

# Introduction to RNA-Seq

---

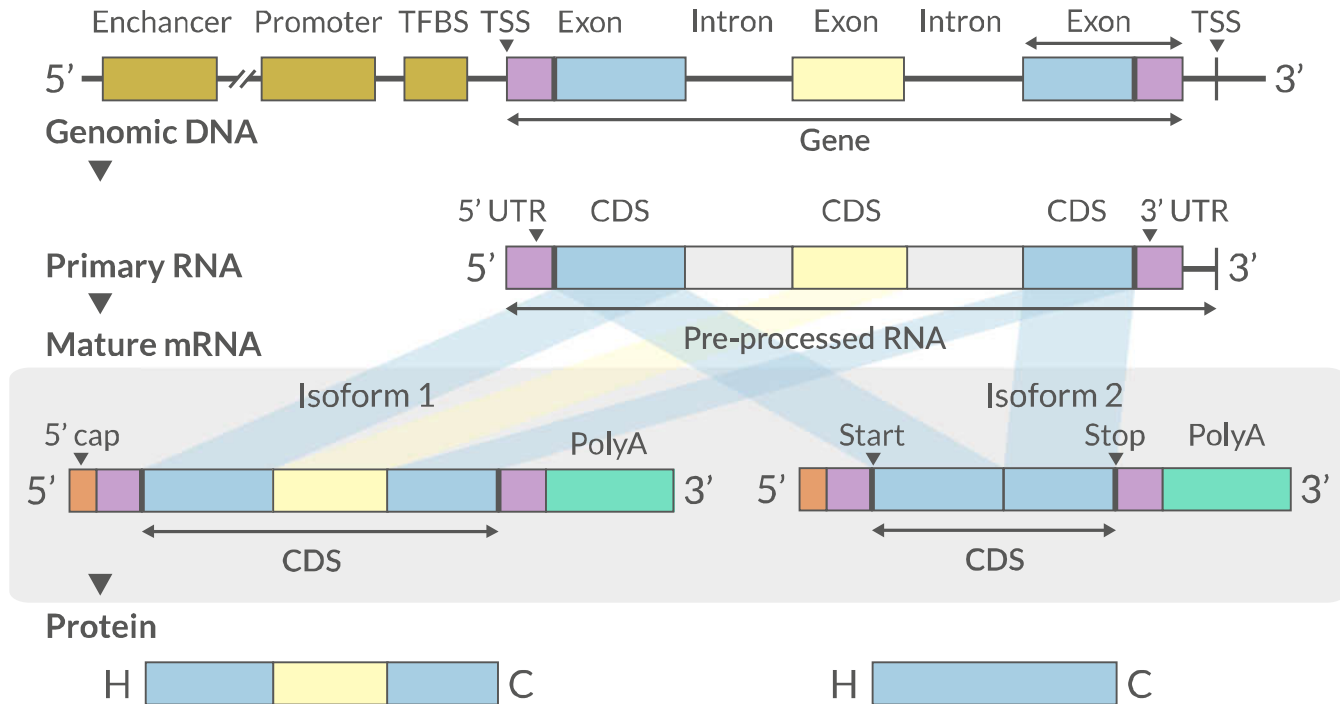
Introduction To Bioinformatics Using NGS Data

Roy Francis | 12-Sep-2018

# Contents

- Why RNA-Seq?
- Workflow
- DGE Workflow
- ReadQC
- Mapping
- Alignment QC
- Quantification
- Normalisation
- Exploratory
- DGE
- Functional analyses
- Single-cell RNA-Seq
- Summary
- Help

# Why sequence RNA?

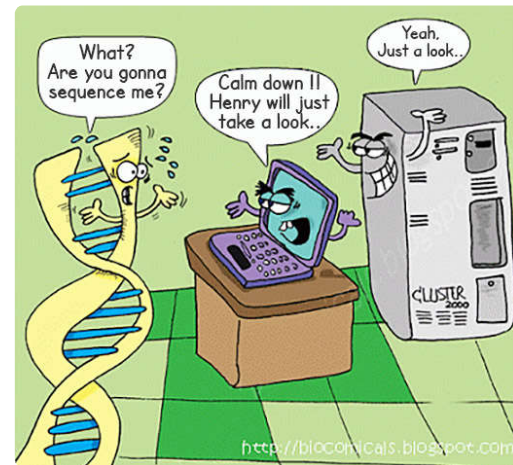
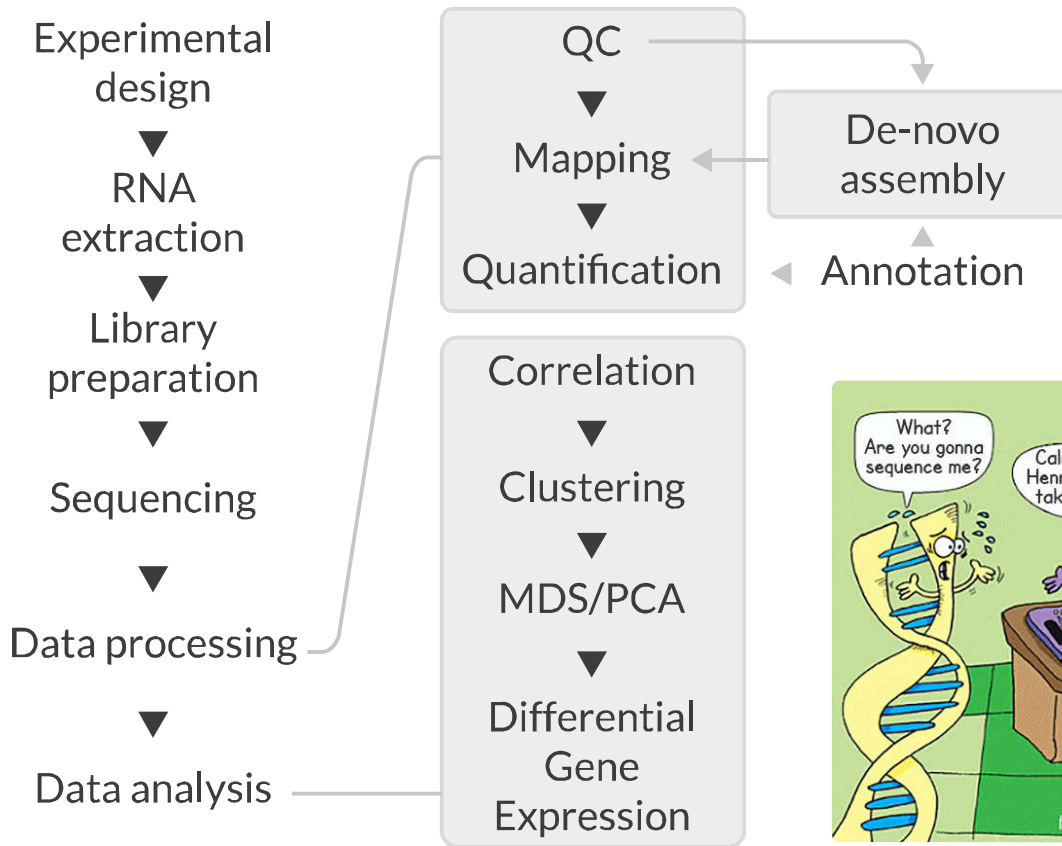


- The transcriptome is spatially and temporally dynamic
- Data comes from functional units (coding regions)
- Only a tiny fraction of the genome

# Applications

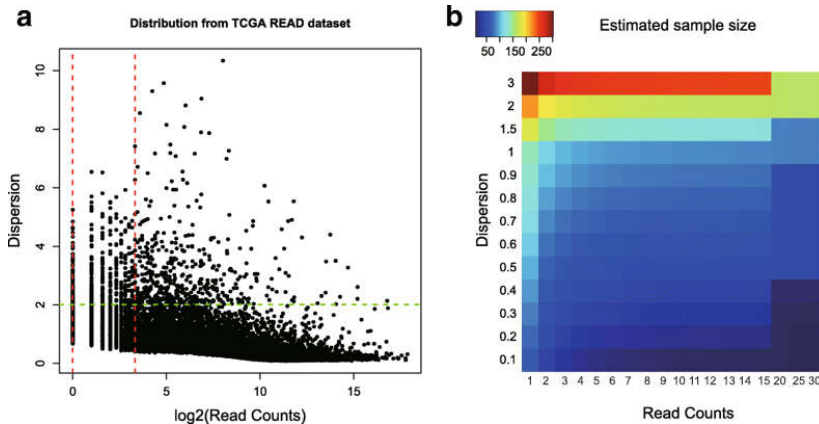
- Identify gene sequences in genomes
- Learn about gene function
- Differential gene expression
- Explore isoform and allelic expression
- Understand co-expression, pathways and networks
- Gene fusion
- RNA editing

# Workflow





# Experimental design


- 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)

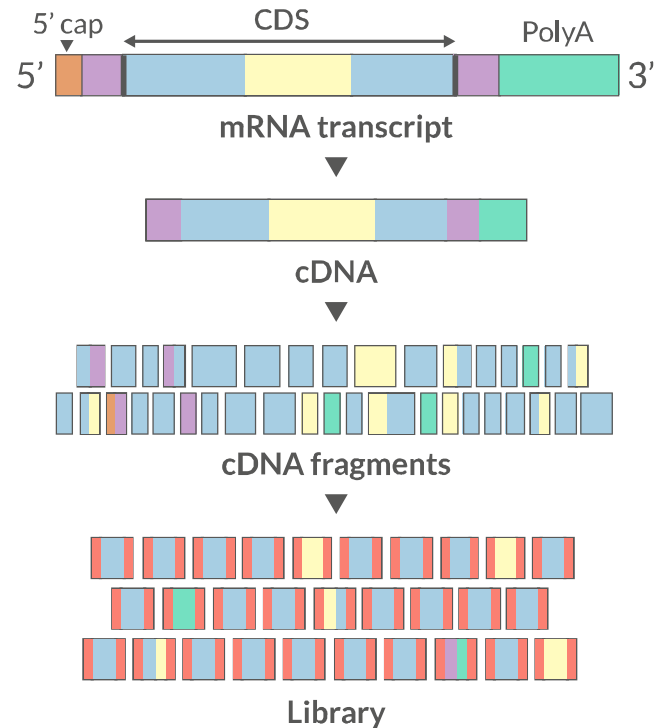
 Zhao, Shilin, *et al.* "RnaSeqSampleSize: real data based sample size estimation for RNA sequencing." *BMC bioinformatics* 19.1 (2018): 191

# RNA extraction

- Sample processing and storage
- Total RNA/mRNA/small RNA
- DNase treatment
- Quantity & quality
- RIN values (Strong effect)
- Batch effect

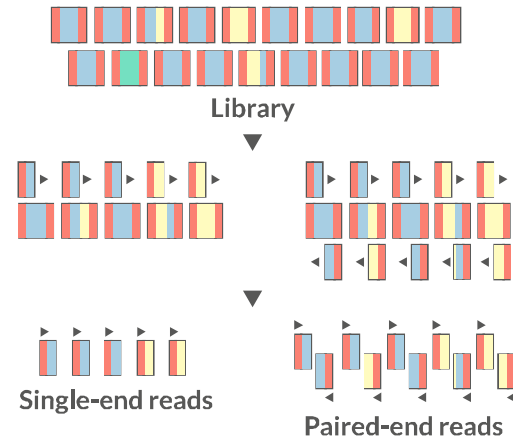
# Library prep

- 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





- 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



🔗 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
- Assess quality of assembly and refine
- Map reads back to assembled transcriptome

 [Trinity](#), [SOAPdenovo-Trans](#), [Oases](#), [rnaSPAdes](#)

# 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



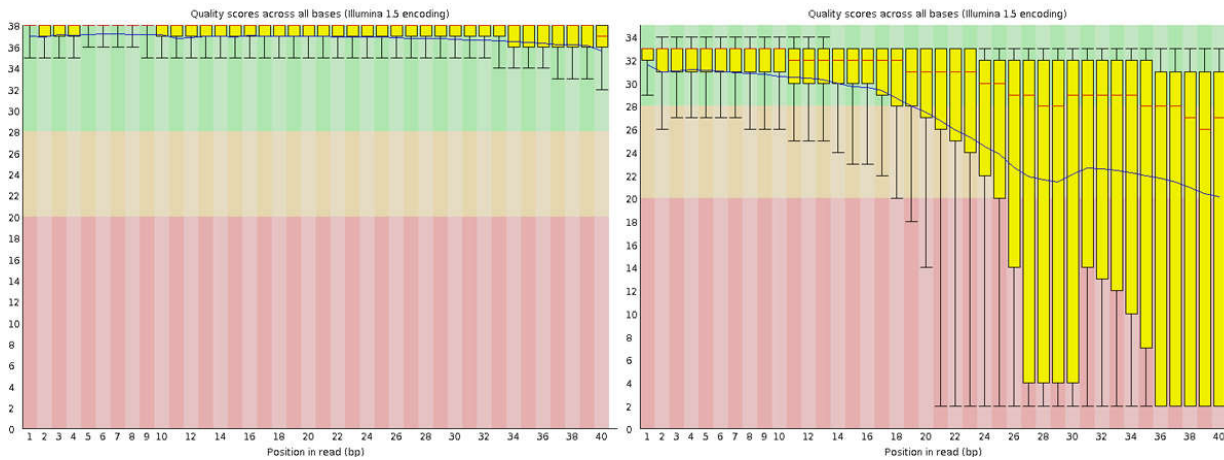
 FastQC, MultiQC

<https://sequencing.qcfail.com/>

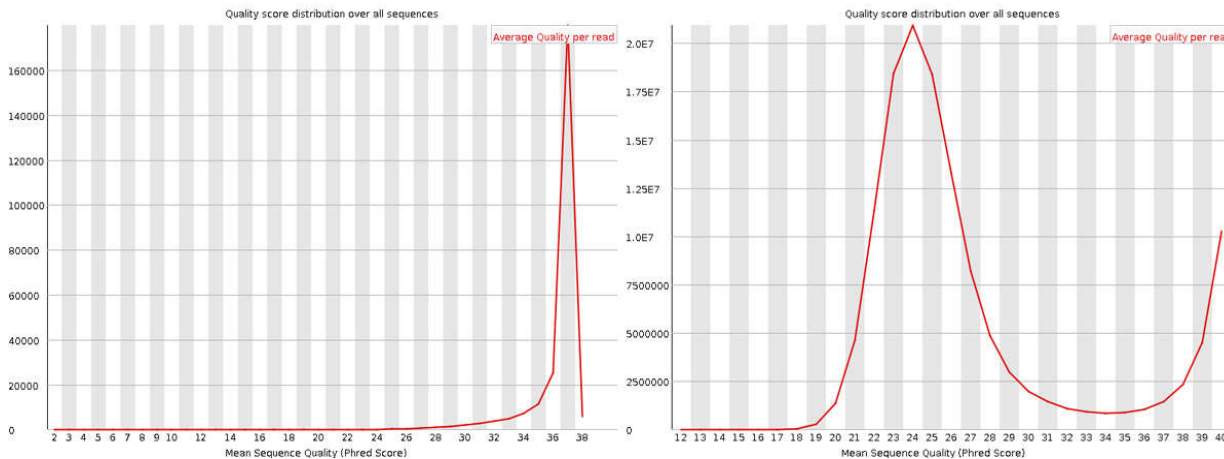
 **QCFAIL.com**

Articles about common next-generation  
sequencing problems

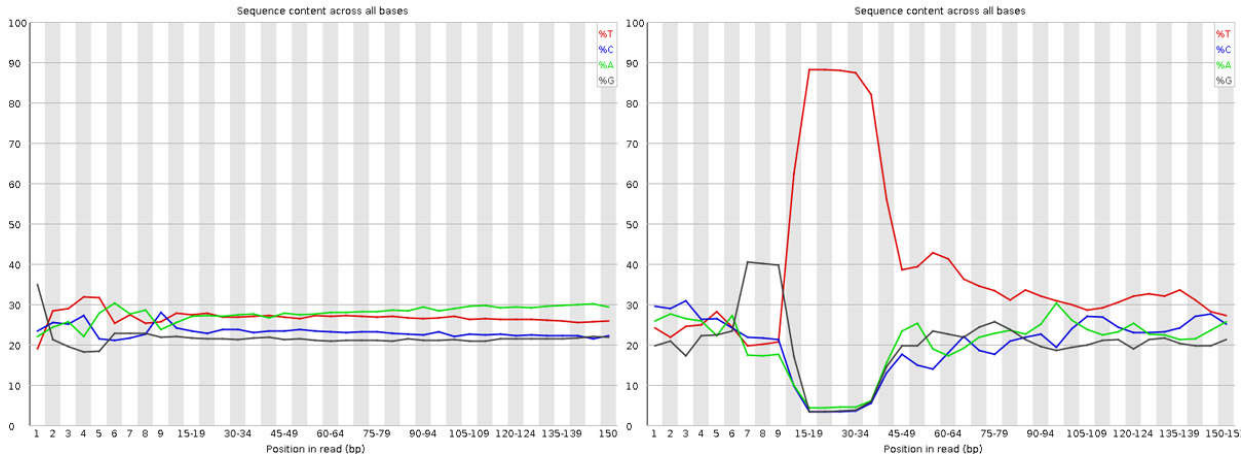
## Per base sequence quality



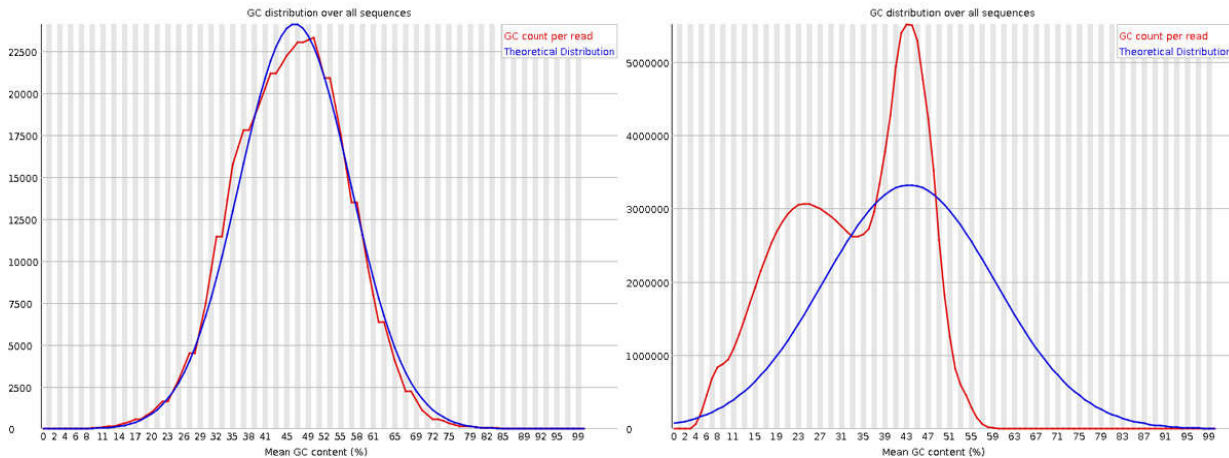
## Per sequence quality scores



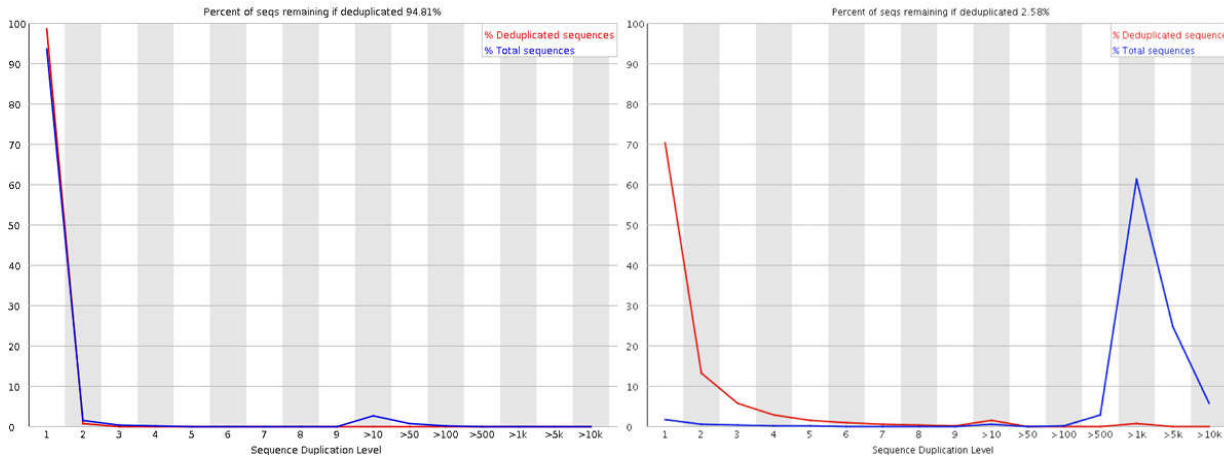
## Per base sequence content



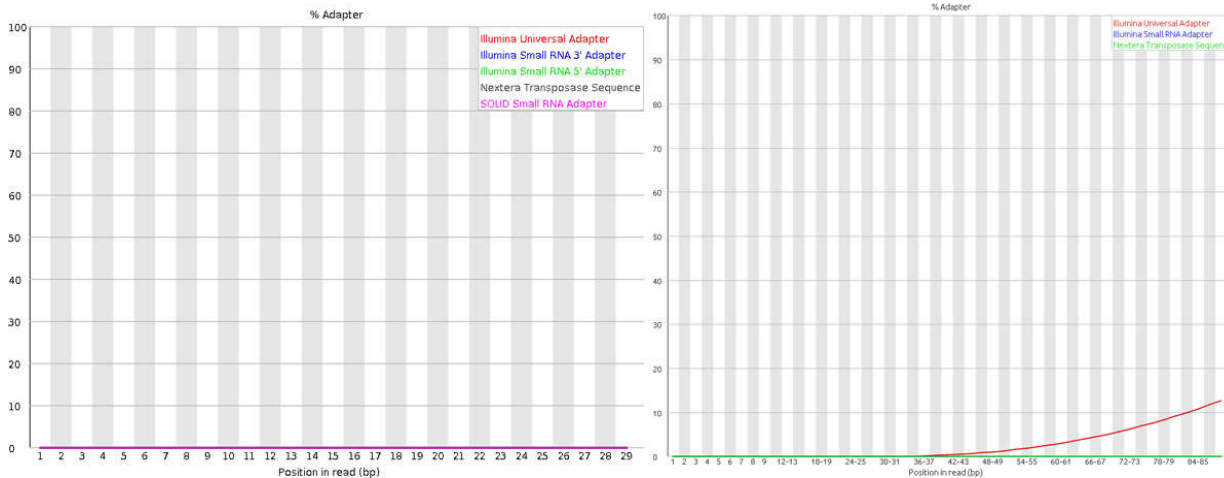
## Per sequence GC content



## Sequence duplication level



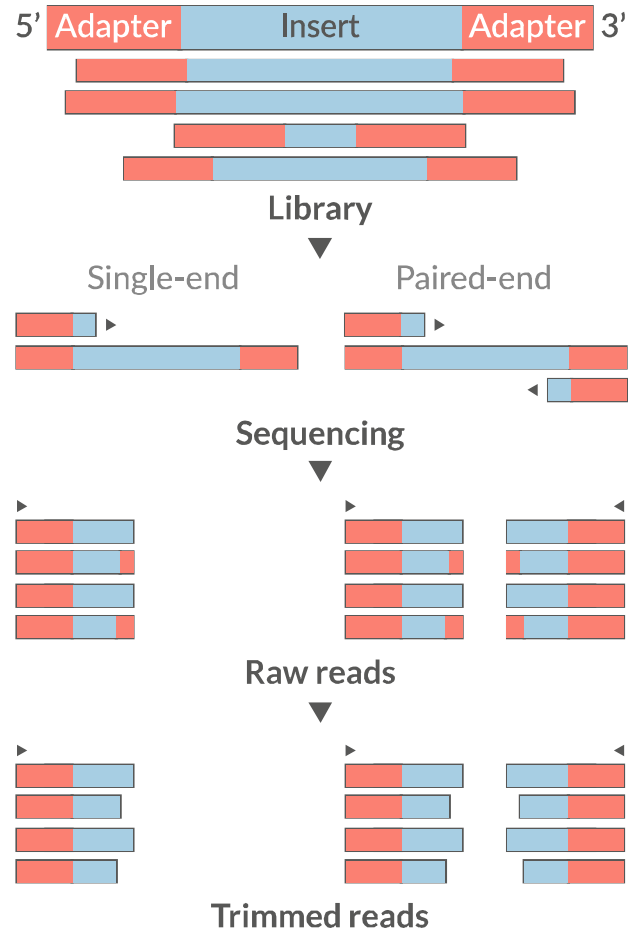
## Adapter content



# Trimming

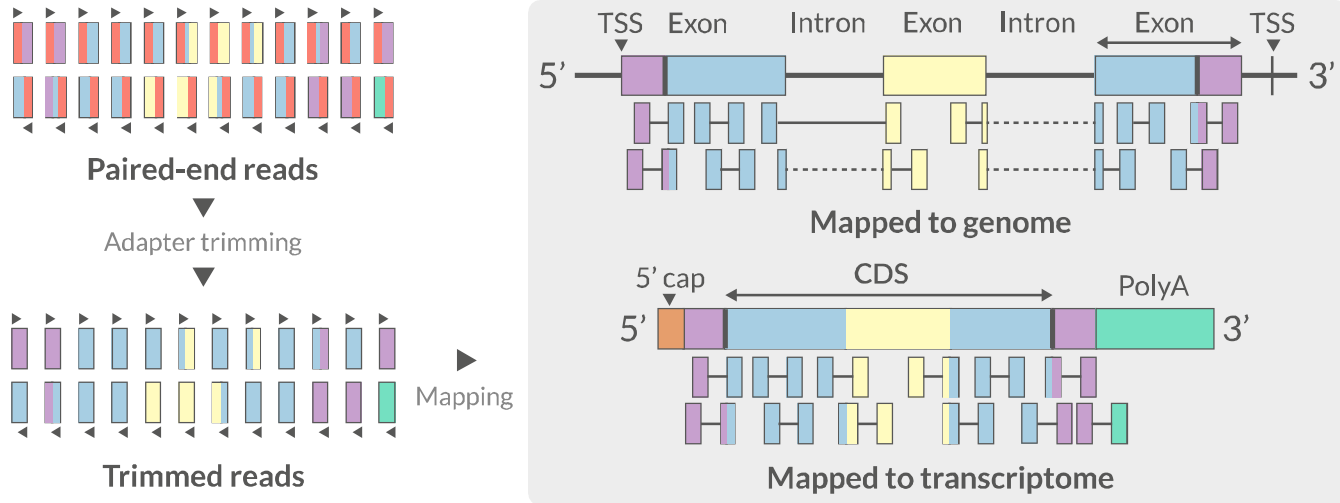
- Trim IF necessary
  - Synthetic bases can be an issue for SNP calling
  - Insert size distribution may be more important for assemblers
- Trim/Clip/Filter reads
- Remove adapter sequences
- Trim reads by quality
- Sliding window trimming
- Filter by min/max read length
  - Remove reads less than ~22nt
- Demultiplexing/Splitting

 [Cutadapt](#), [fastp](#), [Skewer](#), [Prinseq](#)





# Mapping



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

 **STAR, HiSat2, GSNAP, Novoalign** (Commercial)

- Reads (FASTQ)

```
@ST-E00274:179:HHYMLALXX:8:1101:1641:1309 1:N:0:NGATGT
NCATCGTGGTATTTGCACATCTTTTCTTATCAAATAAAAAGTTTAACTACTCAGTTATGCGCATACGTTTTTTGATGGCATTTC
+
#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
CAAATTAATTTAGCCAGAGGCGCACAAACATACGACCTCTAAAAAAGGTGCTGTAACATG
```

- Annotation (GTF/GFF)

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

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

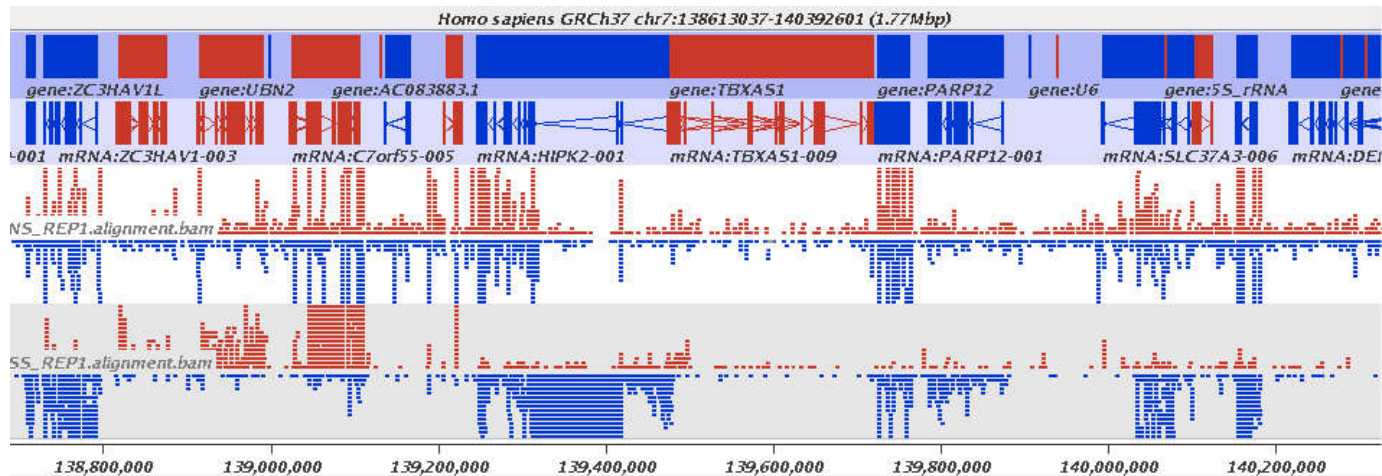
# Alignment

- 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

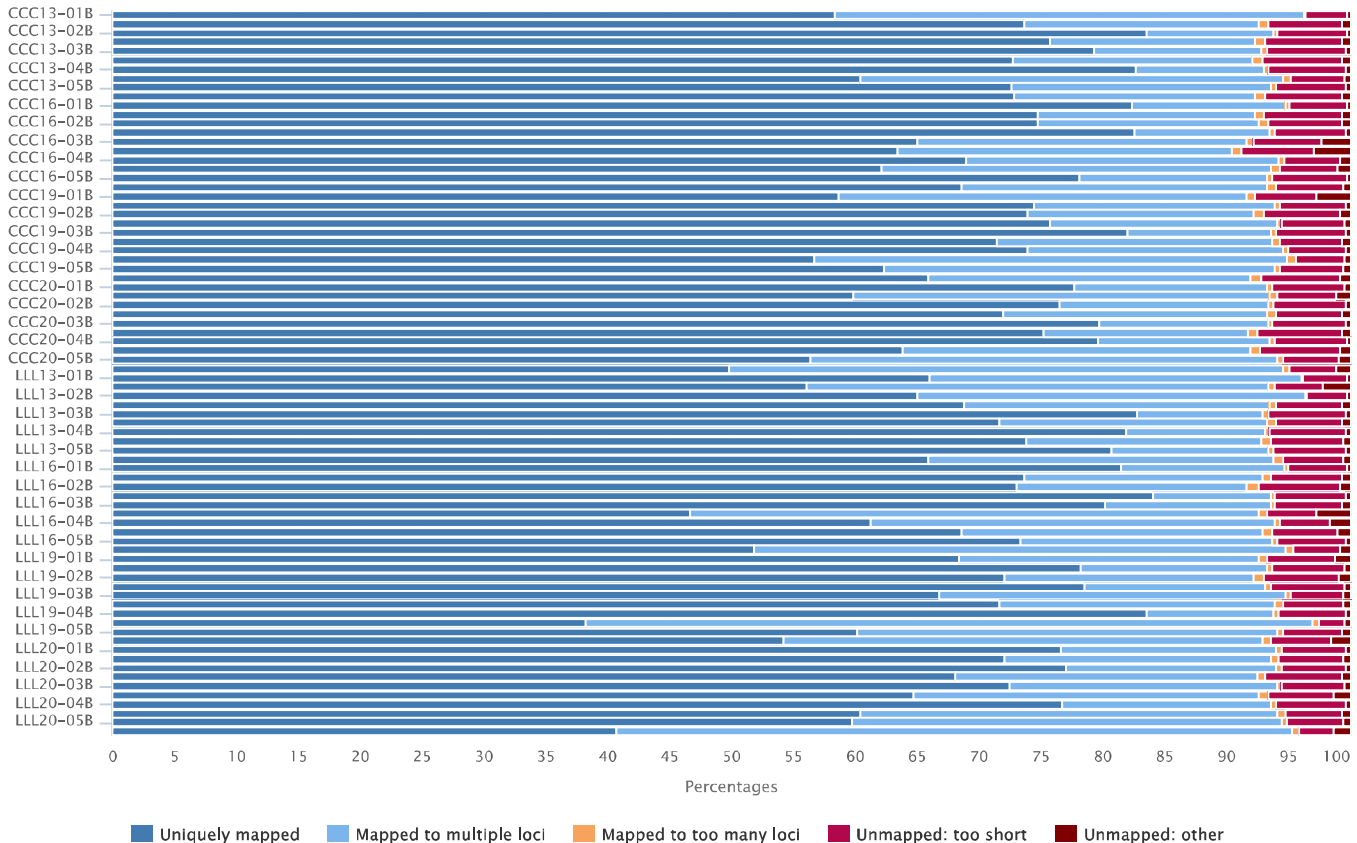
- Number of reads mapped/unmapped/paired etc
- Uniquely mapped
- Insert size distribution
- Gene body coverage
- Biotype counts / Chromosome counts
- Counts by region: gene/intron/non-genic

 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



```
samtools stats file.bam
```

```
SN      raw total sequences:    522095280
SN      filtered sequences:    0
SN      sequences:          522095280
SN      is sorted:          1
SN      1st fragments:      261047640
SN      last fragments:    261047640
SN      reads mapped:      514139025
SN      reads mapped and paired:    510035006
SN      reads unmapped:    7956255
SN      reads properly paired: 460249078
SN      reads paired:      522095280
SN      reads duplicated:    60151694
SN      reads MQ0:         54098384
SN      reads QC failed:    0
SN      non-primary alignments: 15023188
SN      total length:      78437013272
SN      bases mapped:      77238941462
SN      bases mapped (cigar): 74139898333
SN      bases trimmed:     0
SN      bases duplicated:    9022025650
SN      mismatches:        1695194781
SN      error rate:        2.286481e-02
SN      average length:    150
SN      maximum length:    151
SN      average quality:    37.6
```

```
...
```

```
bamtools stats file.bam
```

```
*****
```

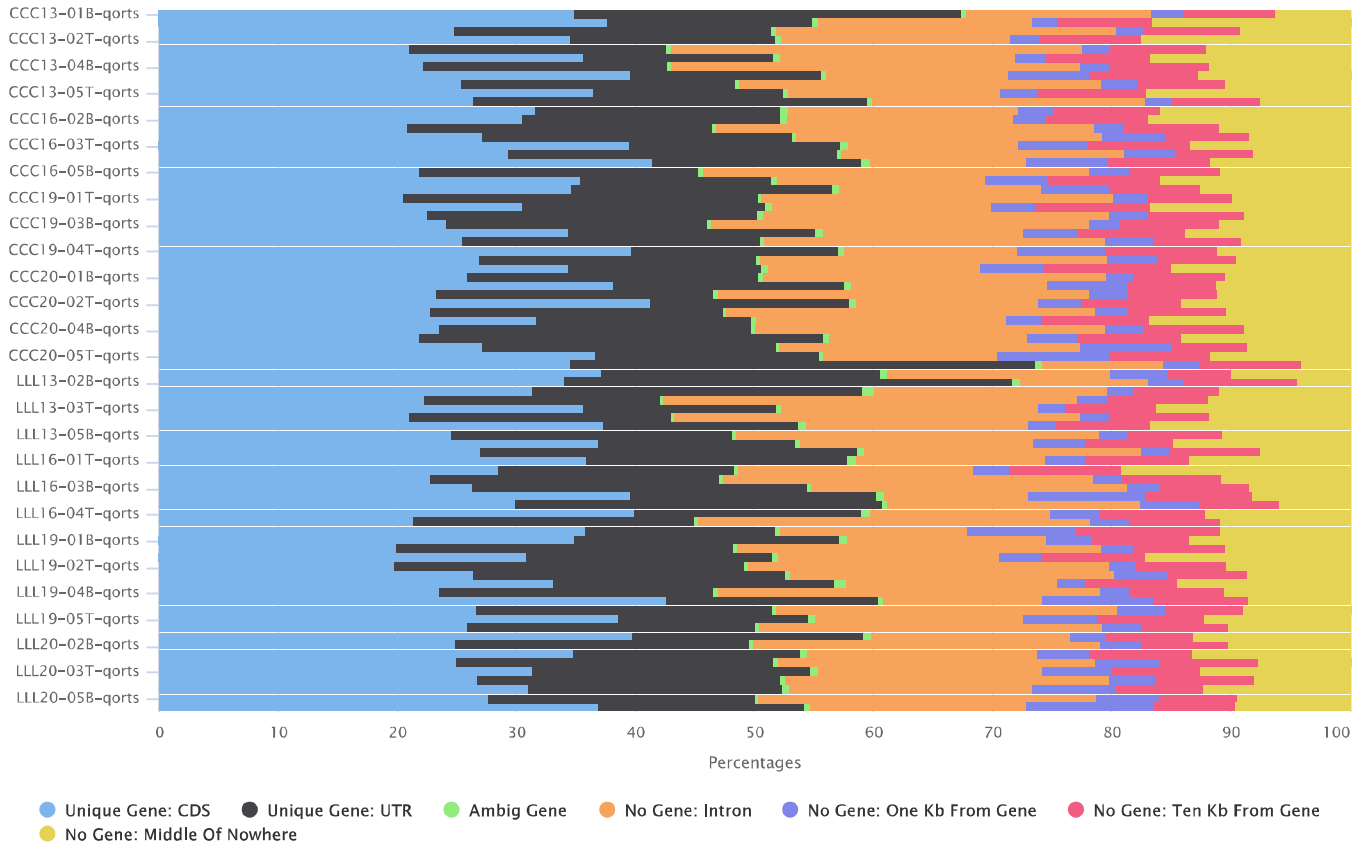
```
Stats for BAM file(s):
```

```
*****
```

Total reads:	537118468	
Mapped reads:	529162213	(98.5187%)
Forward strand:	270376825	(50.3384%)
Reverse strand:	266741643	(49.6616%)
Failed QC:	0	(0%)
Duplicates:	61425418	(11.4361%)
Paired-end reads:	537118468	(100%)
'Proper-pairs':	465991264	(86.7576%)
Both pairs mapped:	524501668	(97.651%)
Read 1:	268374707	
Read 2:	268743761	
Singletons:	4660545	(0.867694%)

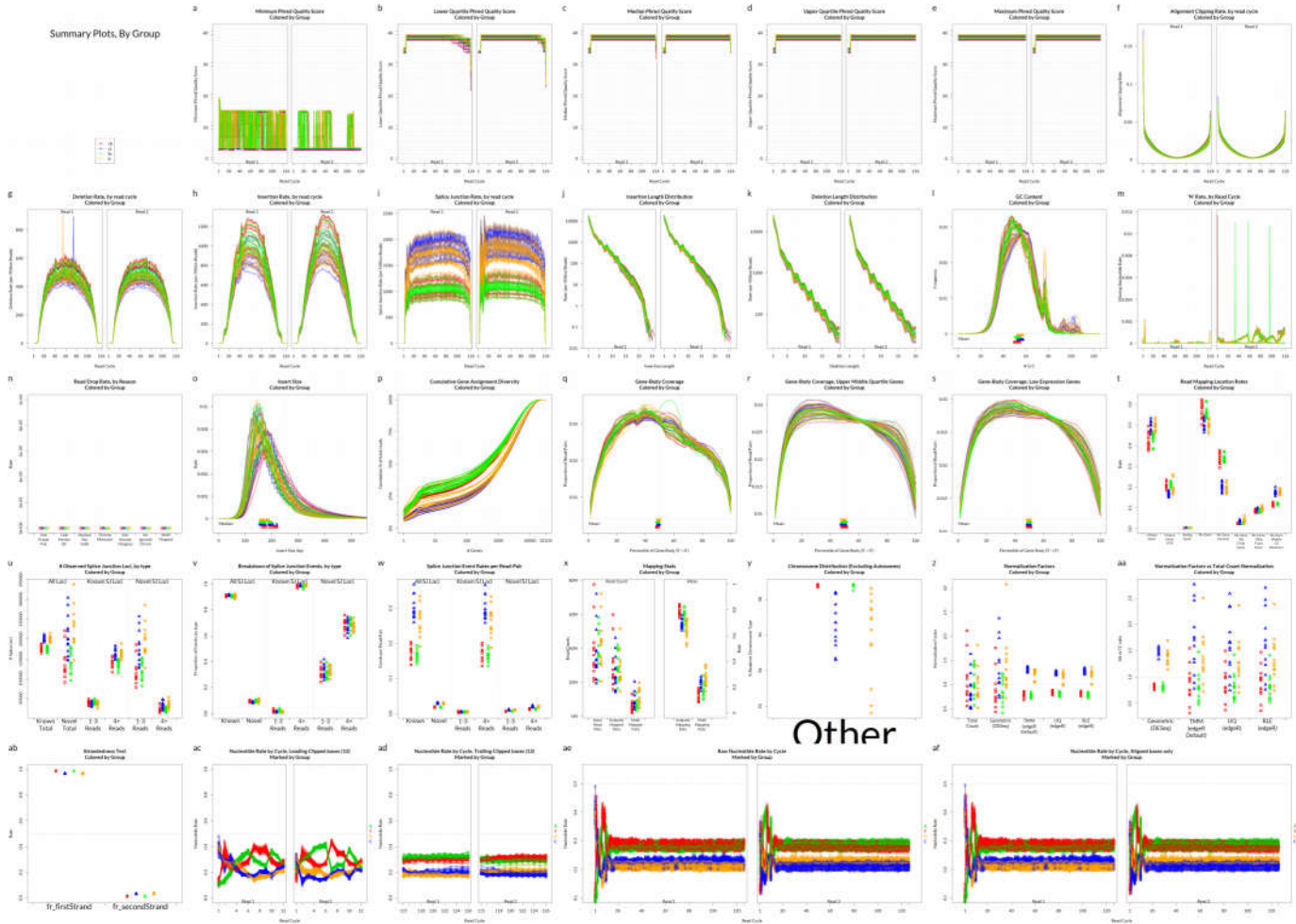
QoRTs was run on all samples and summarised using MultiQC.

QoRTs: Alignment Locations





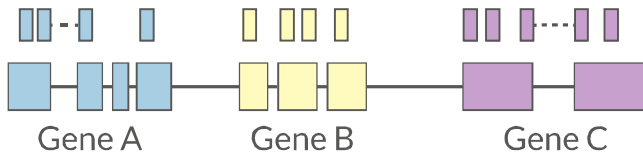
## Summary Plots, By Group



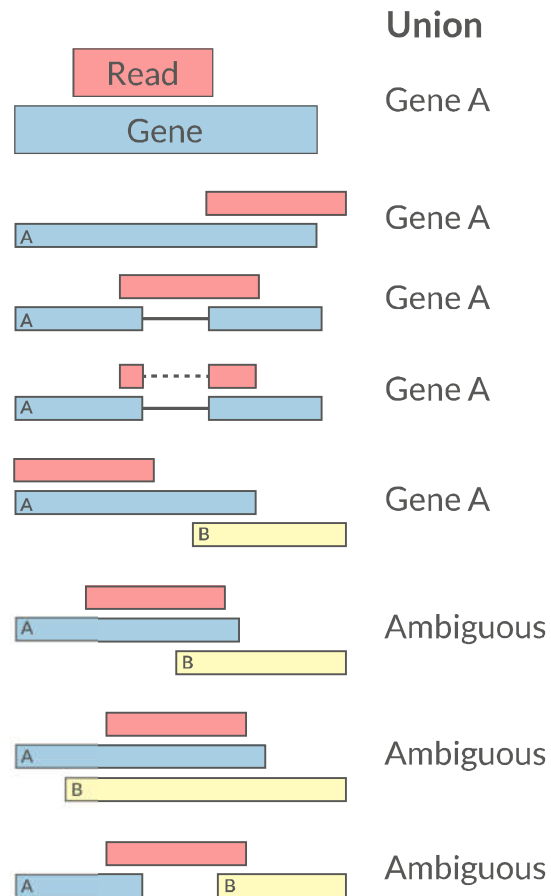
Other

# 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



# Quantification | PCR duplicates

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

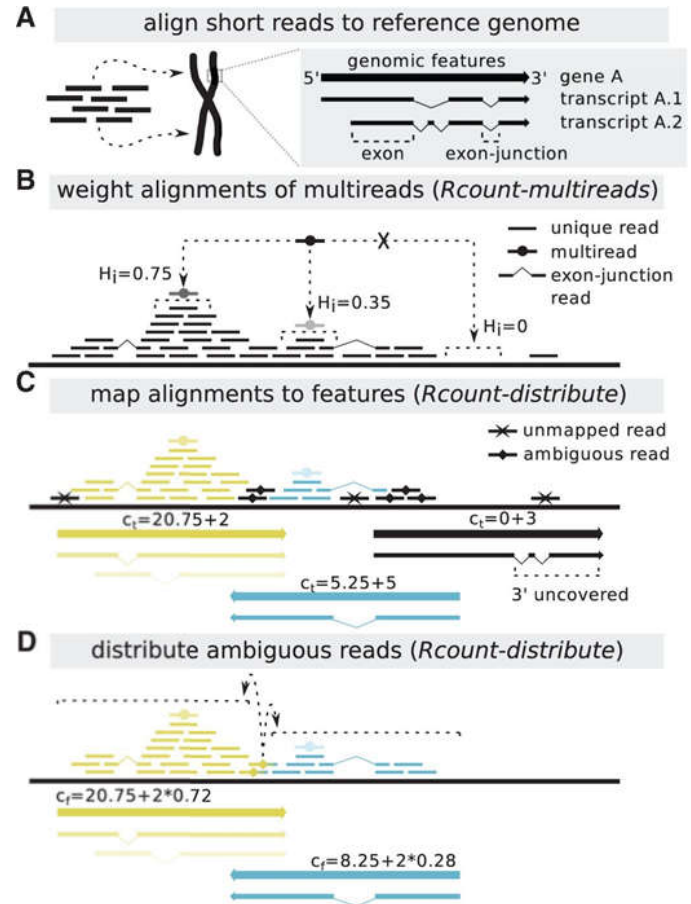
🔗 Fu, Yu, *et al.* "Elimination of PCR duplicates in RNA-seq and small RNA-seq using unique molecular identifiers." *BMC genomics* 19.1 (2018): 531

🔗 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 | Multi-mapping

- Added (BEDTools multicov)
- Discard (featureCounts, HTSeq)
- Distribute counts (Cufflinks)
- Rescue
  - Probabilistic assignment (Rcount, Cufflinks)
  - Prioritise features (Rcount)
  - Probabilistic assignment with EM (RSEM)



- Count methods
  - Provide no inference on isoforms
  - Cannot accurately measure fold change
- Probabilistic assignment
  - Deconvolute ambiguous mappings
  - Transcript-level
  - cDNA reference

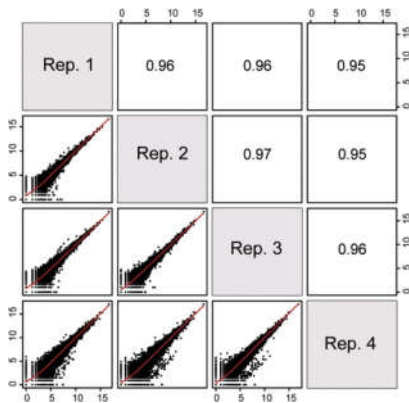
## Kallisto, Salmon

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

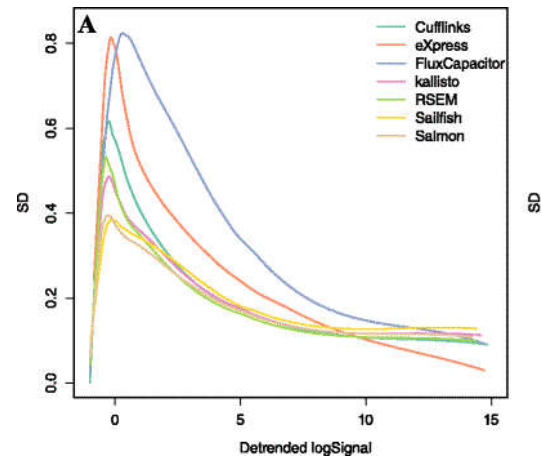
 RSEM, Kallisto, Salmon, Cufflinks2

ENSG00000000003	140	242	188	143	287	344	438	280	253
ENSG00000000005	0	0	0	0	0	0	0	0	0
ENSG000000000419	69	98	77	55	52	94	116	79	69
ENSG000000000457	56	75	104	79	157	205	183	178	153
ENSG000000000460	33	27	23	19	27	42	69	44	40
ENSG000000000938	7	38	13	17	35	76	53	37	24
ENSG000000000971	545	878	694	636	647	216	492	798	323
ENSG000000001036	79	154	74	80	128	167	220	147	72

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

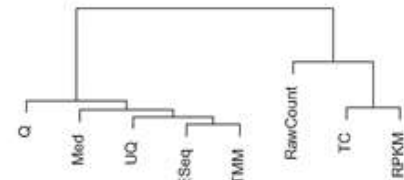
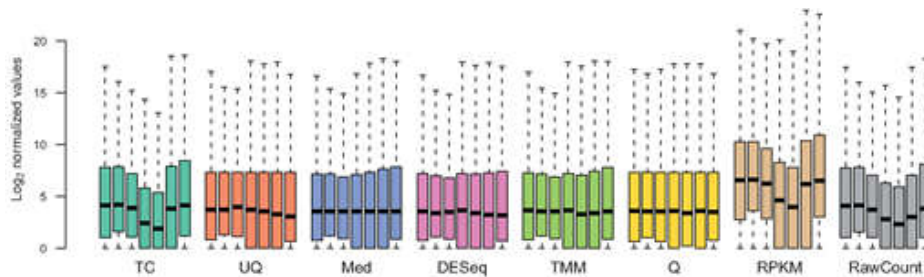


- Count QC using RNASeqComp

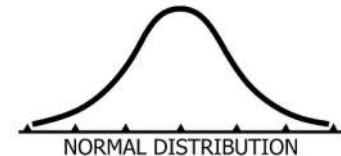


# Normalisation

- 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



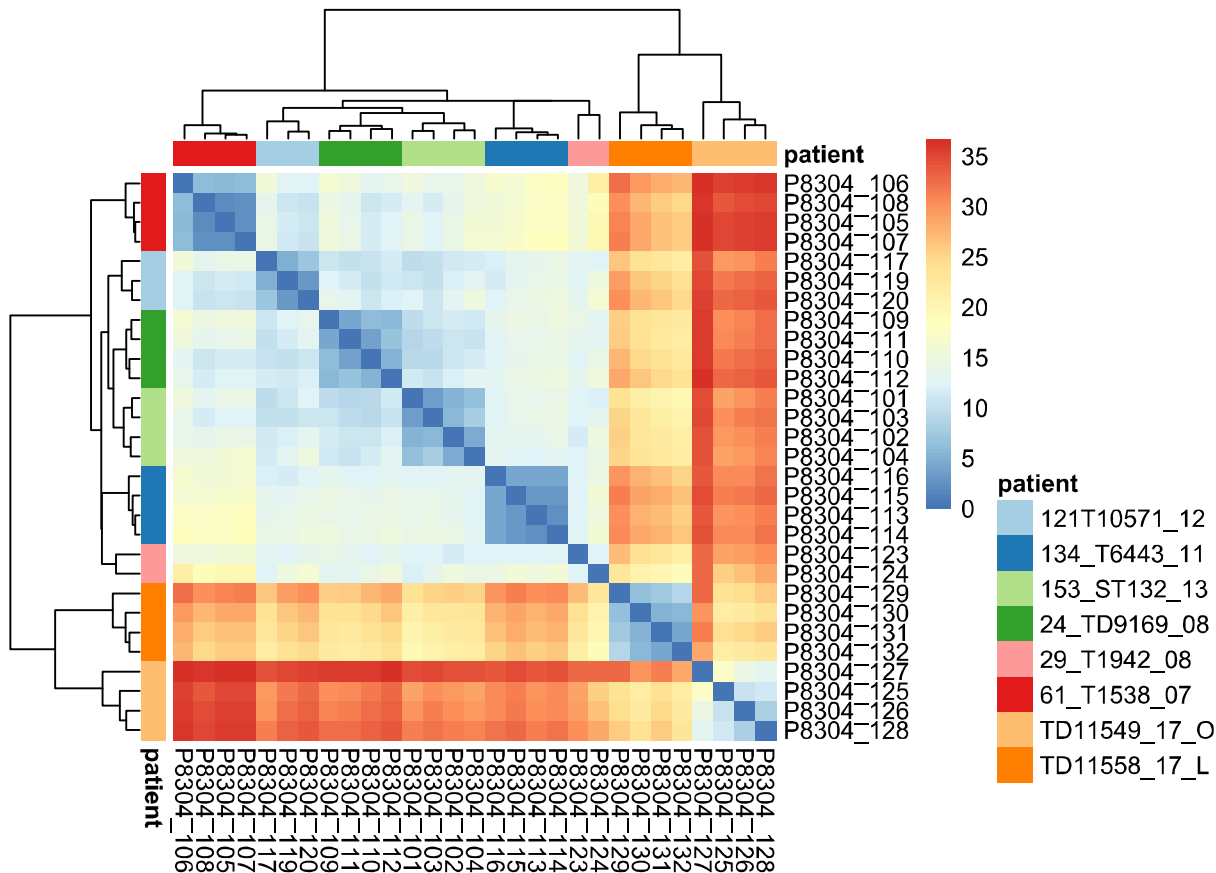
🔗 Dillies, Marie-Agnes, *et al.* "A comprehensive evaluation of normalization methods for Illumina high-throughput RNA sequencing data analysis." [Briefings in bioinformatics](#) 14.6 (2013): 671-683

🔗 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






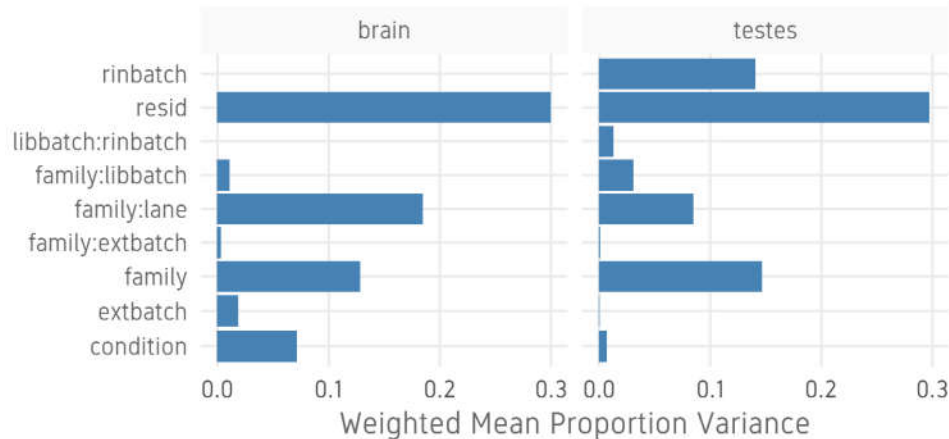
# Exploratory | MDS

NBIS 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`

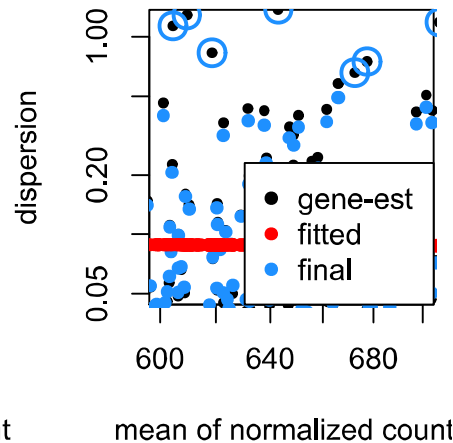
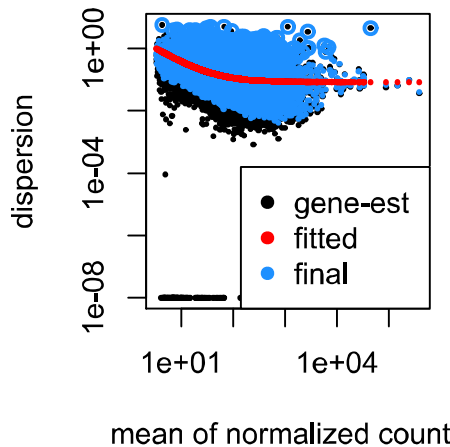
- Estimate variation explained by variables (PVCA)



- 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

- 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

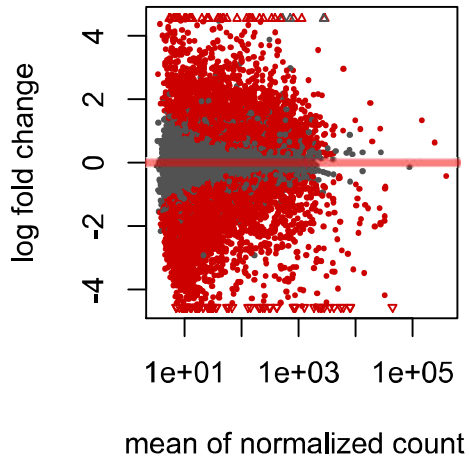
- Results `results()`

```
## 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
```

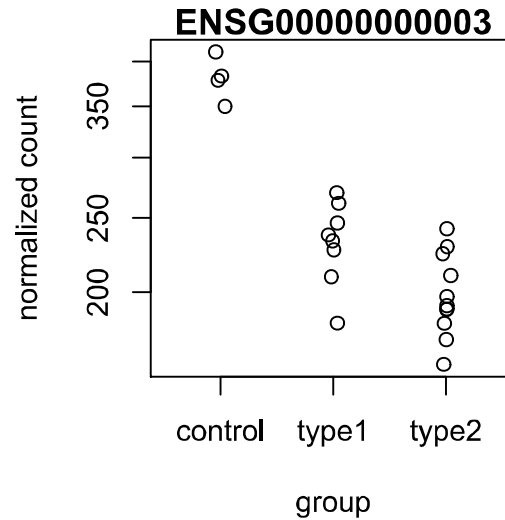
- Summary `summary()`

```
##
## 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
```

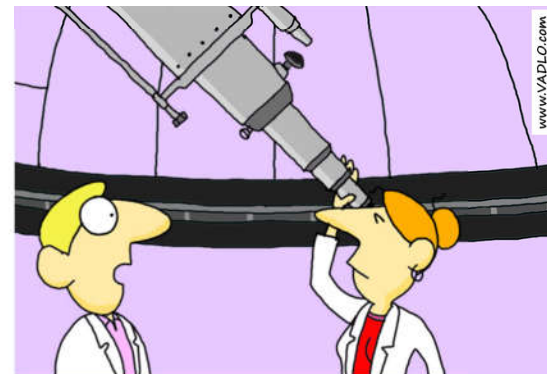
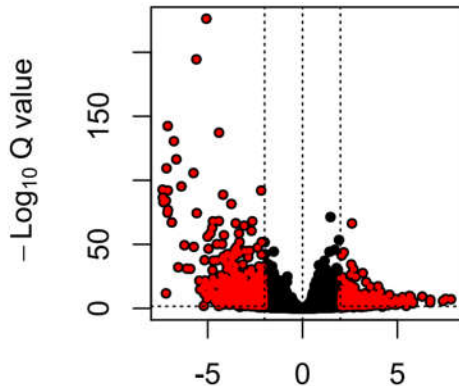
- MA plot `plotMA()`



- Normalised counts `plotCounts()`



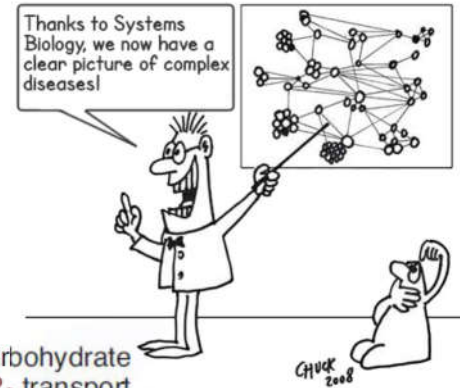
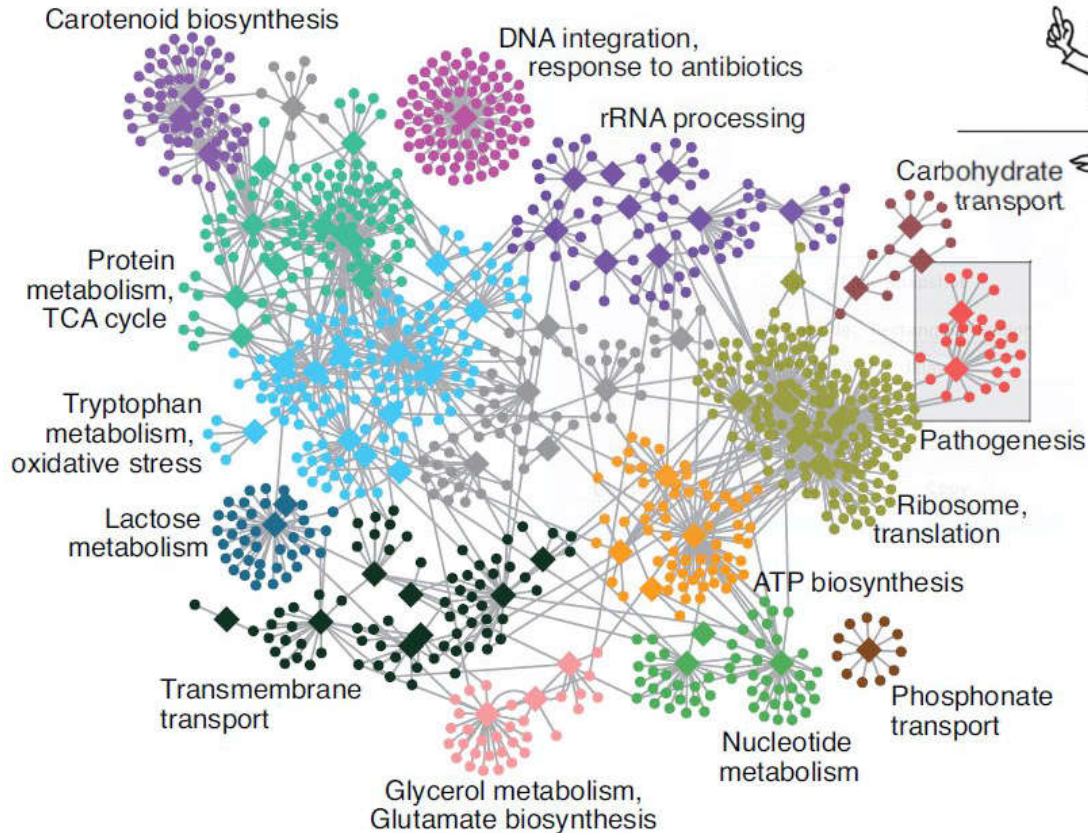
- Volcano plot



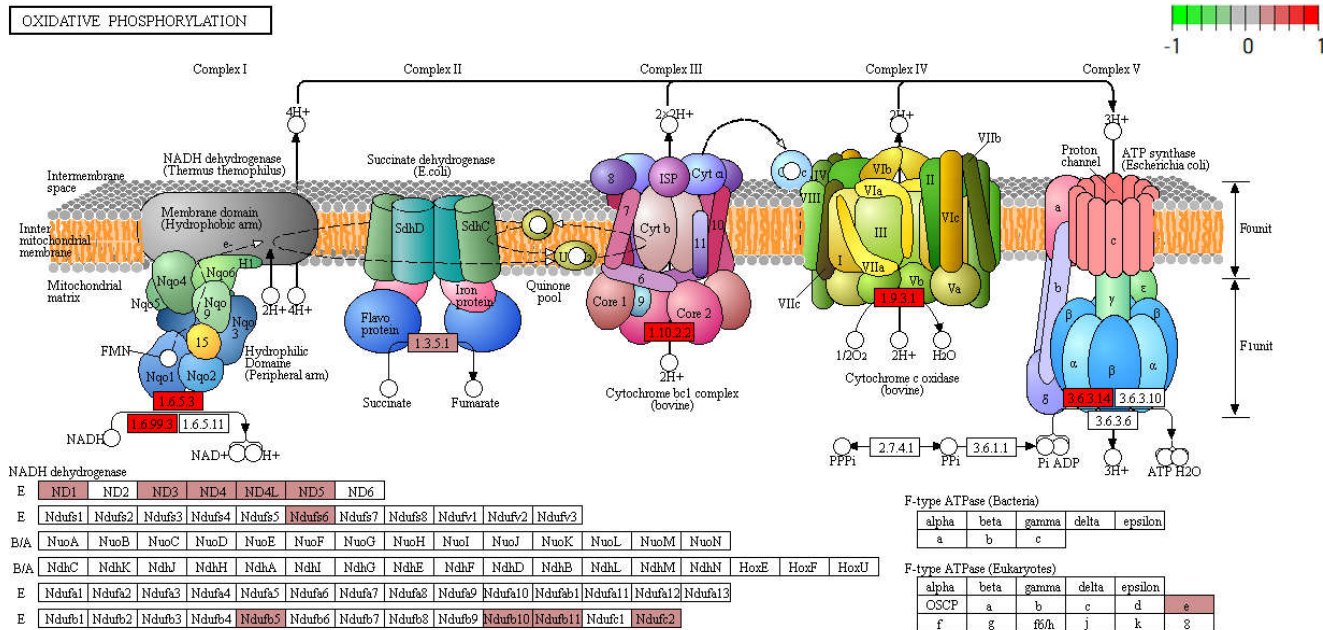
"Can you see the upper points of my scatter plot?"

# Functional analysis | GO

- Gene enrichment analysis
- Gene set enrichment analysis (GSEA)
- Gene ontology / Reactome databases

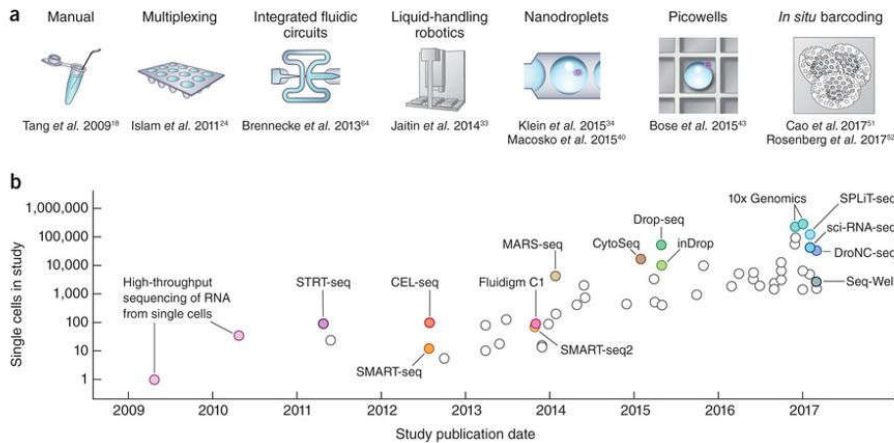


- Pathway analysis (Kegg)



DAVID, clusterProfiler, ClueGO, ErmineJ, pathview

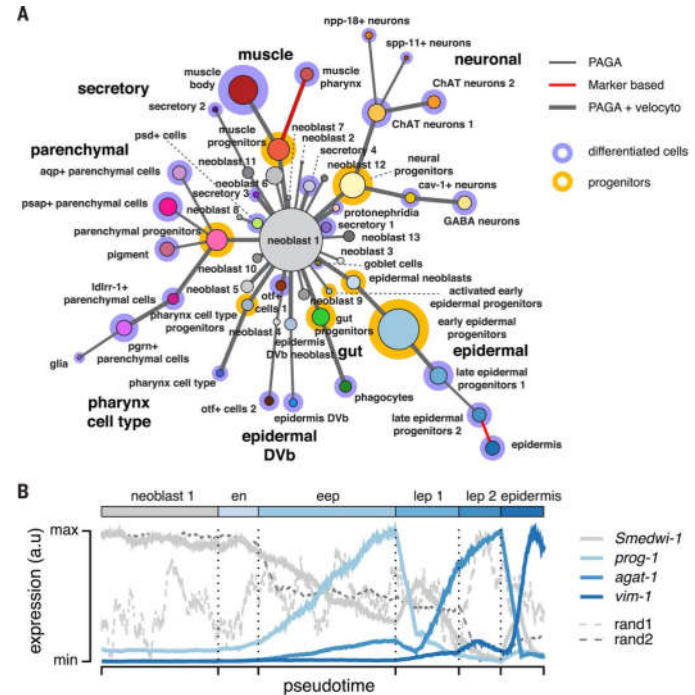
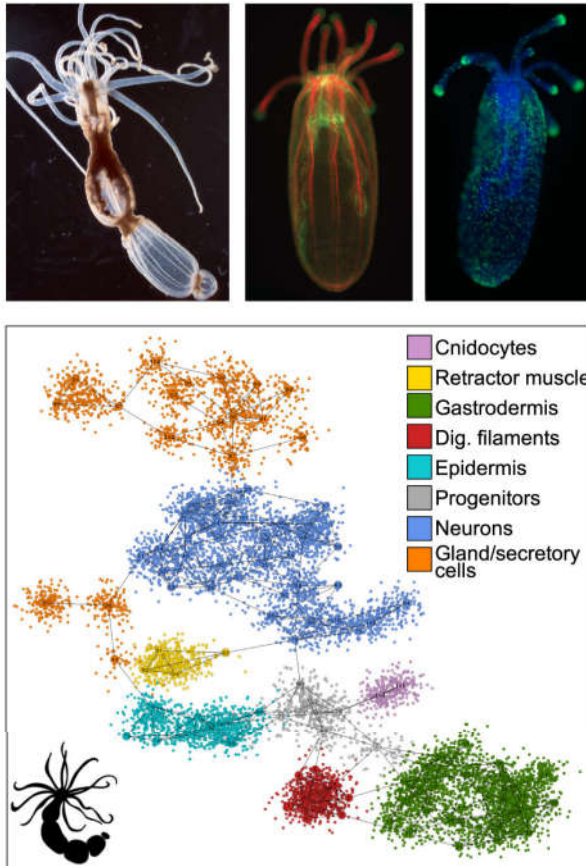
- Bulk RNA-Seq measures mean expression-level over many cells
- Poor resolution for development, differentiation, heterogenous tissues
- Identify cell types in a tissue
- Temporal/spatial/conditional change in cellular state and composition



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



# scRNA-Seq | Example



- 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 (ScISO-Seq) across thousands of cells reveals isoforms of cerebellar cell types.

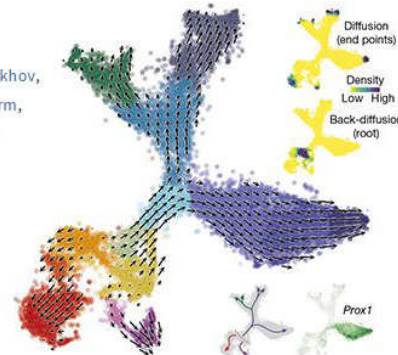
Ishaan Gupta, Paul G Collier, Bettina Haase, Ahmed Mahfouz, Anoushka Joglekar, Taylor Floyd, Frank Koopmans, Ben Barres, August B Smit, Steven Sloan, Wenjie Luo, Olivier Fedrigo, M Elizabeth Ross, Hagen U Tilgner

- 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. Kastriiti, 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](#)



- Nothing can fix a poor experimental design
- 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
- QC! QC everything at every step
- Verify that tools and methods align with data assumptions
- Experiment with multiple pipelines and tools

 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:  12-Sep-2018 at 21:33:47

---

2018 Roy Francis | [SciLifeLab](#) | [NBIS](#)

# Session

This presentation was created in RStudio using `remarkjs` framework through R package `xaringan`.

```
getS3method("print", "sessionInfo")(sessionInfo()[ -7])
```

```
## R version 3.5.1 (2018-07-02)
## Platform: x86_64-w64-mingw32/x64 (64-bit)
## Running under: Windows >= 8 x64 (build 9200)
##
## Matrix products: default
##
## Locale:
## [1] LC_COLLATE=English_United Kingdom.1252
## [2] LC_CTYPE=English_United Kingdom.1252
## [3] LC_MONETARY=English_United Kingdom.1252
## [4] LC_NUMERIC=C
## [5] LC_TIME=English_United Kingdom.1252
##
## attached base packages:
## [1] parallel stats4 stats graphics grDevices utils datasets
## [8] methods base
##
## other attached packages:
## [1] bindrcpp_0.2.2 DESeq2_1.20.0
## [3] SummarizedExperiment_1.10.1 DelayedArray_0.6.5
## [5] BiocParallel_1.14.2 matrixStats_0.54.0
## [7] Biobase_2.40.0 GenomicRanges_1.32.6
## [9] GenomeInfoDb_1.16.0 IRanges_2.14.11
## [11] S4Vectors_0.18.3 BiocGenerics_0.26.0
## [13] plotly_4.8.0 ggplot2_3.0.0
## [15] pheatmap_1.0.10 dplyr_0.7.6
## [17] bookdown_0.7 knitr_1.20
```