# Introduction to RNA-Seq

Introduction To Bioinformatics Using NGS Data

**Roy Francis** | 25-Oct-2018
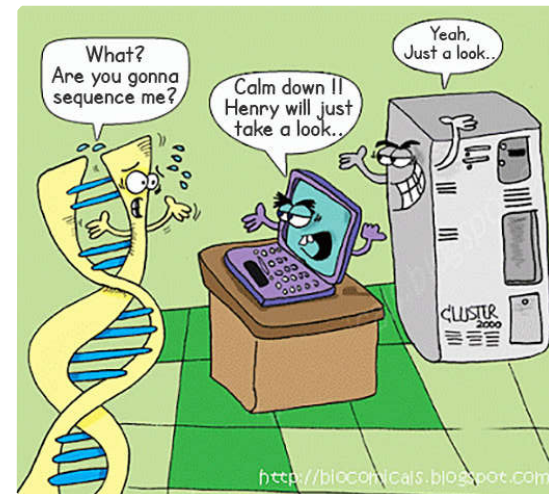
# Contents
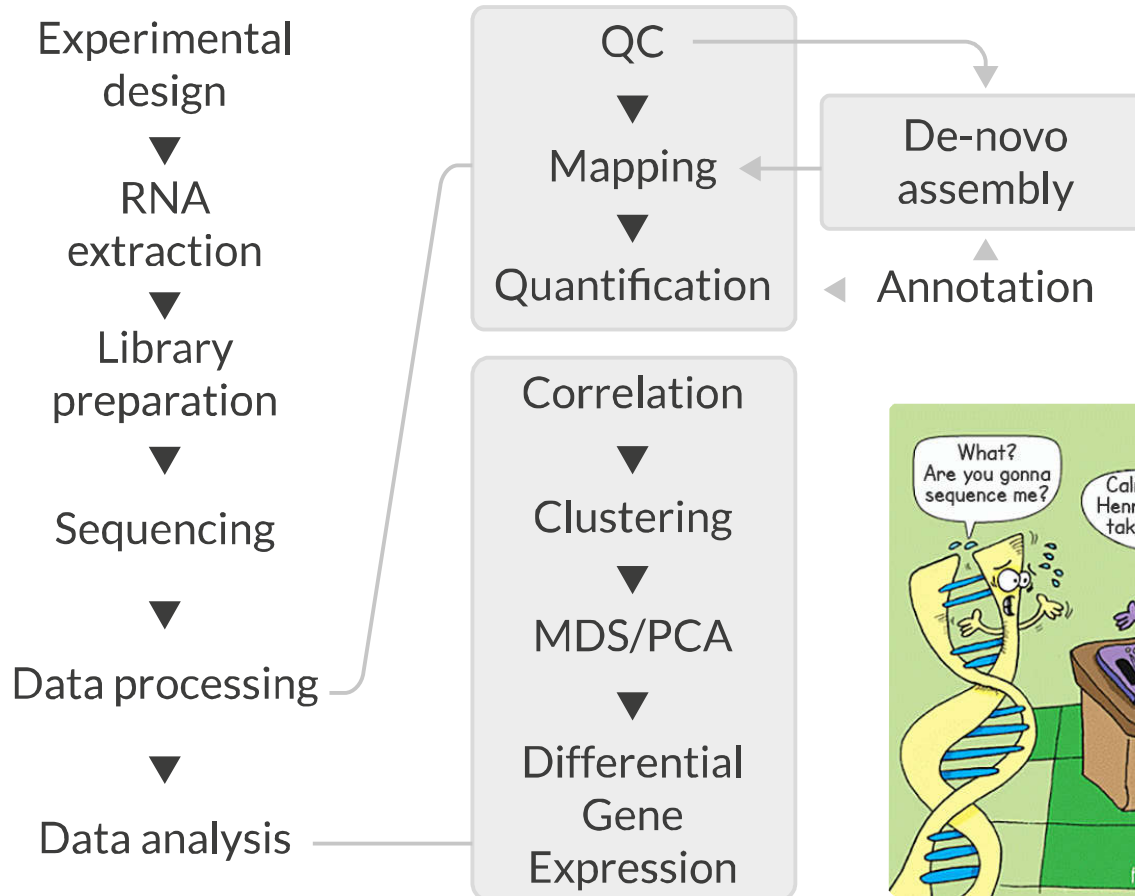
# RNA Sequencing

- The transcriptome is spatially and temporally dynamic
- Data comes from functional units (coding regions)
- Only a tiny fraction of the genome

# Applications

- Identify gene sequences in genomes
- Learn about gene function
- Differential gene expression
- Explore isoform and allelic expression
- Understand co-expression, pathways and networks
- Gene fusion
- RNA editing

# Workflow

Experimental
design
▼
RNA
extraction
▼
Library
preparation
▼
Sequencing
▼
Data processing
▼
Data analysis

QC
▼
Mapping
▼
Quantification

De-novo
assembly

Annotation

Correlation
▼
Clustering
▼
MDS/PCA
▼
Differential
Gene
Expression



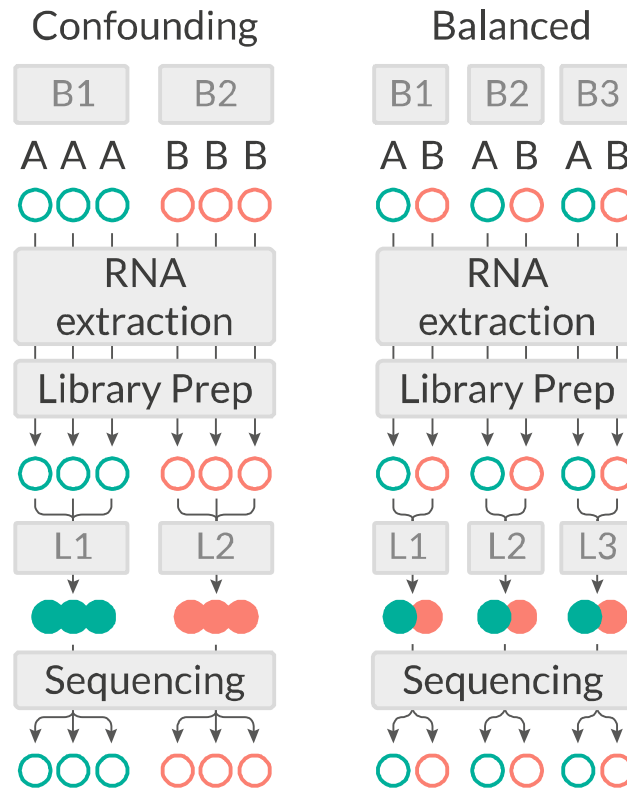🔗 Conesa, Ana, *et al.* "A survey of best practices for RNA-seq data analysis." Genome biology 17.1 (2016): 13

# Experimental design

- Balanced design
- Technical replicates not necessary (Marioni *et al.*, 2008)
- Biological replicates: 6 - 12 (Schurch *et al.*, 2016)
- ENCODE consortium
- Previous publications
- Power analysis

🧰 RnaSeqSampleSize (Power analysis), Scotty (Power analysis with cost)

🔗 Busby, Michele A., *et al.* "Scotty: a web tool for designing RNA-Seq experiments to measure differential gene expression." Bioinformatics 29.5 (2013): 656-657
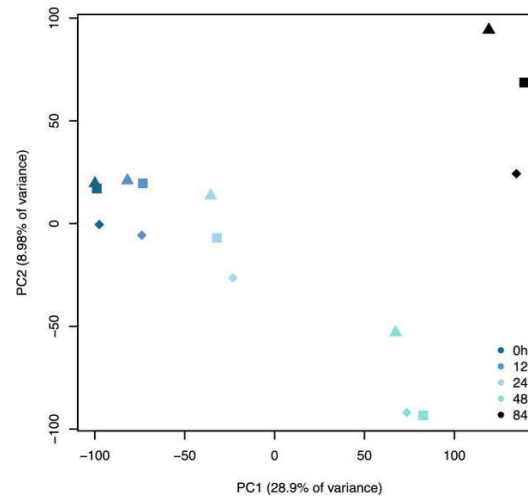🔗 Marioni, John C., *et al.* "RNA-seq: an assessment of technical reproducibility and comparison with gene expression arrays." Genome research (2008)
🔗 Schurch, Nicholas J., *et al.* "How many biological replicates are needed in an RNA-seq experiment and which differential expression tool should you use?" Rna (2016)
🔗 Zhao, Shilin, *et al.* "RnaSeqSampleSize: real data based sample size estimation for RNA sequencing." BMC bioinformatics 19.1 (2018): 191

# RNA extraction

- Sample processing and storage
- Total RNA/mRNA/small RNA
- DNAse treatment
- Quantity & quality
- RIN values (Strong effect)
- Batch effect
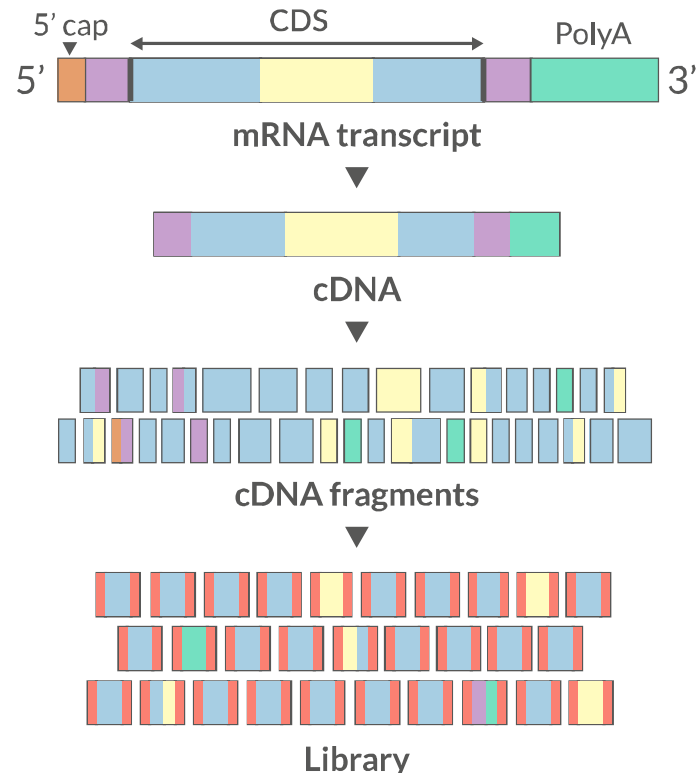
Romero, Irene Gallego, *et al*. "RNA-seq: impact of RNA degradation on transcript quantification." BMC biology 12.1 (2014): 42
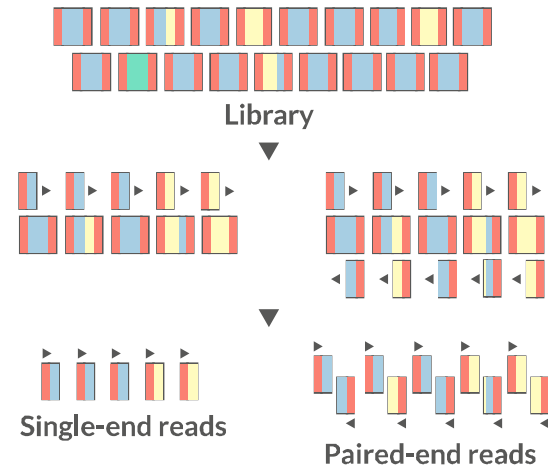
Kim, Young-Kook, *et al*. "Short structured RNAs with low GC content are selectively lost during extraction from a small number of cells." Molecular cell 46.6 (2012): 893-89500481-9).

# Library prep

- PolyA selection
- rRNA depletion
- Size selection
- PCR amplification (See section PCR duplicates)
- Stranded (directional) libraries
    - Accurately identify sense/antisense transcript
    - Resolve overlapping genes
- Exome capture
- Library normalisation
- Batch effect



5' cap     CDS     PolyA

5'                 3'

**mRNA transcript**

**cDNA**

**cDNA fragments**

**Library**

🔗 Zhao, Shanrong, et al. "Comparison of stranded and non-stranded RNA-seq transcriptome profiling and investigation of gene overlap." BMC genomics 16.1 (2015): 675

🔗 Levin, Joshua Z., et al. "Comprehensive comparative analysis of strand-specific RNA sequencing methods." Nature methods 7.9 (2010): 709

# Sequencing

- Sequencer (Illumina/PacBio)
- Read length

  - Greater than 50bp does not improve DGE
  - Longer reads better for isoforms

- Pooling samples
- Sequencing depth (Coverage/Reads per sample)
- Single-end reads (Cheaper)
- Paired-end reads

  - Increased mappable reads
  - Increased power in assemblies
  - Better for structural variation and isoforms
  - Decreased false-positives for DGE



Library

Single-end reads

Paired-end reads

🔗 Chhangawala, Sagar, et al. "The impact of read length on quantification of differentially expressed genes and splice junction detection." Genome biology 16.1 (2015): 131

🔗 Corley, Susan M., et al. "Differentially expressed genes from RNA-Seq and functional enrichment results are affected by the choice of single-end versus paired-end reads and stranded versus non-stranded protocols." BMC genomics 18.1 (2017): 399

🔗 Liu, Yuwen, Jie Zhou, and Kevin P. White. "RNA-seq differential expression studies: more sequence or more replication?." Bioinformatics 30.3 (2013): 301-304

🔗 Comparison of PE and SE for RNA-Seq, SciLifeLab

# Workflow | DGE

| Reads | FastQ | FastQ | FastQ |
|---|---|---|---|
| | ▼ | ▼ | ▼ |
| Mapping | STAR | HiSat2 | [Kallisto/ Salmon] |
| | ▼ | ▼ | |
| Quantification | featureCounts | StringTie | |
| | ▼ | ▼ | ▼ |
| Differential gene expression | DESeq2/ edgeR/ Limma | Ballgown | Sleuth |

# De-Novo assembly

- When no reference genome available
- To identify novel genes/transcripts/isoforms
- Identify fusion genes
- Assemble transcriptome from short reads
- Access quality of assembly and refine
- Map reads back to assembled transcriptome

🧰 Trinity, SOAPdenovo-Trans, Oases, rnaSPAdes

🔗 Hsieh, Ping-Han *et al.*, "Effect of de novo transcriptome assembly on transcript quantification" 2018 bioRxiv 380998

🔗 Wang, Sufang, and Michael Gribskov. "Comprehensive evaluation of de novo transcriptome assembly programs and their effects on differential gene expression analysis." Bioinformatics 33.3 (2017): 327-333

# Read QC

- Number of reads
- Per base sequence quality
- Per sequence quality score
- Per base sequence content
- Per sequence GC content
- Per base N content
- Sequence length distribution
- Sequence duplication levels
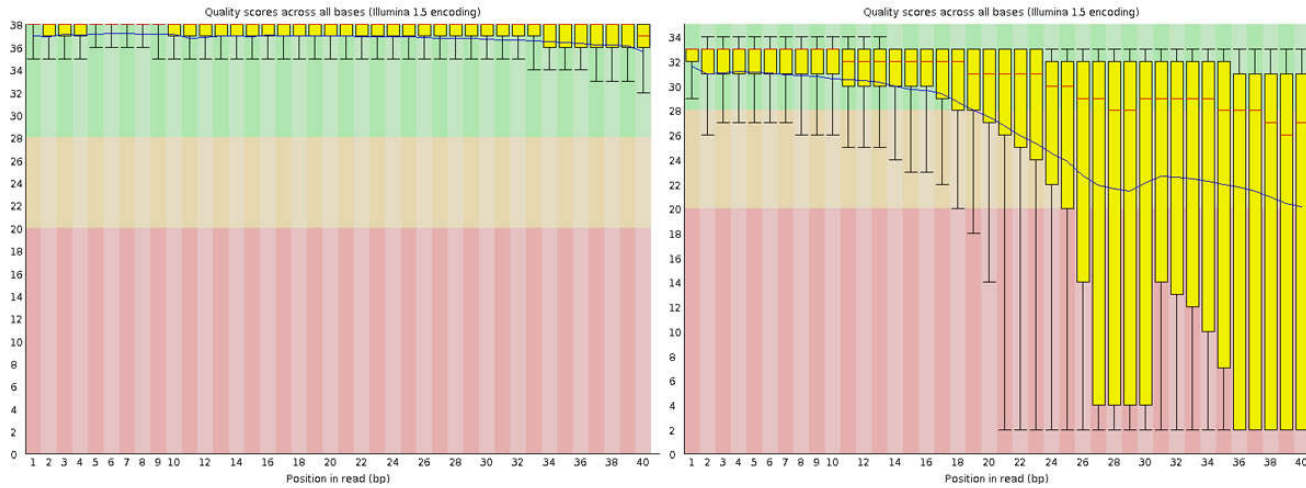- Overrepresented sequences
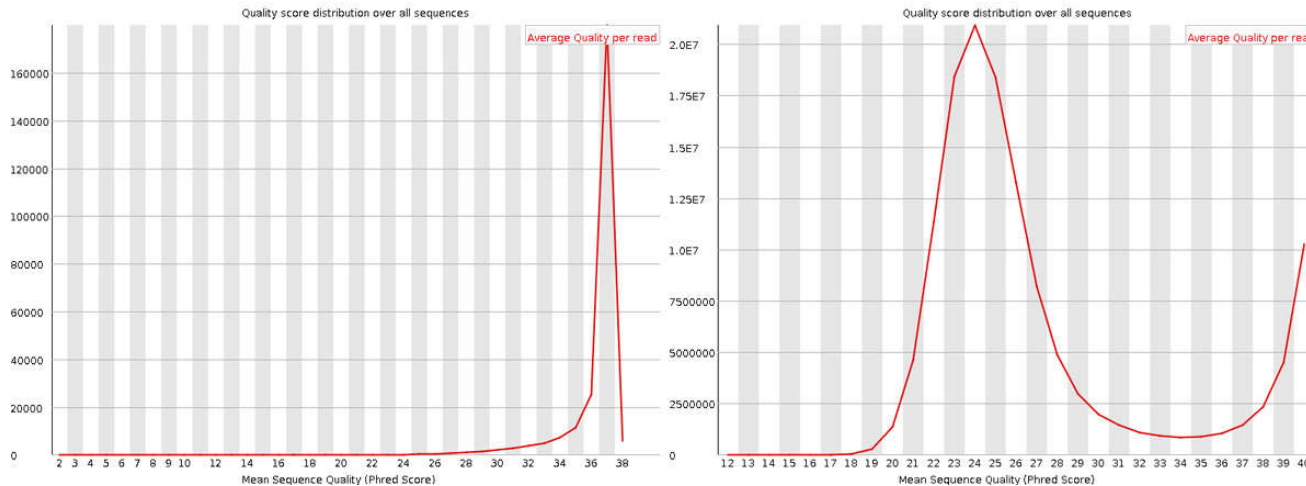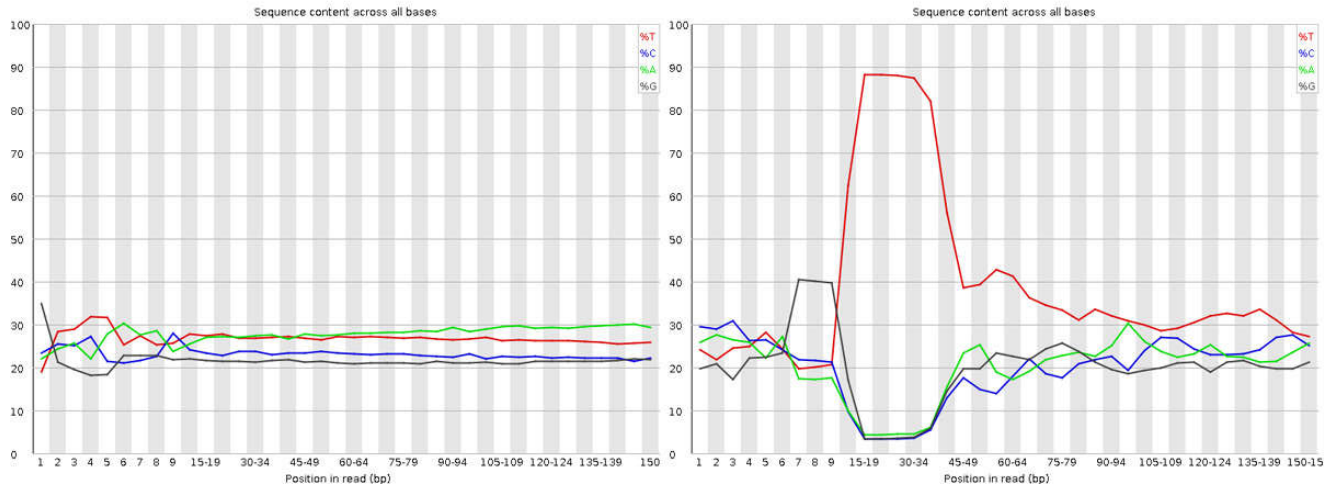- Adapter content
- Kmer content

🧰 FastQC, MultiQC

https://sequencing.qcfail.com/

# Read QC | PBSQ, PSQS

**Per base sequence quality**



**Per sequence quality scores**

# Read QC | PBSC, PSGC

**Per base sequence content**
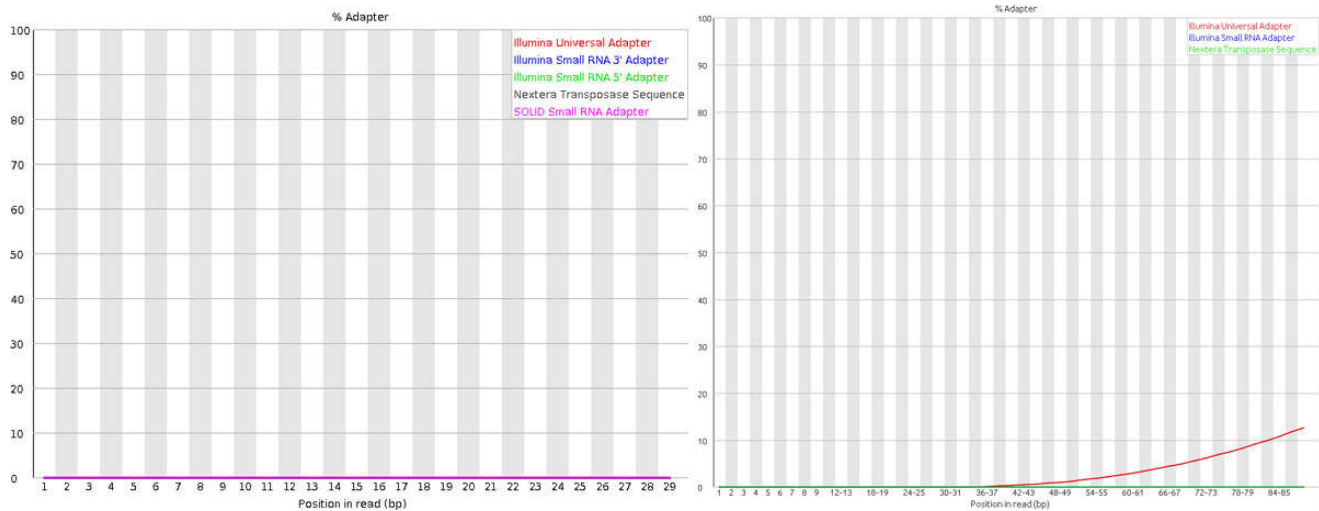


**Per sequence GC content**

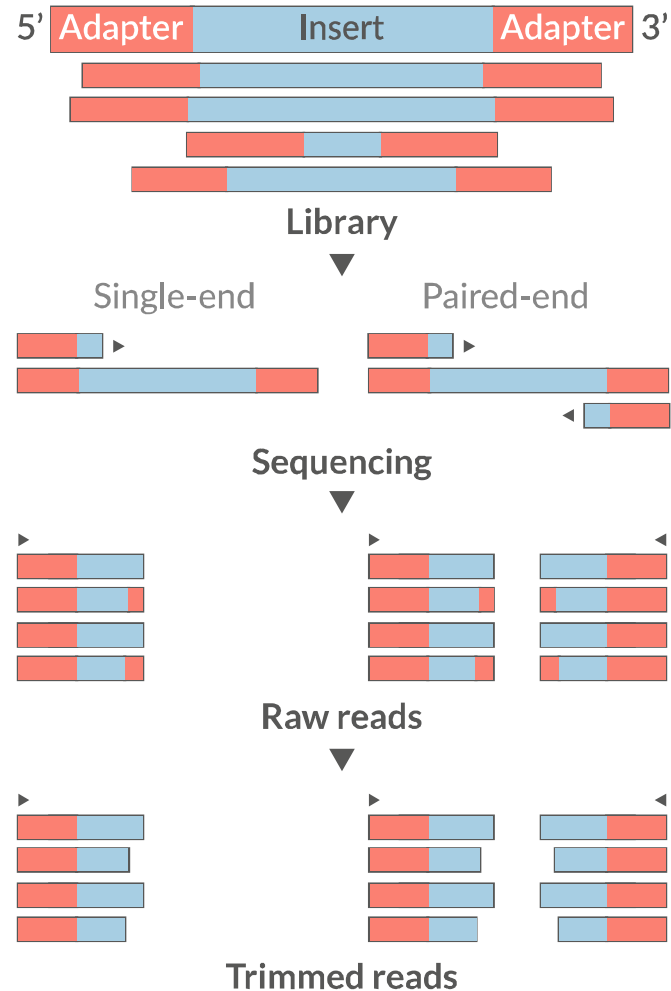# Read QC | SDL, AC

**Sequence duplication level**



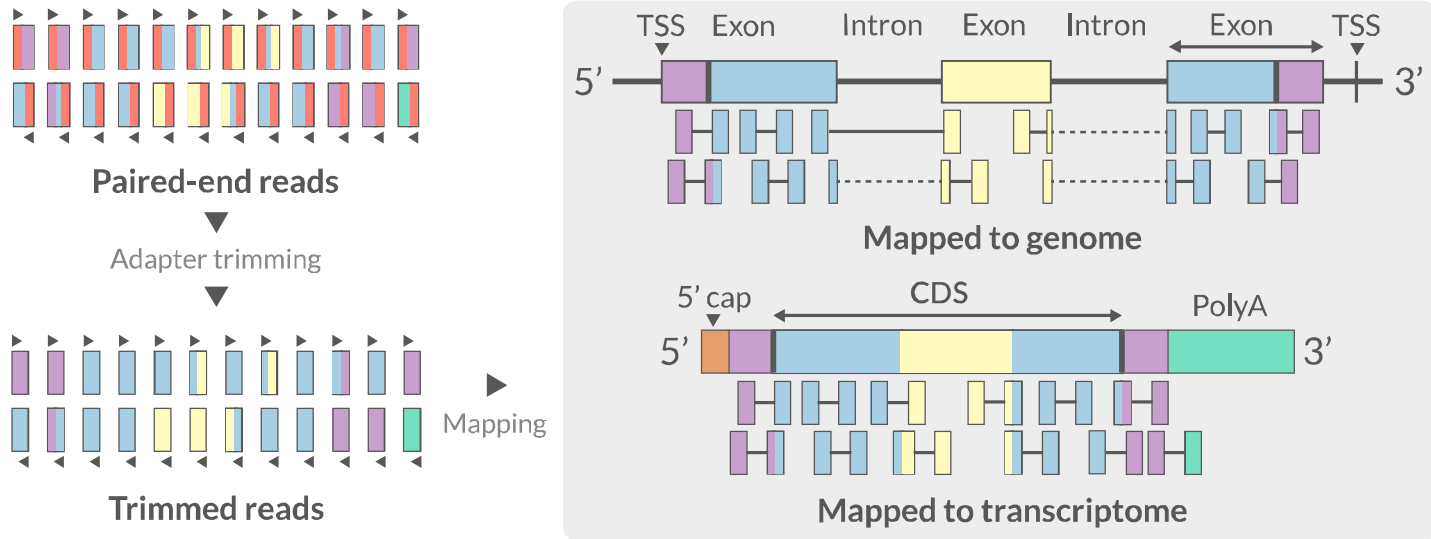**Adapter content**

# Trimming

- Trim IF necessary

  - Synthetic bases can be an issue for SNP calling
  - Insert size distribution may be more important for assemblers

- Trim/Clip/Filter reads
- Remove adapter sequences
- Trim reads by quality
- Sliding window trimming
- Filter by min/max read length

  - Remove reads less than ~22nt

- Demultiplexing/Splitting

🧰 Cutadapt, fastp, Skewer, Prinseq

# Mapping

**Paired-end reads**

Adapter trimming

**Trimmed reads**

Mapping

**Mapped to genome**

**Mapped to transcriptome**

- Aligning reads back to a reference sequence
- Mapping to genome vs transcriptome
- Splice-aware alignment (genome)

STAR, HiSat2, GSNAP, Novoalign (Commercial)

Baruzzo, Giacomo, *et al.* "Simulation-based comprehensive benchmarking of RNA-seq aligners." Nature methods 14.2 (2017): 135

# Mapping

- Reads (FASTQ)

```
@ST-E00274:179:HHYMLALXX:8:1101:1641:1309 1:N:0:NGATGT
NCATCGTGGTATTTGCACATCTTTTCTTATCAAATAAAAAGTTTAACCTACTCAGTTATGCGCATACGTTTTTTGATGGCATTT
+
#AAAFAFA<-AFFJJJAFA-FFJJJJFFFAJJJJ-<FFJJJ-A-F-7--FA7F7-----FFFJFA<FFFFJ<AJ--FF-A<A-<.
```

```
@instrument:runid:flowcellid:lane:tile:xpos:ypos
read:isfiltered:controlnumber:sampleid
```

- Reference Genome/Transcriptome (FASTA)

```
>1 dna:chromosome chromosome:GRCz10:1:1:58871917:1 REF
GATCTTAAACATTTATTCCCCCTGCAAACATTTTCAATCATTACATTGTCATTTCCCCTC
CAAATTAAATTTAGCCAGAGGCGCACAACATACGACCTCTAAAAAAGGTGCTGTAACATG
```

- Annotation (GTF/GFF)

```
#!genome-build GRCz10
#!genebuild-last-updated 2016-11
4       ensembl_havana  gene    6732    52059   .       -       .       gene_id "ENS
```

```
seq source feature start end score strand frame attribute
```
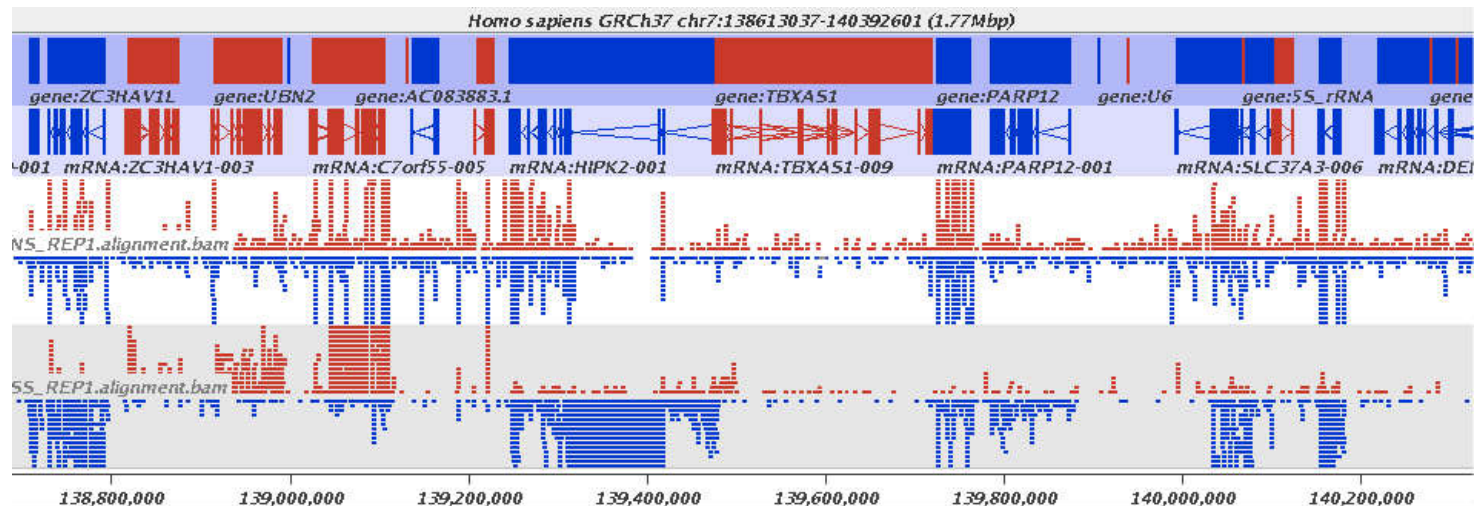
# Alignment

- SAM/BAM (Sequence Alignment Map format)

```
ST-E00274:188:H3JWNCCXY:4:1102:32431:49900        163      1        1        60        8S13
```

query flag ref pos mapq cigar mrnm mpos tlen seq qual opt
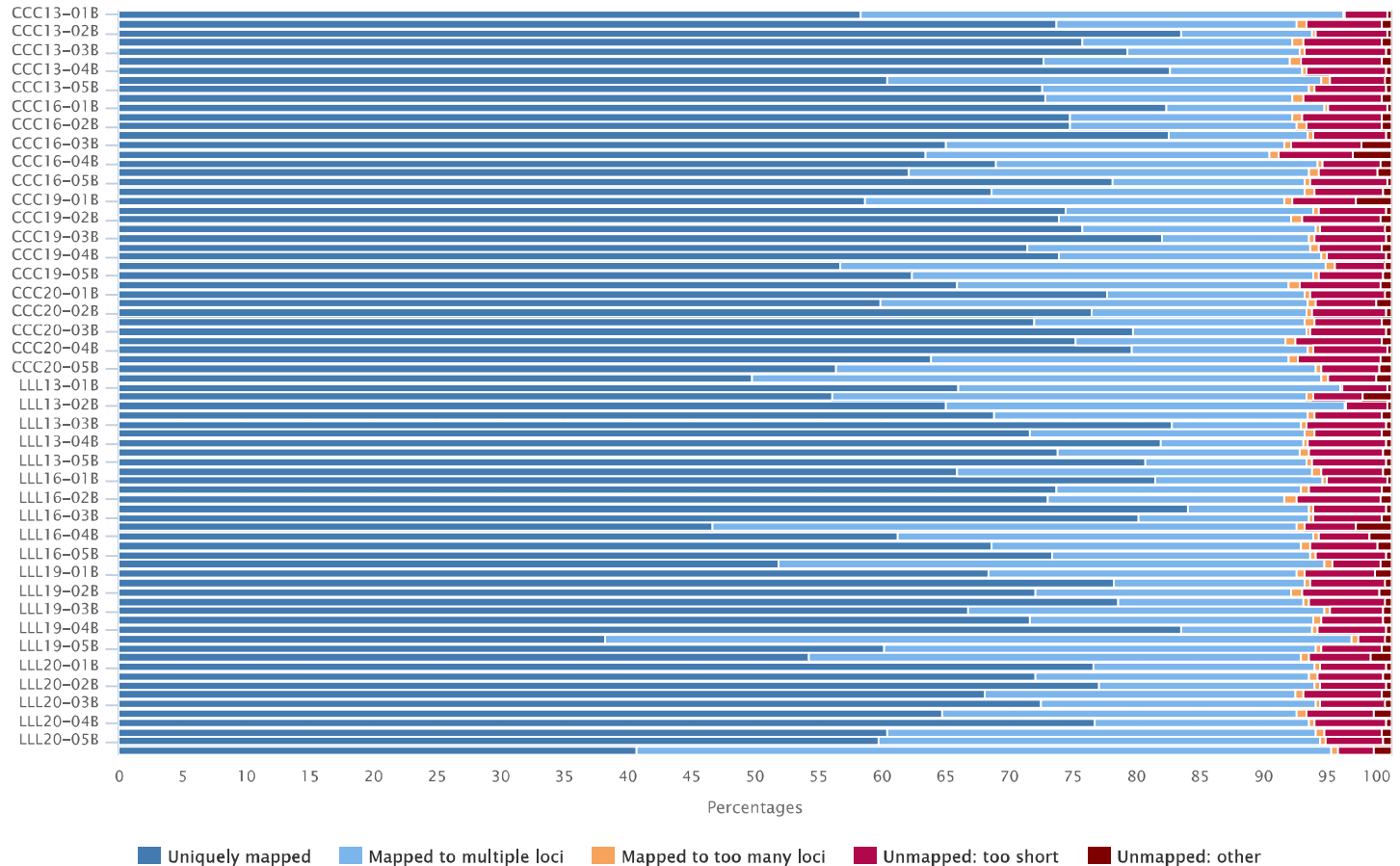
🧰 SeqMonk, IGV, UCSC Genome Browser

# Alignment QC

- Number of reads mapped/unmapped/paired etc
- Uniquely mapped
- Insert size distribution
- Coverage
- Gene body coverage
- Biotype counts / Chromosome counts
- Counts by region: gene/intron/non-genic
- Sequencing saturation
- Strand specificity

🧰 STAR (final log file), samtools > stats, bamtools > stats, QoRTs, RSeQC, Qualimap

# Alignment QC | STAR Log
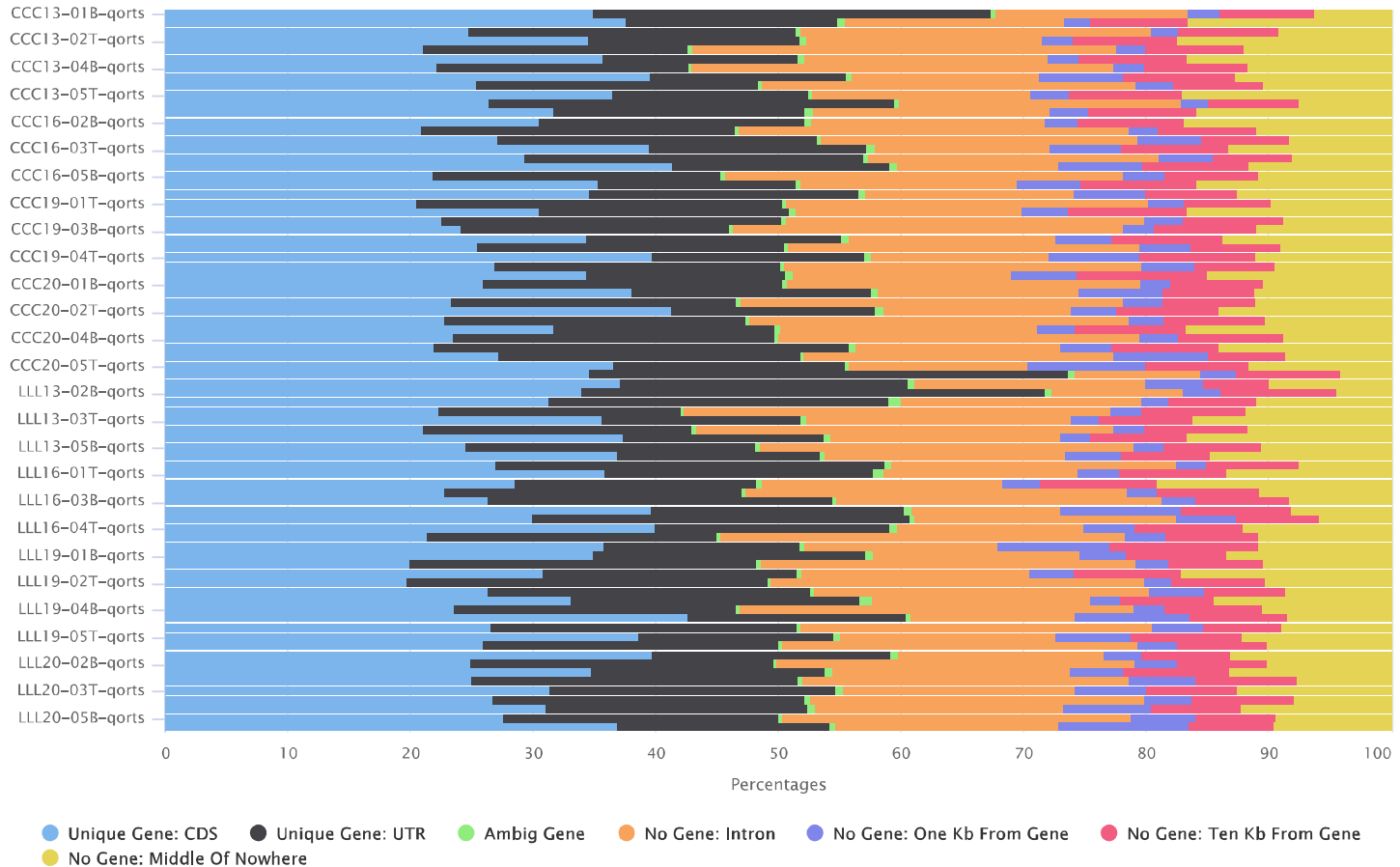
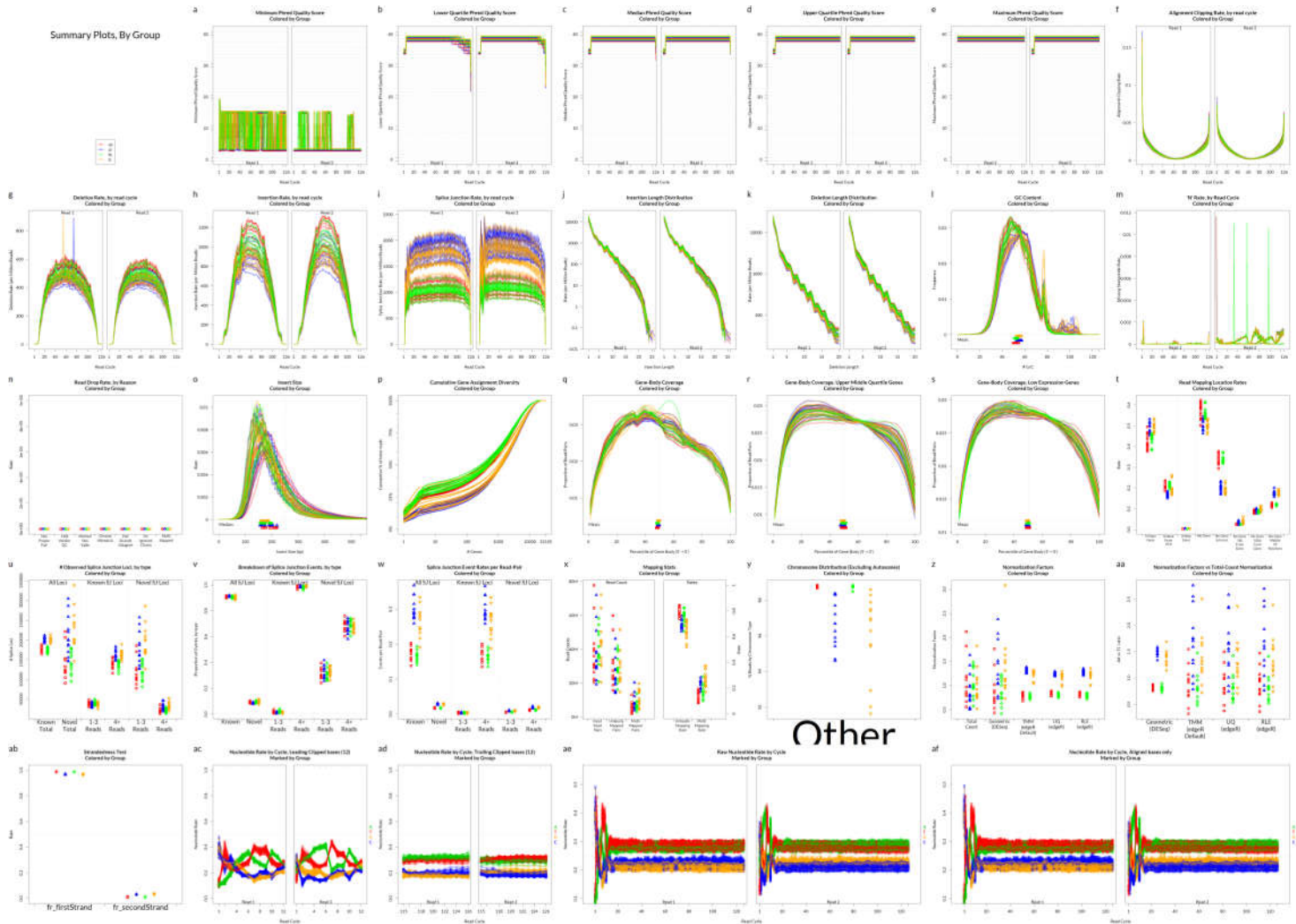MultiQC can be used to summarise and plot STAR log files.



STAR Alignment Scores

# Alignment QC | Features

QoRTs was run on all samples and summarised using MultiQC.
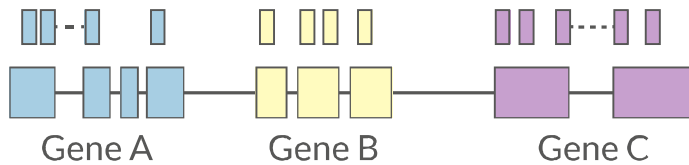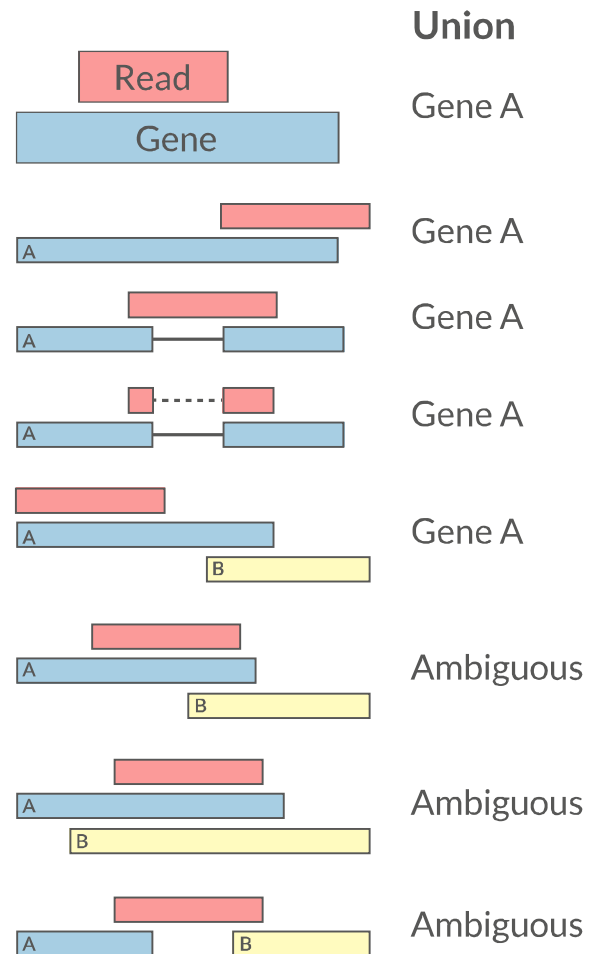


QoRTs: Alignment Locations

# QoRTs

# Quantification | Counts

- Read counts = gene expression
- Reads can be quantified on any feature (gene, transcript, exon etc)
- Intersection on gene models
- Gene/Transcript level

**featureCounts**, **HTSeq**



Gene A       Gene B       Gene C

**Union**

Read

Gene

Gene A

Gene A

Gene A

Gene A

Gene A

Ambiguous

Ambiguous

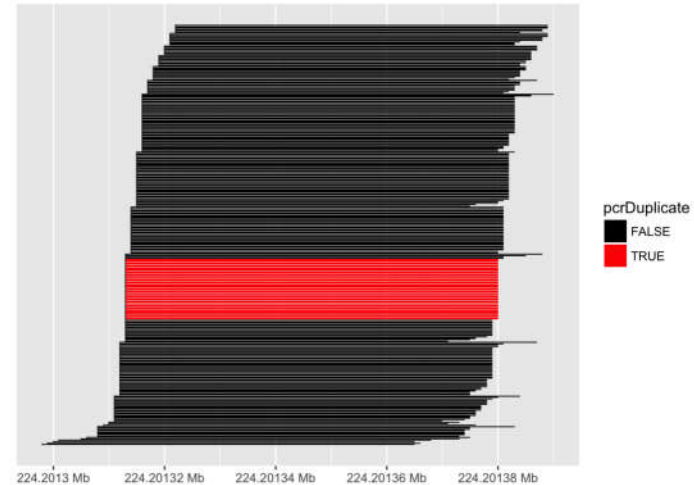Ambiguous

# Quantification

## PCR duplicates

- Ignore for RNA-Seq data
- Computational deduplication (Don't!)
- Use PCR-free library-prep kits
- Use UMIs during library-prep

## Multi-mapping

- Added (BEDTools multicov)
- Discard (featureCounts, HTSeq)
- Distribute counts (Cufflinks)
- Rescue

  - Probabilistic assignment (Rcount, Cufflinks)
  - Prioritise features (Rcount)
  - Probabilistic assignment with EM (RSEM)

🔗 Fu, Yu, *et al*. "Elimination of PCR duplicates in RNA-seq and small RNA-seq using unique molecular identifiers." BMC genomics 19.1 (2018): 531
🔗 Parekh, Swati, *et al*. "The impact of amplification on differential expression analyses by RNA-seq." Scientific reports 6 (2016): 25533
🔗 Klepikova, Anna V., *et al*. "Effect of method of deduplication on estimation of differential gene expression using RNA-seq." PeerJ 5 (2017): e3091

# Quantification | Abundance

- Count methods

    - Provide no inference on isoforms
    - Cannot accurately measure fold change

- Probabilistic assignment

    - Deconvolute ambiguous mappings
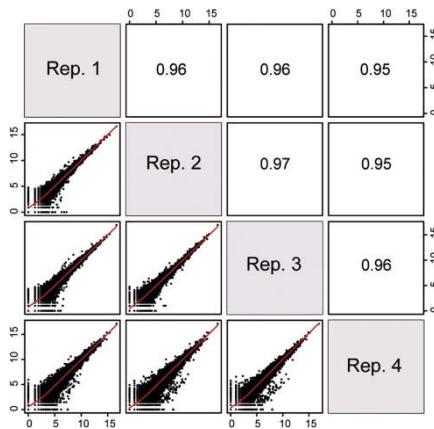    - Transcript-level
    - cDNA reference

**Kallisto, Salmon**

- Ultra-fast & alignment-free
- Subsampling & quantification confidence
- Transcript-level estimates improves gene-level estimates
- Kallisto/Salmon > transcript-counts > `tximport()` > gene-counts

🧰 RSEM, Kallisto, Salmon, Cufflinks2

🔗 Soneson, Charlotte, *et al.* "Differential analyses for RNA-seq: transcript-level estimates improve gene-level inferences." F1000Research 4 (2015)
🔗 Zhang, Chi, *et al.* "Evaluation and comparison of computational tools for RNA-seq isoform quantification." BMC genomics 18.1 (2017): 583

# Quantification QC

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| ENSG00000000003 | 140 | 242 | 188 | 143 | 287 | 344 | 438 | 280 | 253 |
| ENSG00000000005 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| ENSG00000000419 | 69 | 98 | 77 | 55 | 52 | 94 | 116 | 79 | 69 |
| ENSG00000000457 | 56 | 75 | 104 | 79 | 157 | 205 | 183 | 178 | 153 |
| ENSG00000000460 | 33 | 27 | 23 | 19 | 27 | 42 | 69 | 44 | 40 |
| ENSG00000000938 | 7 | 38 | 13 | 17 | 35 | 76 | 53 | 37 | 24 |
| ENSG00000000971 | 545 | 878 | 694 | 636 | 647 | 216 | 492 | 798 | 323 |
| ENSG00000001036 | 79 | 154 | 74 | 80 | 128 | 167 | 220 | 147 | 72 |

- Pairwise correlation between samples must be high (>0.9)

- Count QC using RNASeqComp





RNASeqComp

Teng, Mingxiang, *et al*. "A benchmark for RNA-seq quantification pipelines." Genome biology 17.1 (2016): 74
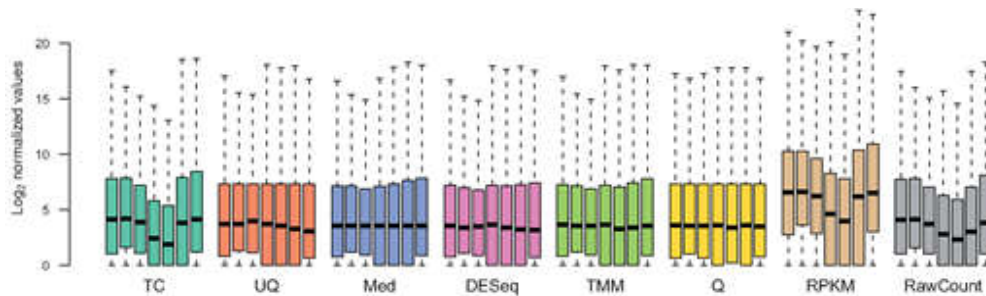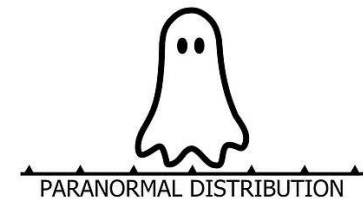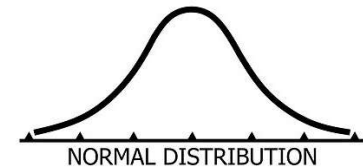
# Normalisation

- Control for Sequencing depth & compositional bias
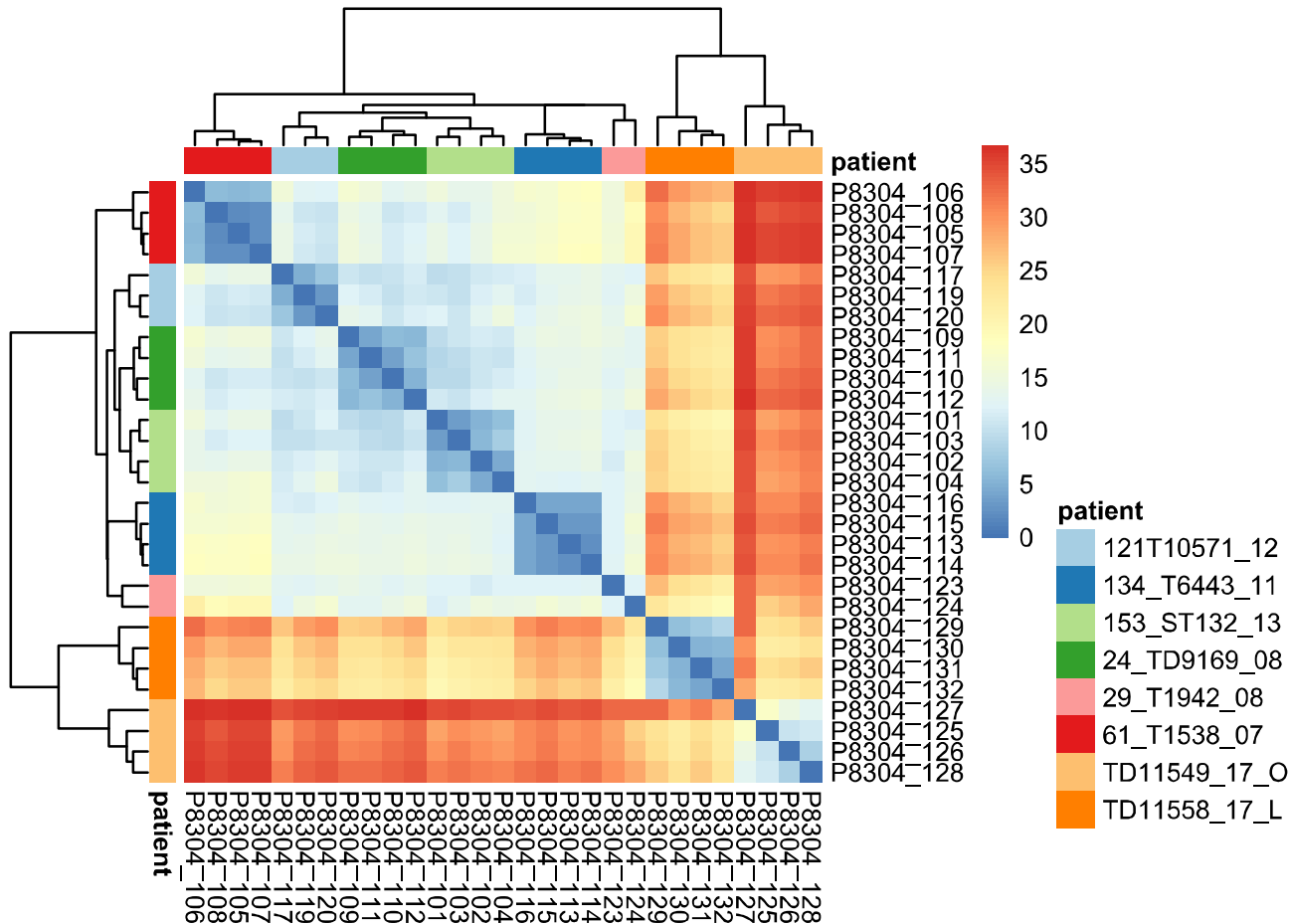- Median of Ratios (DESeq2) and TMM (edgeR) perform the best



- FoR DGE using DGE packages, use raw counts
- For clustering, heatmaps etc use VST, VOOM or RLOG
- For own analysis, plots etc, use TPM
- Other solutions: spike-ins/house-keeping genes



NORMAL DISTRIBUTION

PARANORMAL DISTRIBUTION

Dillies, Marie-Agnes, *et al.* "A comprehensive evaluation of normalization methods for Illumina high-throughput RNA sequencing data analysis." Briefings in bioinformatics 14.6 (2013): 671-683

Evans, Ciaran, Johanna Hardin, and Daniel M. Stoebel. "Selecting between-sample RNA-Seq normalization methods from the perspective of their assumptions." Briefings in bioinformatics (2017)

Wagner, Gunter P., Koryu Kin, and Vincent J. Lynch. "Measurement of mRNA abundance using RNA-seq data: RPKM measure is inconsistent among samples." Theory in biosciences 131.4 (2012): 281-285

# Exploratory | Heatmap

- Remove lowly expressed genes
- Transform raw counts to VST, VOOM, RLOG, TPM etc
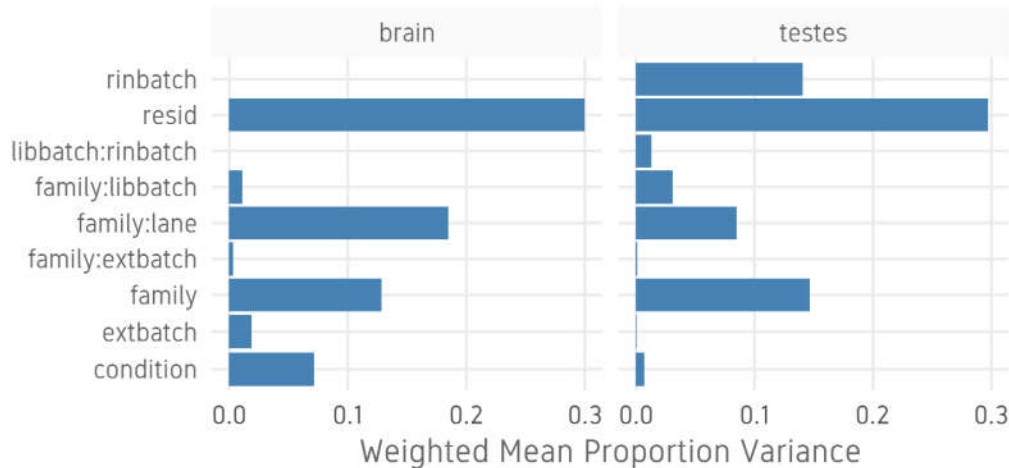- Sample-sample clustering heatmap

NB*S SciLifeLab

- 121T10571_12
- 134_T6443_11
- 153_ST132_13
- 24_TD9169_08
- 29_T1942_08
- 61_T1538_07
- TD11549_17_O
- TD11558_17_L

`cmdscale()` , plotly
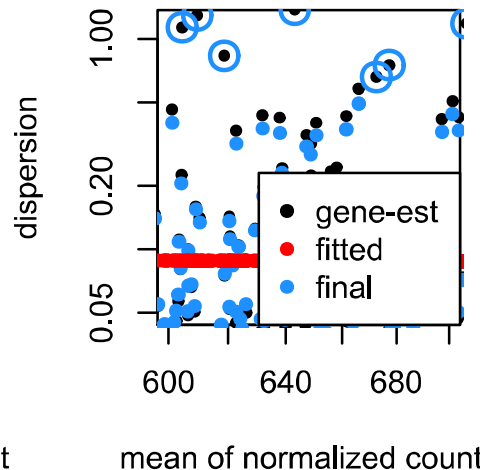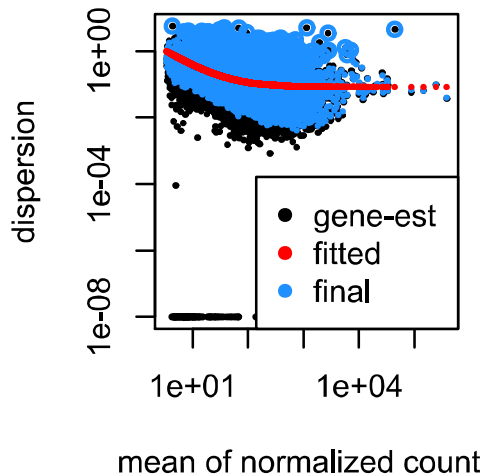
# Batch correction

- Estimate variation explained by variables (PVCA)



- Find confounding effects as surrogate variables (SVA)
- Model known batches in the LM/GLM model
- Correct known batches (ComBat)(Harsh!)
- Interactively evaluate batch effects and correction (BatchQC)

🧰 SVA, PVCA, BatchQC

🔗 Liu, Qian, and Marianthi Markatou. "Evaluation of methods in removing batch effects on RNA-seq data." Infectious Diseases and Translational Medicine 2.1 (2016): 3-9
🔗 Manimaran, Solaiappan, et al. "BatchQC: interactive software for evaluating sample and batch effects in genomic data." Bioinformatics 32.24 (2016): 3836-3838

# DGE

- DESeq2, edgeR (Neg-binom > GLM > Test), Limma-Voom (Neg-binom > Voom-transform > LM > Test)
- DESeq2 `~age+condition`

  - Estimate size factors `estimateSizeFactors()`
  - Estimate gene-wise dispersion `estimateDispersions()`
  - Fit curve to gene-wise dispersion estimates
  - Shrink gene-wise dispersion estimates
  - GLM fit for each gene
  - Wald test `nbinomWaldTest()`



🧰 DESeq2, edgeR, Limma-Voom

🔗 Seyednasrollah, Fatemeh, *et al*. "Comparison of software packages for detecting differential expression in RNA-seq studies." Briefings in bioinformatics 16.1 (2013): 59-70
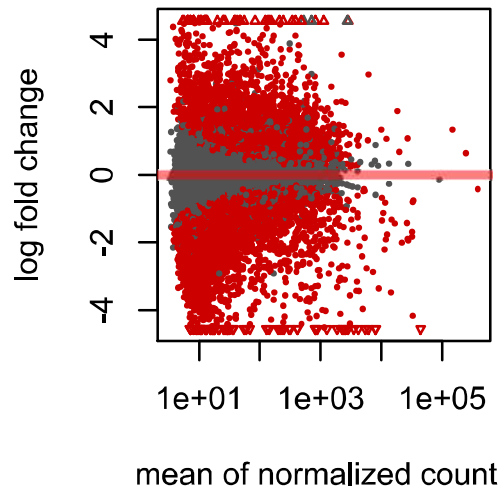
# DGE

- Results `results()`

```
## log2 fold change (MLE): type type2 vs control
## Wald test p-value: type type2 vs control
## DataFrame with 1 row and 6 columns
##                        baseMean      log2FoldChange              lfcSE
##                       <numeric>           <numeric>          <numeric>
## ENSG00000000003 242.307796723287 -0.932926089608558 0.114285150312647
##                            stat             pvalue
##                       <numeric>          <numeric>
## ENSG00000000003 -8.16314356726468 3.26416150312236e-16
##                            padj
##                       <numeric>
## ENSG00000000003 1.36240610027518e-14
```
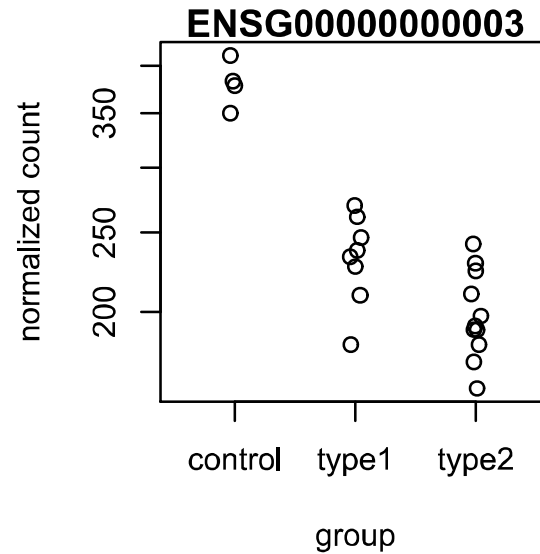
- Summary `summary()`

```
##
## out of 17889 with nonzero total read count
## adjusted p-value < 0.1
## LFC > 0 (up)       : 4526, 25%
## LFC < 0 (down)     : 5062, 28%
## outliers [1]       : 25, 0.14%
## low counts [2]     : 0, 0%
## (mean count < 3)
## [1] see 'cooksCutoff' argument of ?results
## [2] see 'independentFiltering' argument of ?results
```
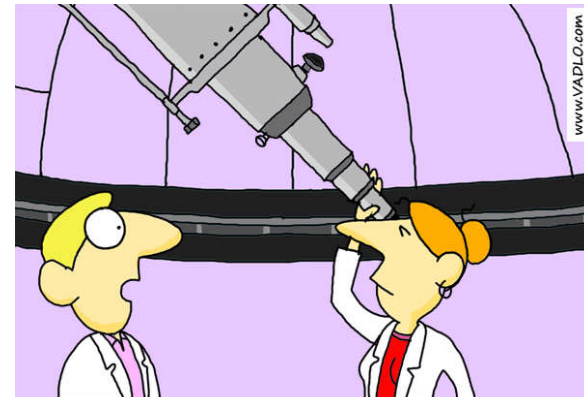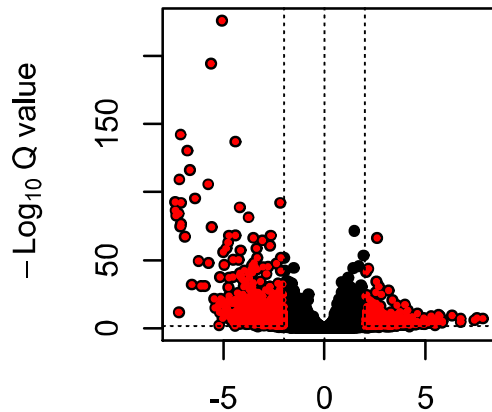
- MA plot `plotMA()`



mean of normalized count

- Volcano plot



- Normalised counts `plotCounts()`



**ENSG00000000003**

group
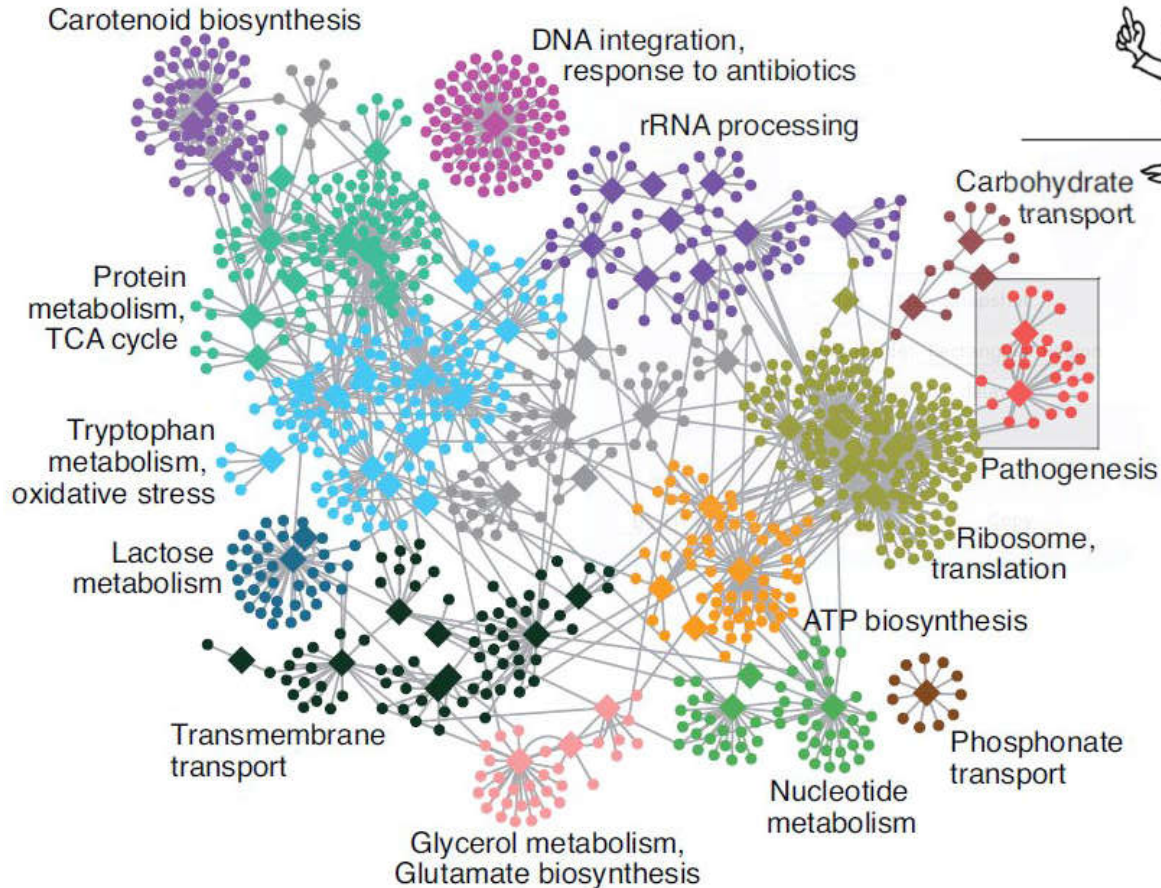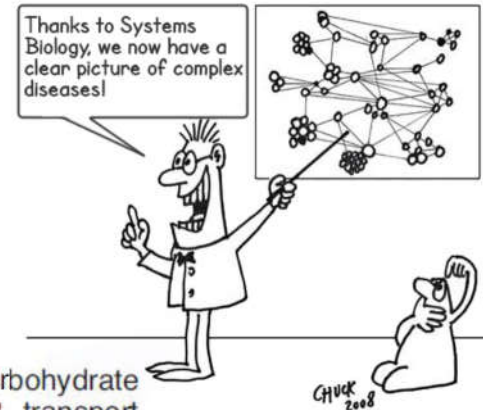


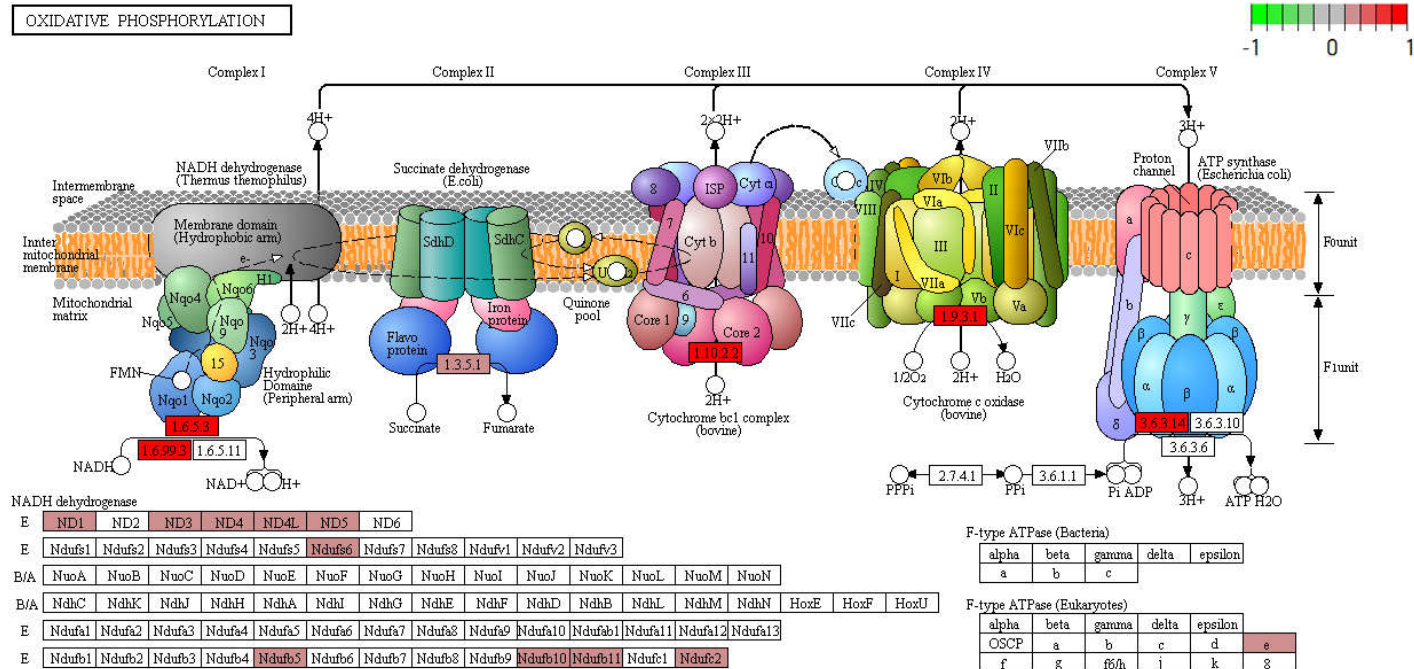"Can you see the upper points of my scatter plot?"

# Functional analysis | GO

- Gene enrichment analysis
- Gene set enrichment analysis (GSEA)
- Gene ontology / Reactome databases
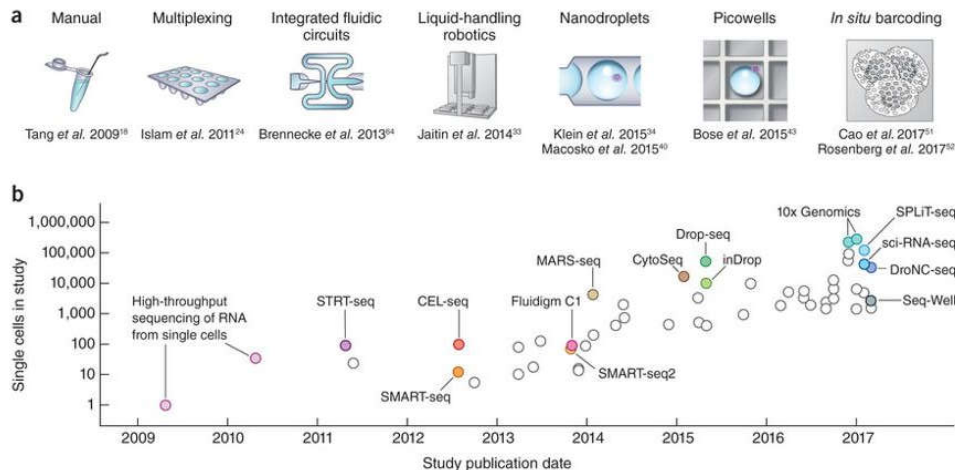
# Functional analysis | Kegg

- Pathway analysis (Kegg)



🧰 DAVID, clusterProfiler, ClueGO, ErmineJ, pathview

# Single cell RNA-Seq

- Bulk RNA-Seq measures mean expression-level over many cells
- Poor resolution for development, differentiation, heterogenous tissues
- Identify cell types in a tissue
- Temporal/spatial/conditional change in cellular state and composition



- Zero-inflated data (~80% missing data)
- Transcriptional bursting, drop-out
- Low RNA, Poor capture efficiency
- Amplification bias and background noise

Svensson, Valentine, Roser Vento-Tormo, and Sarah A. Teichmann. "Exponential scaling of single-cell RNA-seq in the past decade." Nature protocols 13.4 (2018): 599

🔗 Sebe-Pedros, Arnau, et al. "Cnidarian Cell Type Diversity and Regulation Revealed by Whole-Organism Single-Cell RNA-Seq." Cell 173.6 (2018): 1520-1534
🔗 Plass, Mireya, et al. "Cell type atlas and lineage tree of a whole complex animal by single-cell transcriptomics." Science 360.6391 (2018): eaaq1723

# New Advances

- **Long read single molecule RNA-Seq** (Zuo *et al.*, 2018)

Research | Open Access

### Revealing the transcriptomic complexity of switchgrass by PacBio long-read sequencing

Chunman Zuo, Matthew Blow, Avinash Sreedasyam, Rita C. Kuo, Govindarajan Kunde Ramamoorthy, Ivone Torres-Jerez, Guifen Li, Mei Wang, David Dilworth, Kerrie Barry, Michael Udvardi, Jeremy Schmutz, Yuhong Tang ✉ and Ying Xu ✉

*Biotechnology for Biofuels* 2018 **11**:170

- **Single-cell isoform RNA-Seq** (Ishaan *et al.*, 2018)

### Single-cell isoform RNA sequencing (ScISOr-Seq) across thousands of cells reveals isoforms of cerebellar cell types.
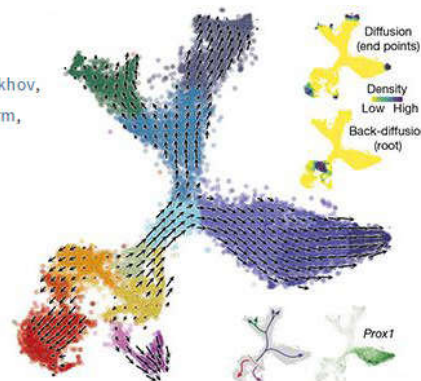
Ishaan Gupta, Paul G Collier, Bettina Haase, Ahmed Mahfouz, Anoushka Joglekar, Taylor Floyd, Frank Koopmans, Ben Barres, August B Smit, Steven Sloan, Wenjie Luo, Olivier Fedrigo, M Elizabeth Ross, Hagen U Tilgner

- **Single-cell lineage tracing** (Manno *et al.*, 2018)

### RNA velocity of single cells

Gioele La Manno, Ruslan Soldatov, Amit Zeisel, Emelie Braun, Hannah Hochgerner, Viktor Petukhov, Katja Lidschreiber, Maria E. Kastriti, Peter Lönnerberg, Alessandro Furlan, Jean Fan, Lars E. Borm, Zehua Liu, David van Bruggen, Jimin Guo, Xiaoling He, Roger Barker, Erik Sundström, Gonçalo Castelo-Branco, Patrick Cramer, Igor Adameyko, Sten Linnarsson ✉ & Peter V. Kharchenko ✉

*Nature* **560**, 494–498 (2018) | Download Citation ⤓

# Summary

- Sound experimental design to avoid confounding
- Plan carefully about lib prep, sequencing etc based on experimental objective
- Biological replicates may be more important than paired-end reads or long reads
- Discard low quality bases, reads, genes and samples
- Verify that tools and methods align with data assumptions
- Experiment with multiple pipelines and tools
- QC! QC everything at every step

🔗 Conesa, Ana, *et al.* "A survey of best practices for RNA-seq data analysis." Genome biology 17.1 (2016): 13

# Further learning

- Griffith lab RNA-Seq using HiSat & StringTie tutorial
- SciLifeLab courses
- HBC Training DGE using DeSeq2 tutorial
- Hemberg lab scRNA-Seq tutorial
- RNA-Seq Blog

# Thank you! Questions?

Built on : 📅 25-Oct-2018 at 🕐 17:34:39

**2018** Roy Francis | SciLifeLab | NBIS

# Hands-On tutorial

NB<sup>IS</sup> SciLifeLab

**Main exercise**

- 01 Check the quality of the raw reads with **FastQC**
- 02 Map the reads to the reference genome using **Star**
- 03 Assess the post-alignment quality using **QualiMap**
- 04 Count the reads overlapping with genes using **featureCounts**
- 05 Build a statistical model to find DE genes using **edgeR** from a prepared R script

**Bonus exercises**

- 01 Functional annotation of DE genes using **GO/Reactome** databases
- 02 Visualisation of RNA-seq BAM files using **IGV** genome browser
- 03 RNA-Seq figures and plots using **R**
- 04 De-novo transcriptome assembly using **Trinity**

**Data directory:** `/sw/share/compstore/courses/ngsintro/rnaseq/`

**Work directory:** `/proj/g2018028/nobackup/<user>/rnaseq/`

# Hands-On tutorial

- Your work directory

`/proj/g2018028/nobackup/[user]/`

```
[user]/
rnaseq/
   +-- 1_raw/
   +-- 2_fastqc/
   +-- 3_mapping/
   +-- 4_qualimap/
   +-- 5_dge/
   +-- 6_multiqc/
   +-- reference/
   |    +-- mouse/
   |    +-- mouse_chr11/
   +-- scripts/
   +-- funannot/
   +-- assembly/
```

- Course data directory

`/sw/share/compstore/courses/ngsintro/rnaseq/`

```
rnaseq/
+-- bonus/
|    +-- assembly/
|    +-- exon/
|    +-- funannot/
|    +-- visual/
+-- documents/
+-- main/
|    +-- 1_raw/
|    +-- 2_fastqc/
|    +-- 3_mapping/
|    +-- 4_qualimap/
|    +-- 5_dge/
|    +-- 6_multiqc/
+-- reference/
|    +-- mouse/
|    +-- mouse_chr11/
+-- scripts/
```