)

39 experiments

Baseline: 4  
Differential: 35[Caenorhabditis elegans](#)30 experiments  
Baseline: 1  
Differential: 29

UNIVERSIDAD  
COMPLUTENSE  
MADRID

[armandorp@ucm.es](mailto:armandorp@ucm.es)



HUMAN  
CELL  
ATLAS

# The Human Cell Atlas

[www.humancellatlas.org](http://www.humancellatlas.org) / [data.humancellatlas.org](http://data.humancellatlas.org)

Explore Guides Metadata Pipelines Analysis Tools Contribute APIs

Update: The DCP 2.0 Data View is now available. | [Learn More](#)

## Mapping the Human Body at the Cellular Level

Community generated, multi-omic,  
open data processed by uniform pipelines



14.9M  
CELLS



76  
ORGANS



1.6k  
DONORS



151  
PROJECTS



294

### About HCA:

- From 2016, > 1,000 Institutions, 75 countries
- Catalog all cell types and sub-types
- Map cell types to their location within tissues
- Distinguish cell states
- Study transitions, such as activation or differentiation.



UNIVERSIDAD  
COMPLUTENSE  
MADRID

armandorp@ucm.es

# The Cancer Genome Atlas

*portal.gdc.cancer.gov*

NATIONAL CANCER INSTITUTE  
GDC Data Portal

Home Projects Exploration Analysis Repository

Manage Sets

Harmonized Cancer Datasets

## Genomic Data Commons Data Portal

Get Started by Exploring:

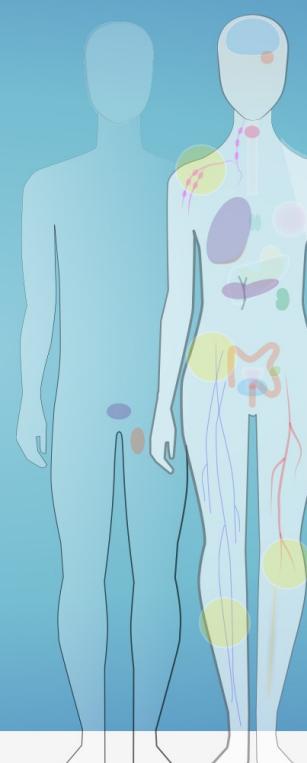
Projects Exploration Analysis Repository

e.g. BRAF, Breast, TCGA-BLCA, TCGA-A5-A0G2

Data Portal Summary Data Release 27.0 - October 29, 2020

PROJECTS	PRIMARY SITES	CASES
67	68	84,392

FILES	GENES	MUTATIONS
596,758	23,399	3,287,299

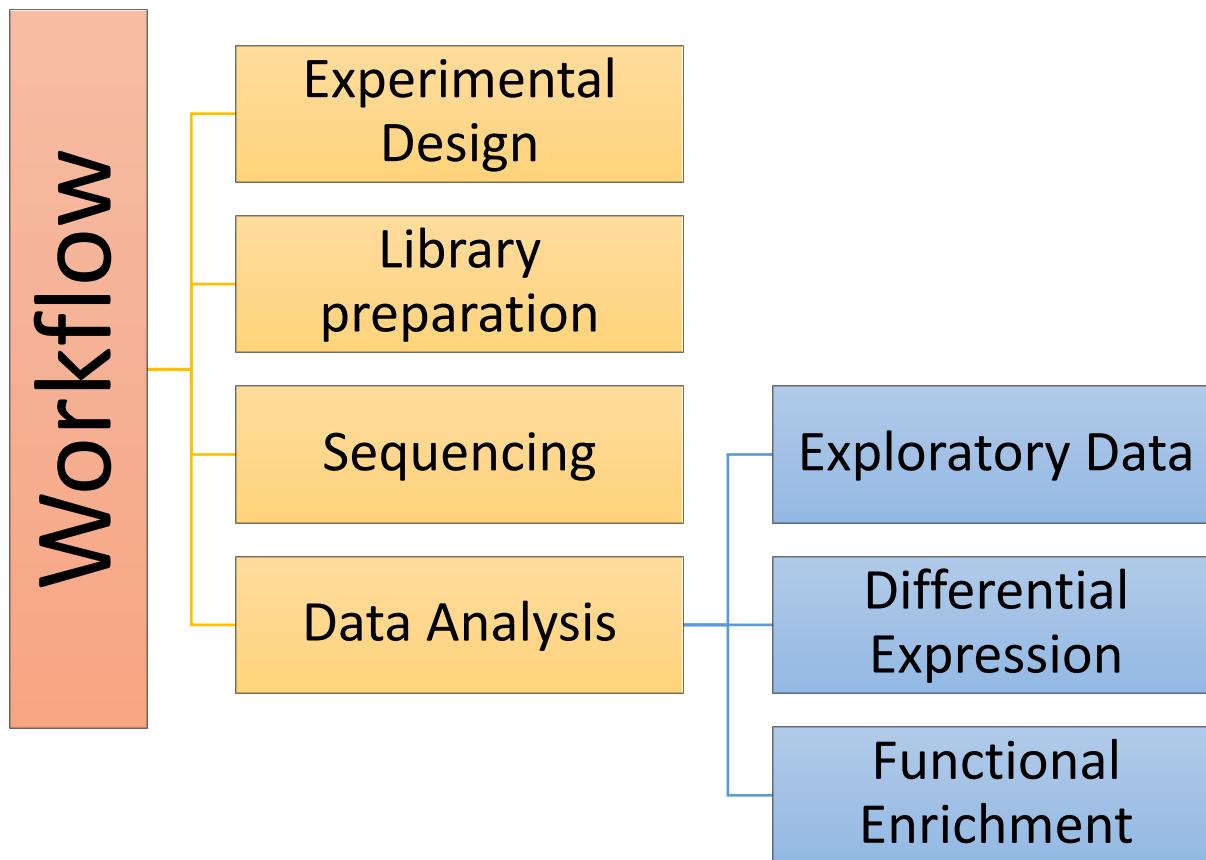


Cases by Major Primary Site

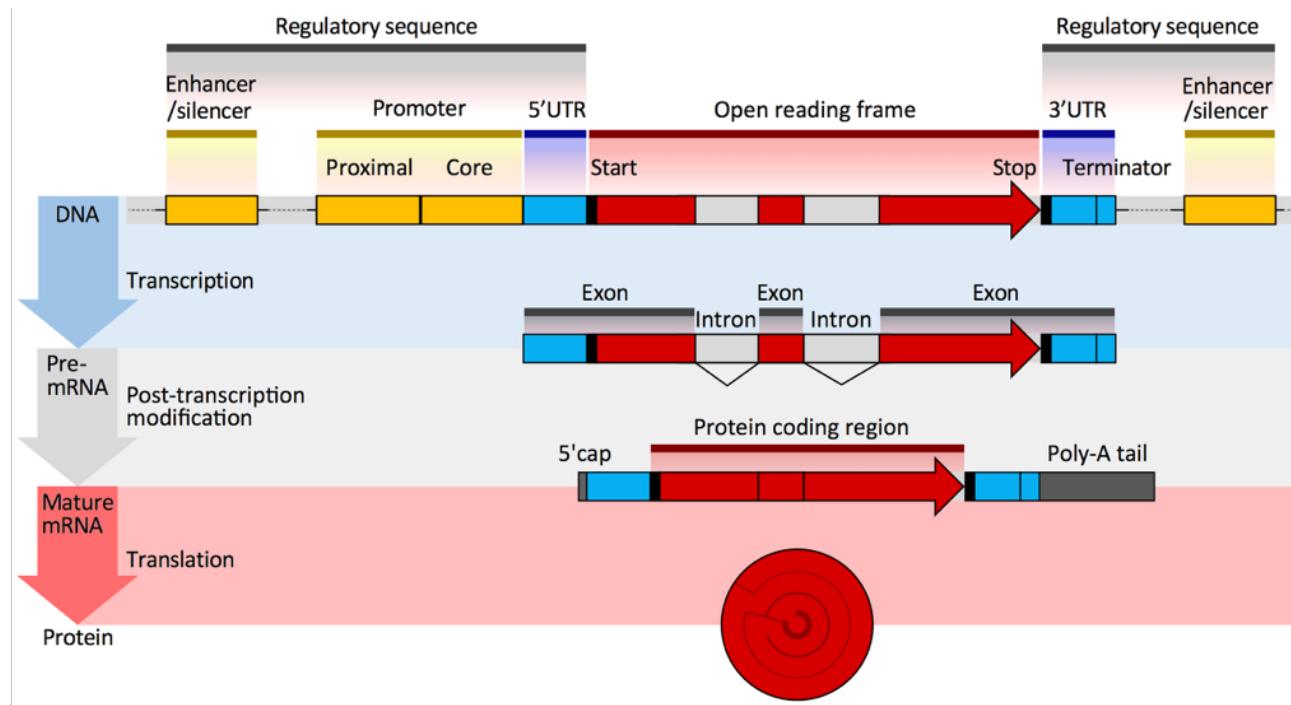
Primary Site	Cases
Adrenal Gland	1
Bile Duct	1
Bladder	2
Bone	1
Bone Marrow	9
Brain	2
Breast	9
Cervix	1
Colorectal	8
Esophagus	1
Eye	1
Head and Neck	2
Kidney	3
Liver	1
Lung	10
Lymph Nodes	1
Nervous System	3
Ovary	3
Pancreas	2
Pleura	1
Prostate	2
Skin	3
Soft Tissue	1
Stomach	1
Testis	1
Thymus	1
Thyroid	2
Uterus	3

Data Portal Website API Data Transfer Tool Documentation Data Submission Portal Legacy Archive Publications

# RNA-seq Workflow



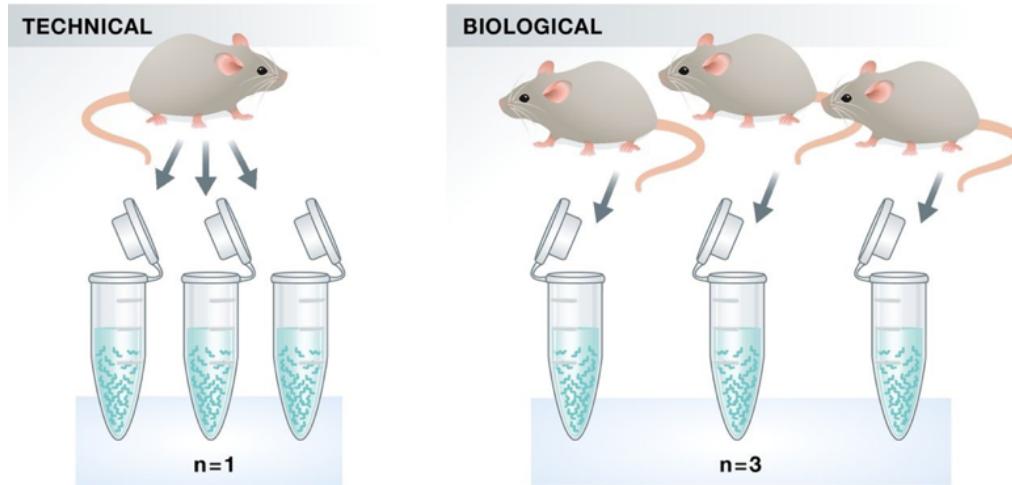
# Experimental Design: what do we want?



Wikimedia Commons [Gene structure eukaryote 2 annotated.svg](#) by Thomas Shafee, used under Creative Commons Attribution 4.0 International

While **mRNA transcripts have a polyA tail**, many of the non-coding RNA transcripts do not as the post-transcriptional processing is different for these transcripts.

# Experimental Design: number and type of replicates



*Image credit: [Klaus B., EMBO J \(2015\) 34: 2727-2730](#)*

- **Technical replicates:** use the same biological sample to repeat the technical or experimental steps in order to accurately measure technical variation and remove it during analysis (**Microarray**)
- **Biological replicates** use different biological samples of the same condition to measure the biological variation between samples.

# Experimental Design: number and type of replicates

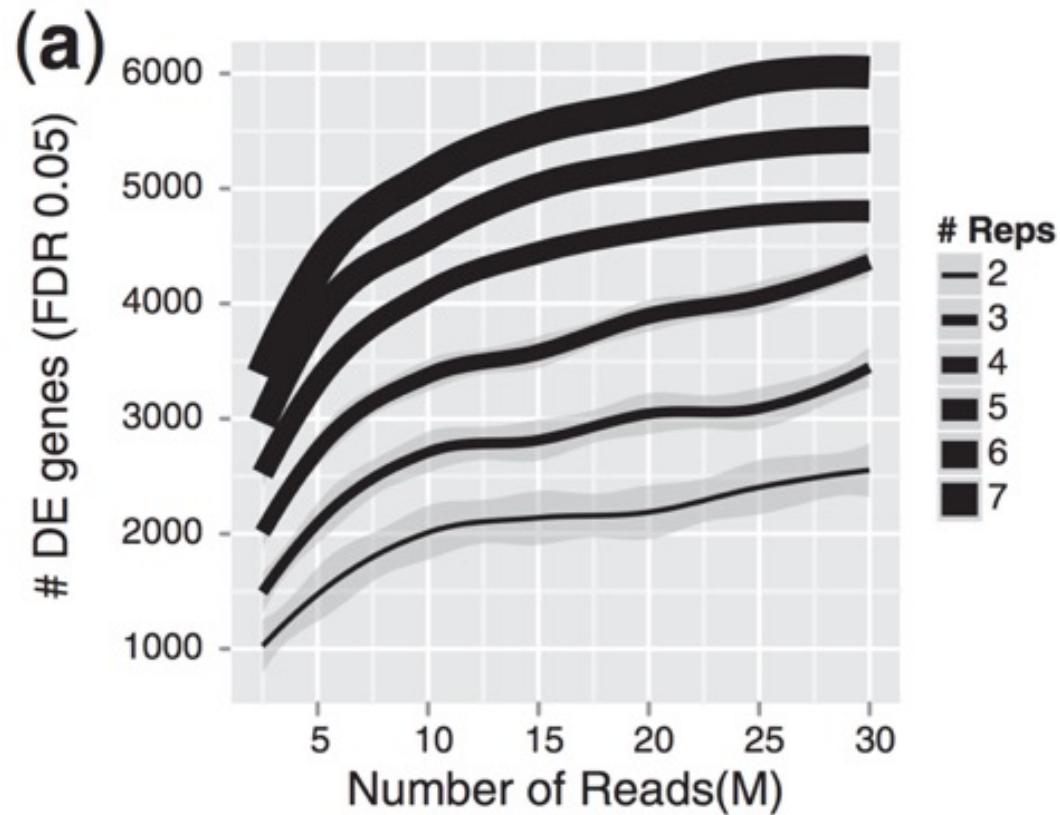
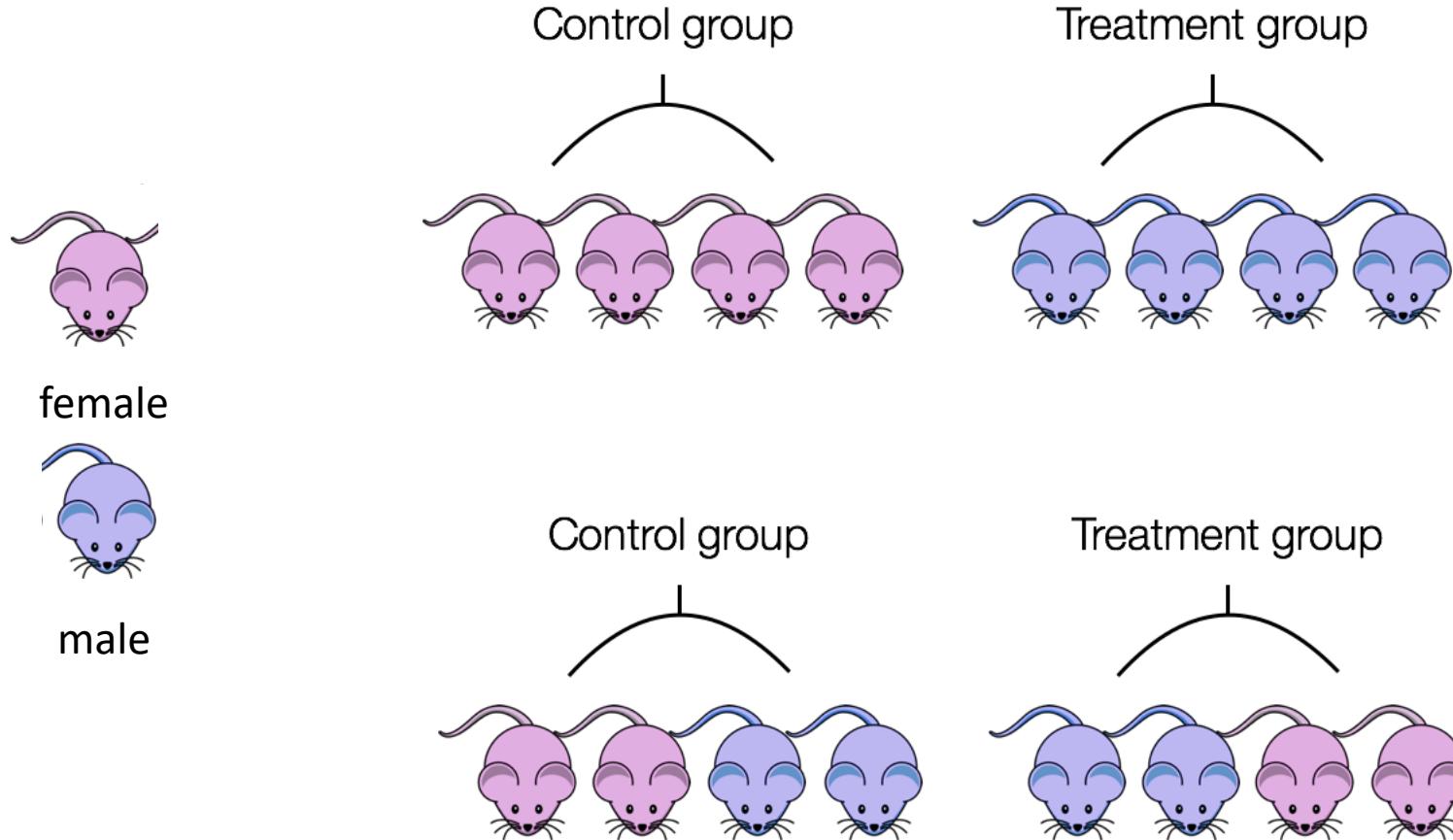
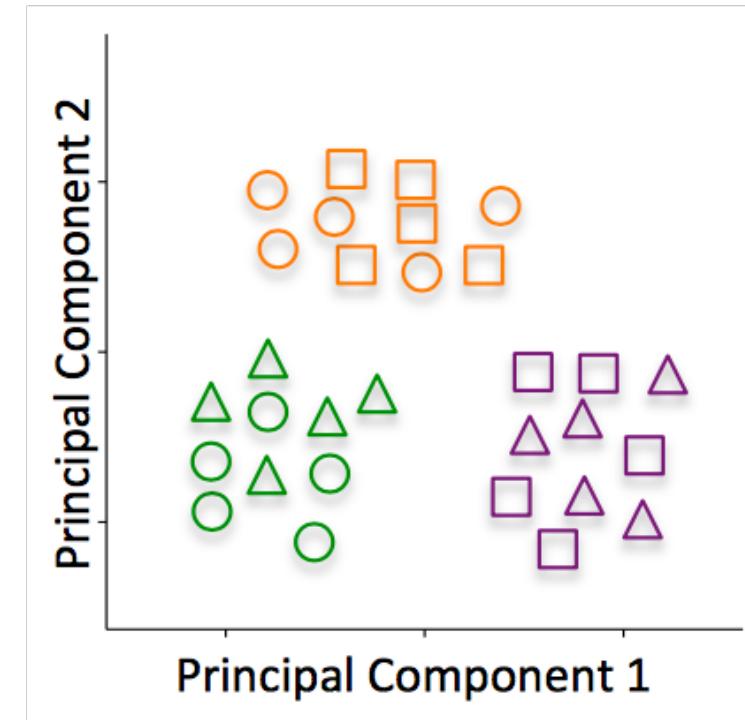
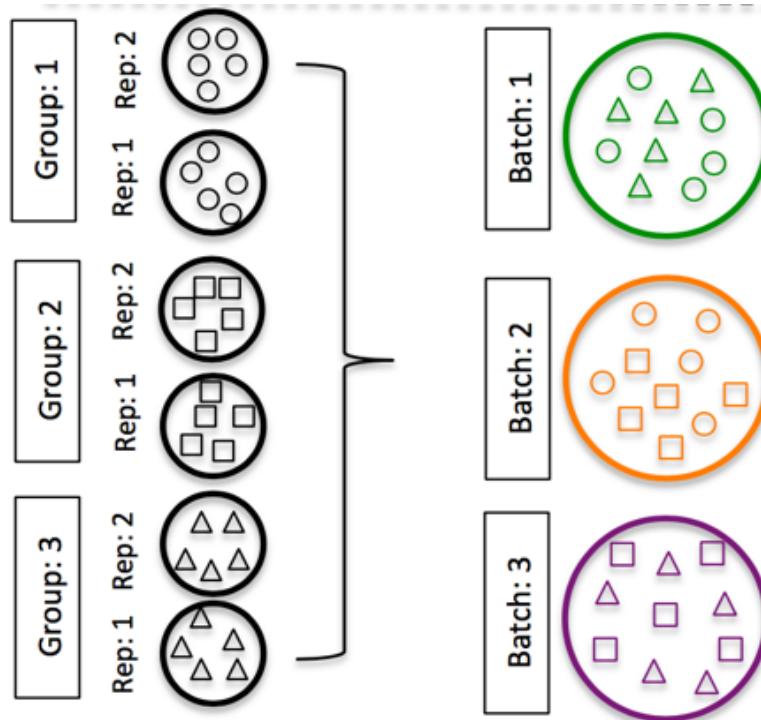


Image credit: [Liu, Y., et al., Bioinformatics \(2014\) 30\(3\): 301–304](#)

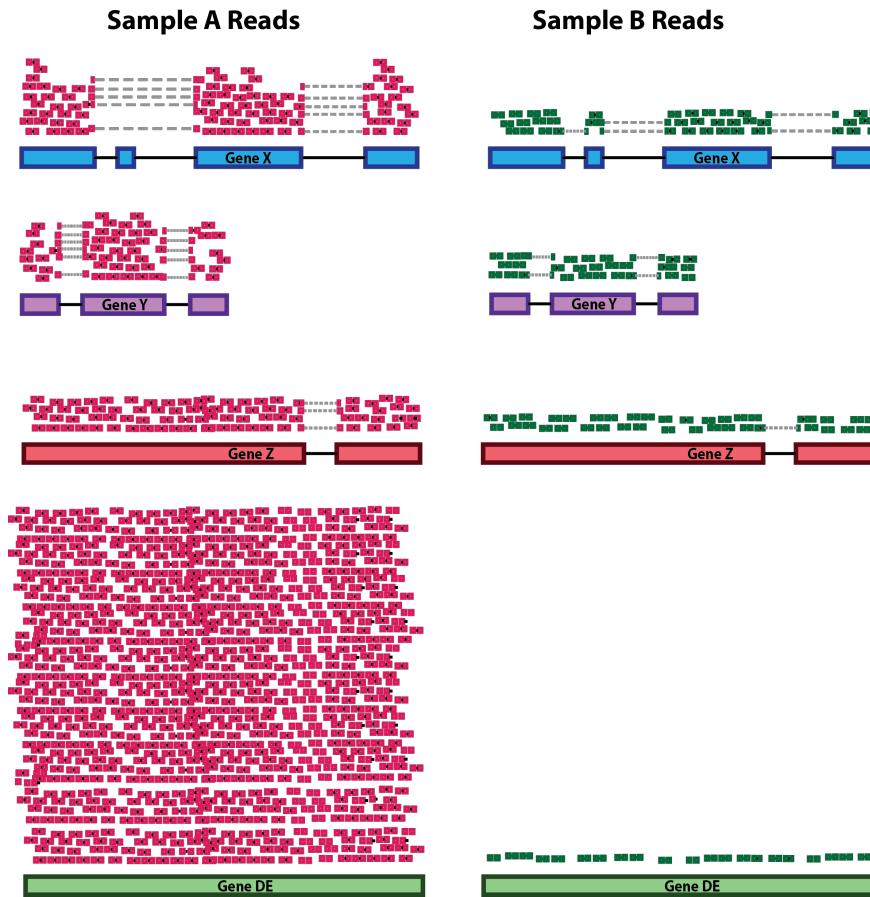
# Experimental Design: confounding factors



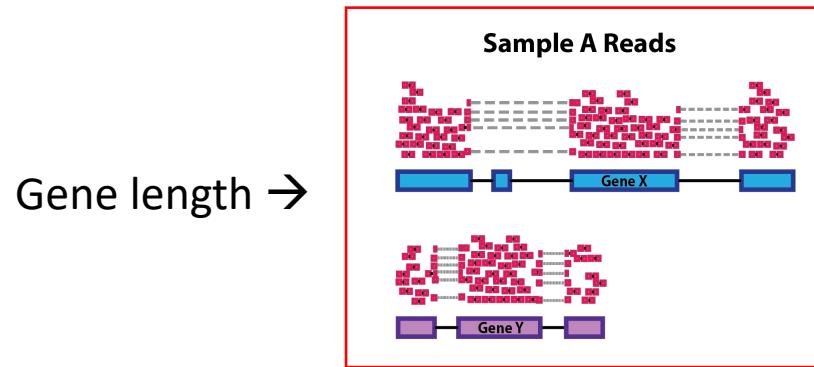
# Experimental Design: Batch effects



# Technical variation (Library preparation) Normalization factors

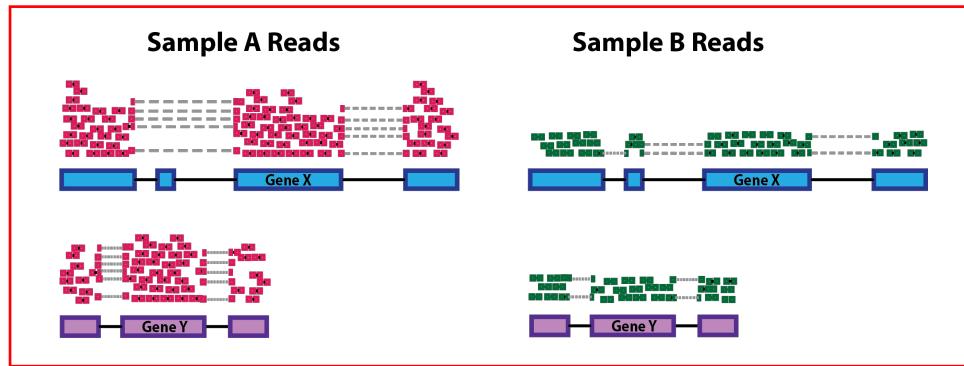


# Technical variation (Library preparation) Normalization factors

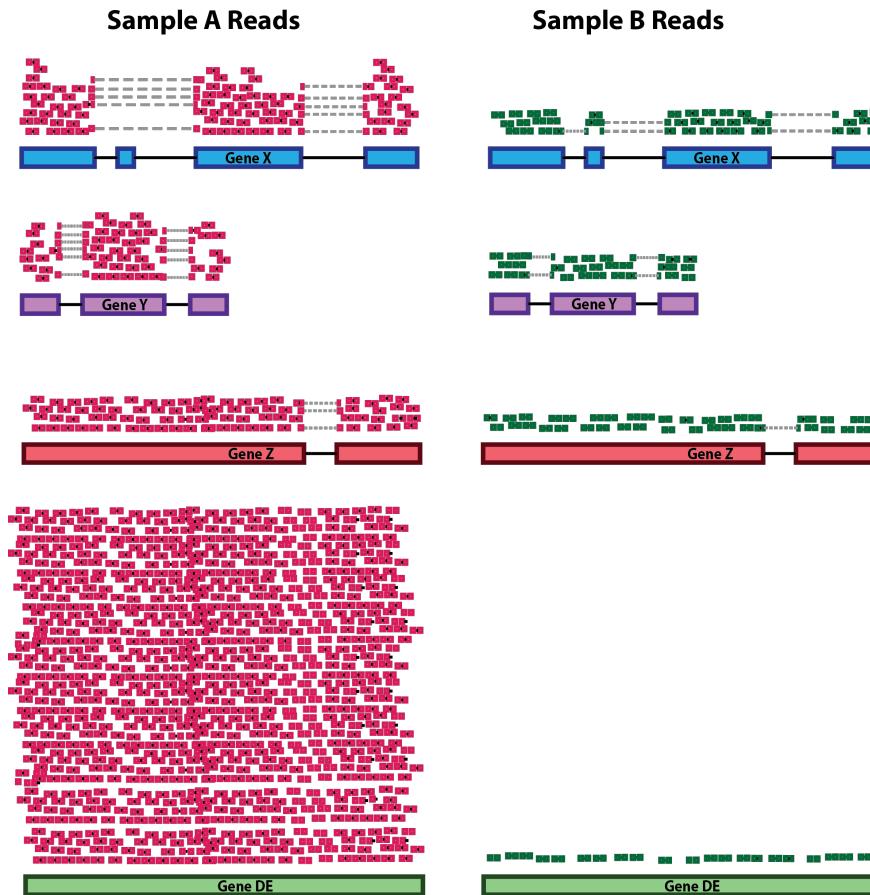


# Technical variation (Library preparation) Normalization factors

Sequence depth  
(Library size) →



# Technical variation (Library preparation) Normalization factors



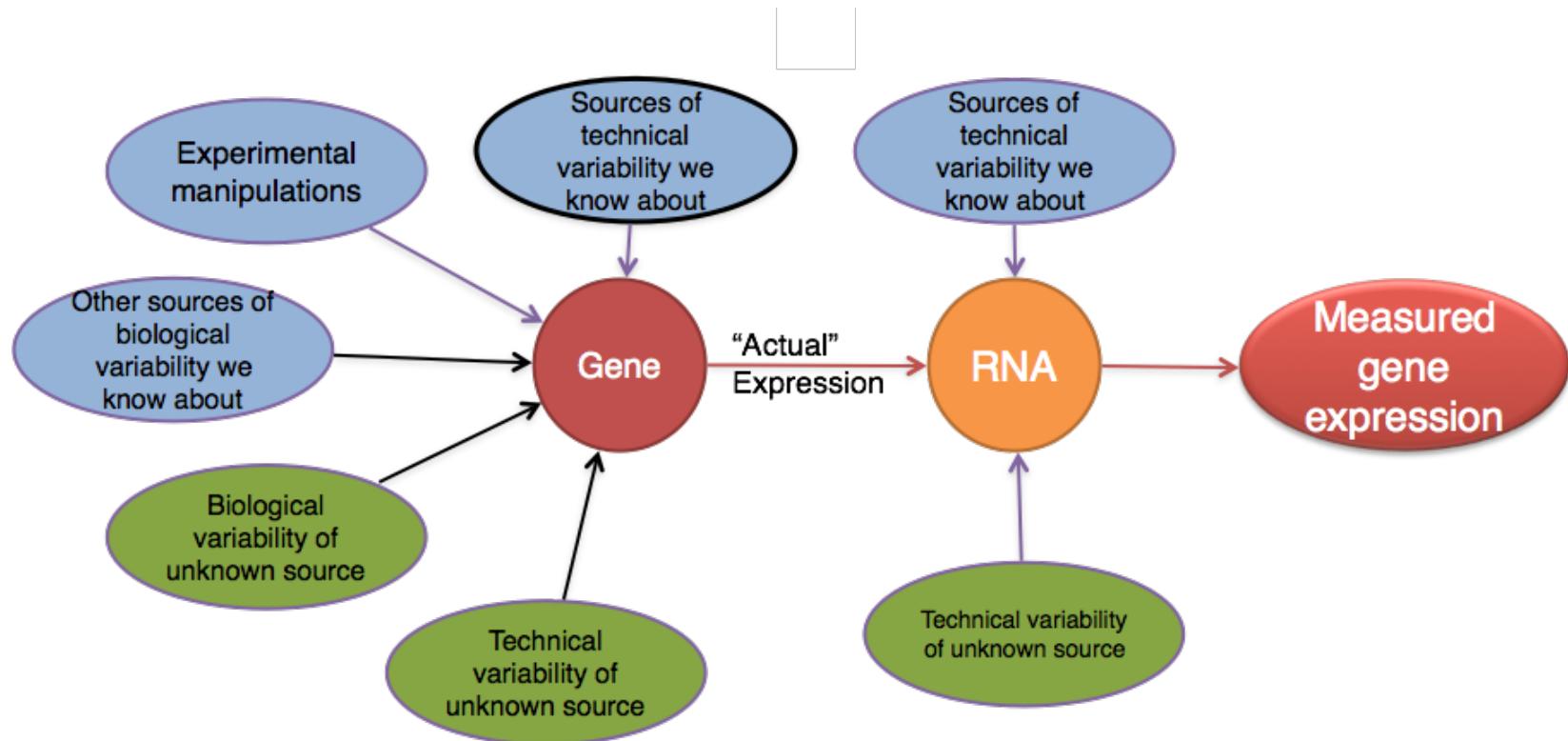
Minimizing the  
effect of highly  
expressed genes →



# Normalization factors in RNAseq data

Normalization method	Description	Accounted factors	Recommendations for use
CPM (counts per million)	counts scaled by total number of reads	sequencing depth	gene count comparisons between replicates of the same samplegroup; <b>NOT for within sample comparisons or DE analysis</b>
TPM (transcripts per kilobase million)	counts per length of transcript (kb) per million reads mapped	sequencing depth and gene length	gene count comparisons within a sample or between samples of the same sample group; <b>NOT for DE analysis</b>
RPKM/FPKM (reads/fragments per kilobase of exon per million reads/fragments mapped)	similar to TPM	sequencing depth and gene length	gene count comparisons between genes within a sample; <b>NOT for between sample comparisons or DE analysis</b>
DESeq2's median of ratios [1]	counts divided by sample-specific size factors determined by median ratio of gene counts relative to geometric mean per gene	sequencing depth and RNA composition	gene count comparisons between samples and for <b>DE analysis</b> ; <b>NOT for within sample comparisons</b>
EdgeR's trimmed mean of M values (TMM) [2]	uses a weighted trimmed mean of the log expression ratios between samples	sequencing depth, RNA composition, and gene length	gene count comparisons between and within samples and for <b>DE analysis</b>

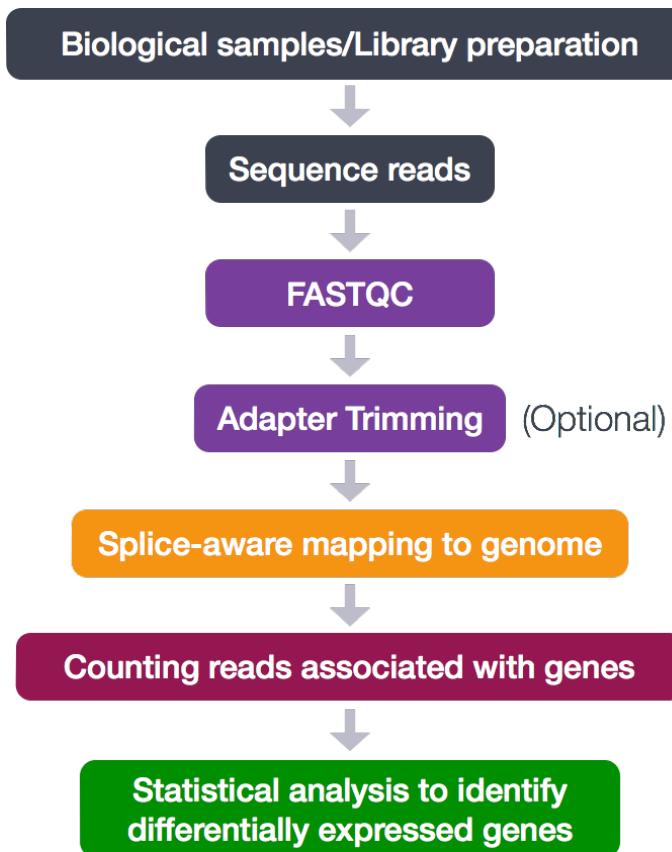
# Sources of variation in assays



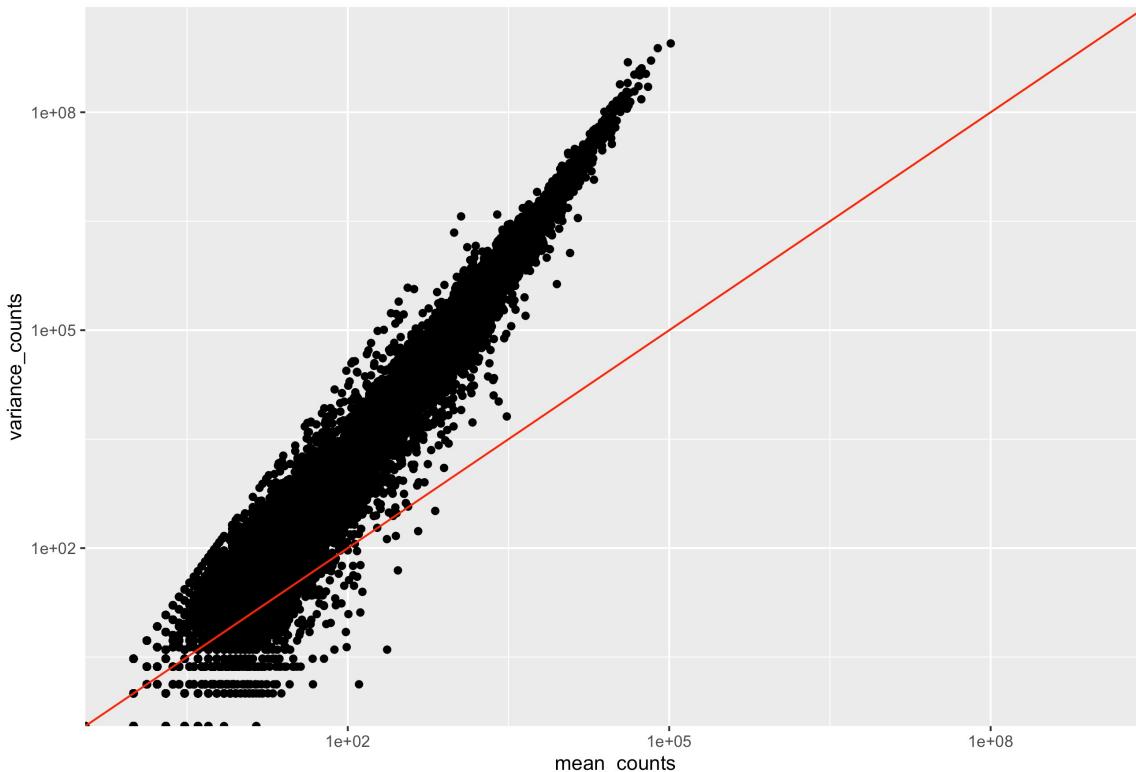
Courtesy of Paul Pavlidis, UBC



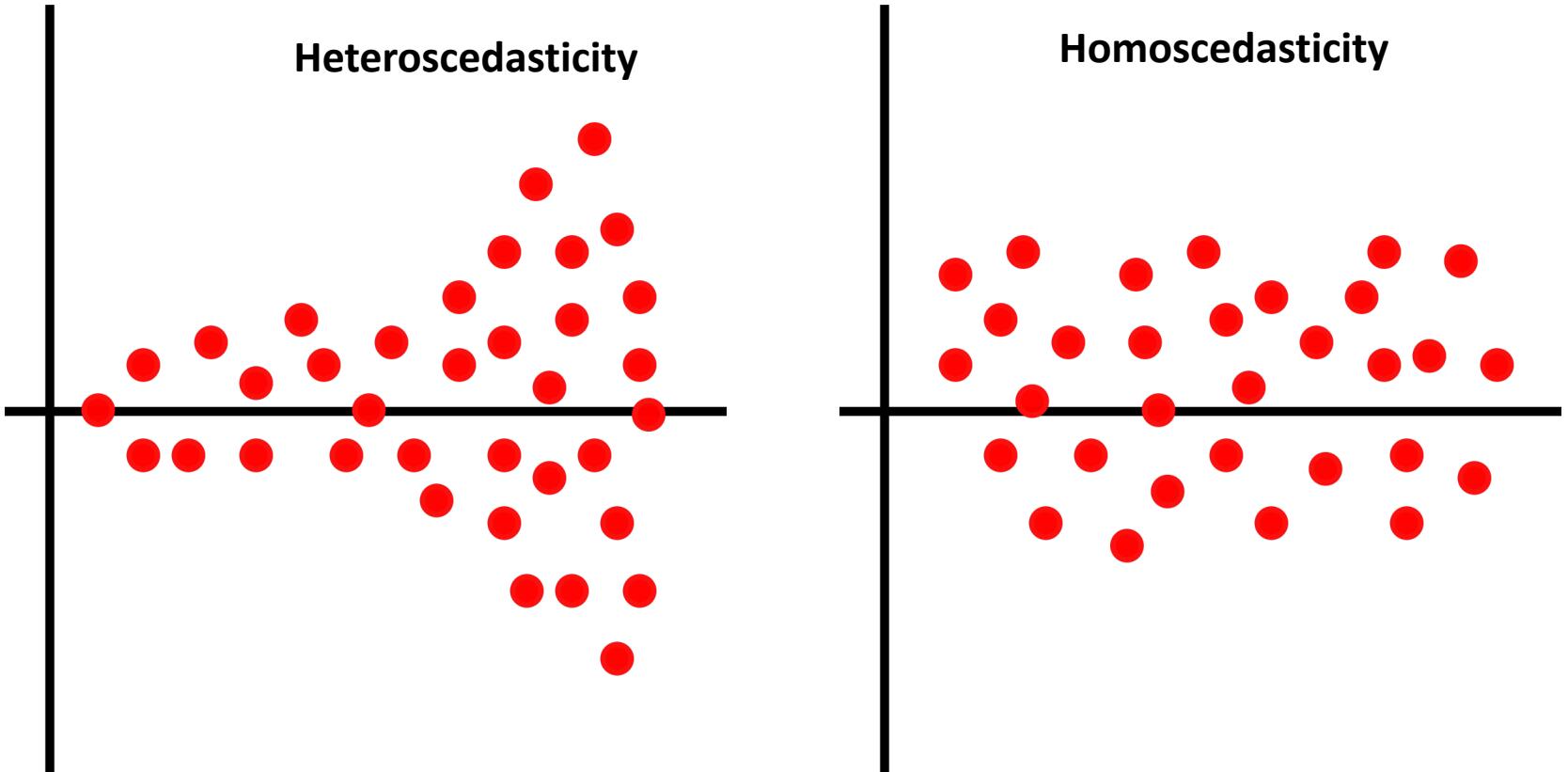
# Raw RNAseq data and quantifications



# Variance increases with mean expression

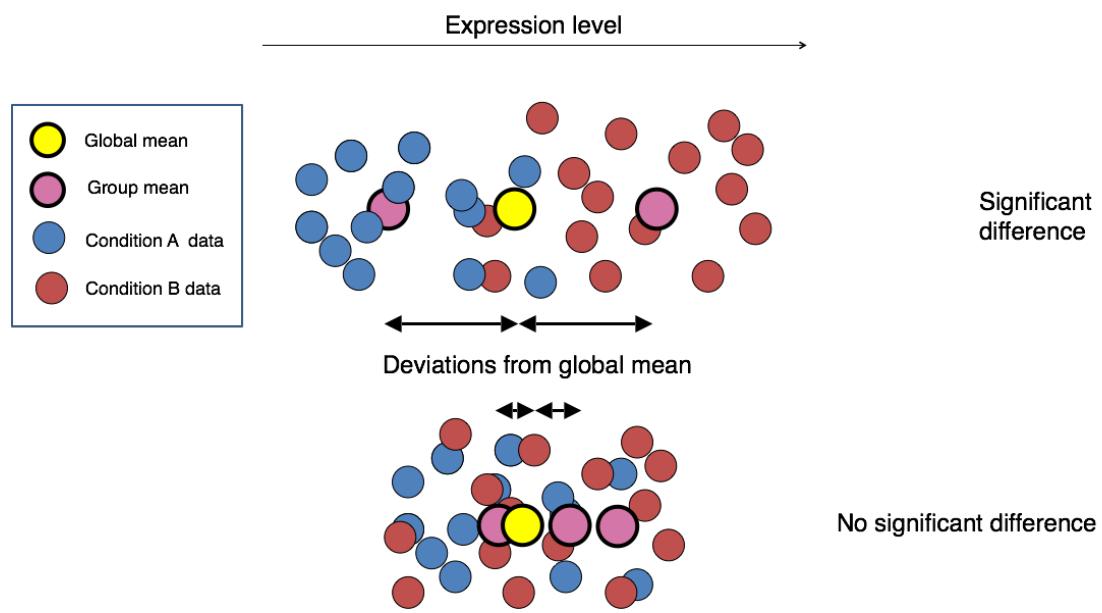


# Raw counts in RNAseq



# Differential Expression Analysis (DEA) steps with DESeq2

- 1 Estimate size factors
- 2 Estimate gene-wise dispersions
- 3 Fit curve to gene-wise dispersion estimates
- 4 Shrink gene-wise dispersion estimates
- 5 GLM fit for each gene



# DESeq2: (1) Size Factors

$$\sqrt{1489 \cdot 906} = 1161.5$$

Estimate size factors

↓  
Estimate gene-wise dispersion

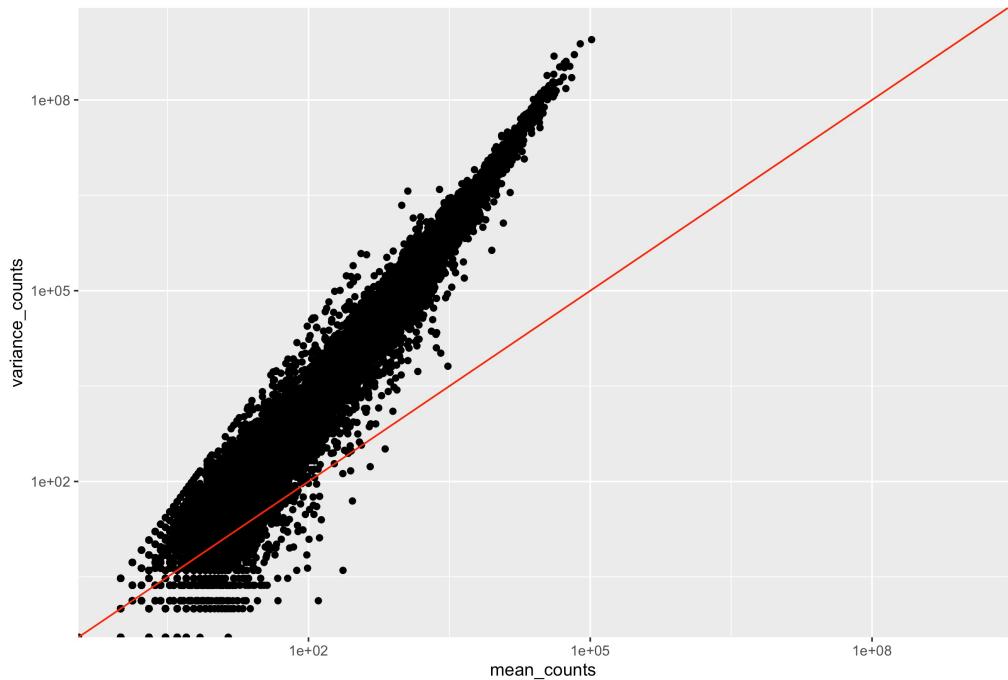
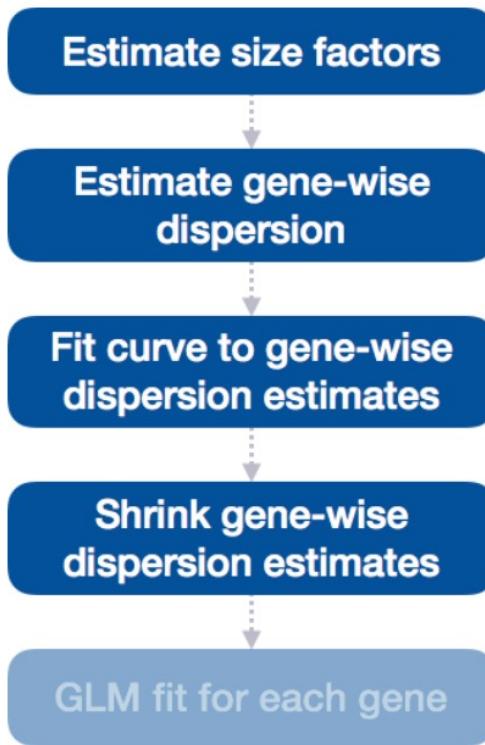
↓  
Fit curve to gene-wise dispersion estimates

↓  
Shrink gene-wise dispersion estimates

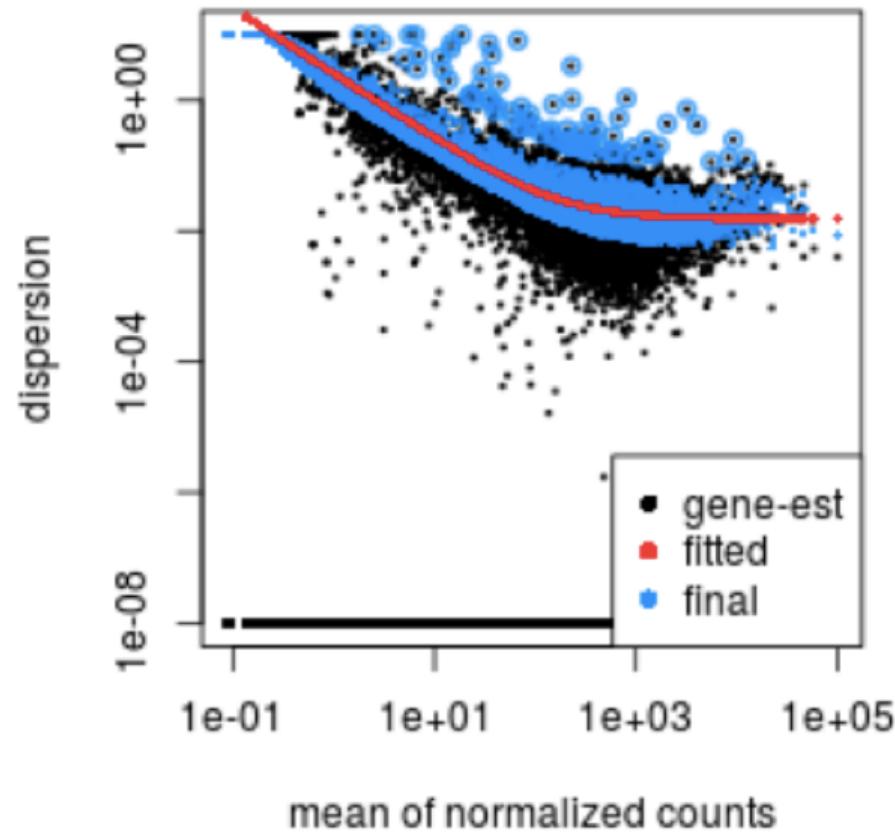
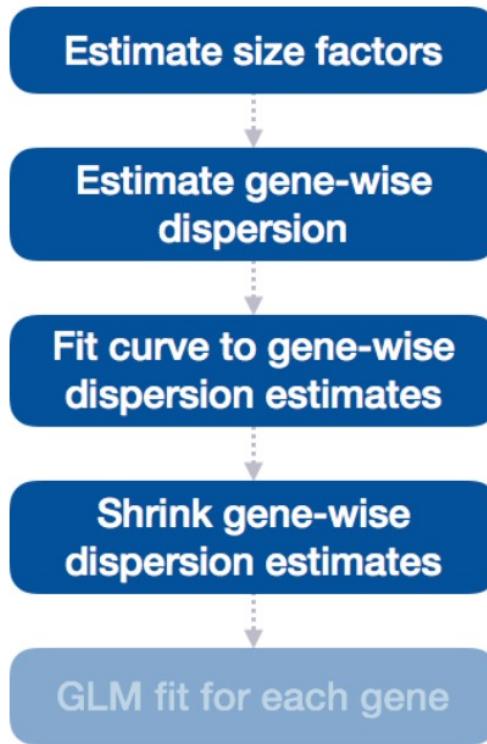
↓  
GLM fit for each gene

gene	sampleA	sampleB	pseudo-reference sample	ratio of sampleA/ref	ratio of sampleB/ref
EF2A	1489	906	1161.5	1489/1161.5 = <b>1.28</b>	906/1161.5 = <b>0.78</b>
ABCD1	22	13	16.9	22/16.9 = <b>1.30</b>	13/16.9 = <b>0.77</b>
MEFV	793	410	570.2	793/570.2 = <b>1.39</b>	410/570.2 = <b>0.72</b>
BAG1	76	42	56.5	76/56.5 = <b>1.35</b>	42/56.5 = <b>0.74</b>
MOV10	521	1196	883.7	521/883.7 = <b>0.590</b>	1196/883.7 = <b>1.35</b>
...	...	...	...		

# DESeq2: (2-4) Gene wise dispersions



# DESeq2: (2-4) Gene wise dispersions



# DESeq2: (5) Negative Binomial Model

raw count for gene i, sample j

The mean is taken as “normalized counts” scaled by a normalization factor

one dispersion per gene

$$K_{ij} \sim \text{NB}(s_{ij}q_{ij}, \alpha_i)$$

**The NB model is a good approximation for data where the mean < variance, as is the case with RNA-Seq count data.**



# DESeq2: Negative Binomial Model

gene\_id CAF0006876

	sample1	sample2	sample3	sample4	sample5	sample6	sample7	sample8
Condition A	23171	22903	29227	24072	23151	26336	25252	24122
Condition B	Sample9	sample10	sample11	sample12	sample13	sample14	sample15	sample16
	19527	26898	18880	24237	26640	22315	20952	25629

