

RNA-seq to study HIV Infection in cells

Rebecca Batorsky

Sr Data Scientist

Tufts Data Intensive Studies Center

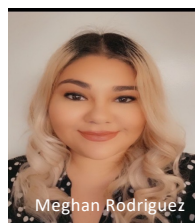
Feb 2024

People at DISC

DISC Faculty



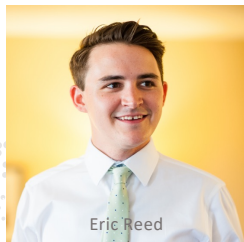
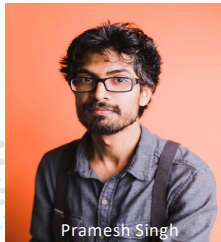
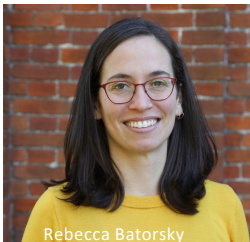
Administration



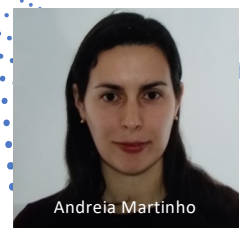
Data Intensive
Studies Center

40+ faculty and **many students**
across different disciplines
partner with us on research,
teaching and learning

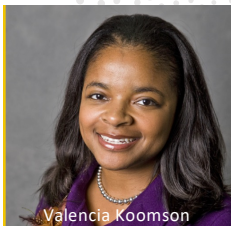
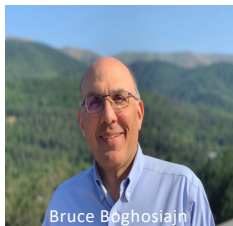
DISC Data Scientists



Postdoctoral



DISC Faculty Fellows



People at DISC

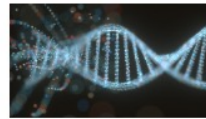


Data Intensive
Studies Center

DISC Faculty



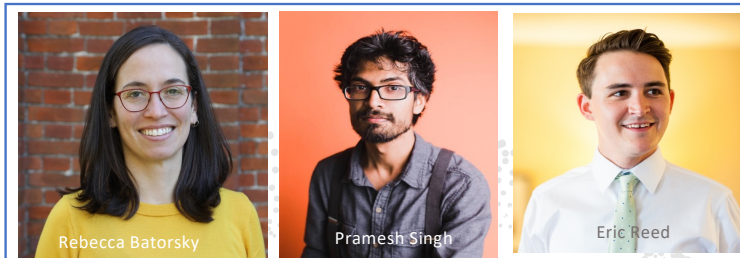
Bioinformatics & Computational Biology



DISC conducts research by developing new data science methods for bioinformatics and computational biology research, in collaboration with faculty, staff, and stakeholders within and outside Tufts University. Some of our current work focuses on:

- Single-cell Transcriptomics (Single cell profiling of Hofbauer cells and fetal brain microglia)
- Biological Networks (System-Level analysis of 'omics data)
- Proteomics (Proteomics profiling to study longevity)

DISC Data Scientists

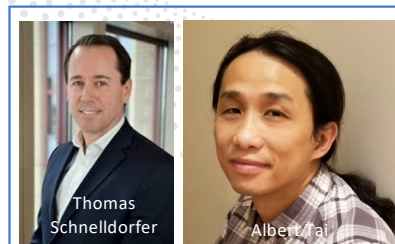


<https://disc.tufts.edu/disc-research/research-projects>

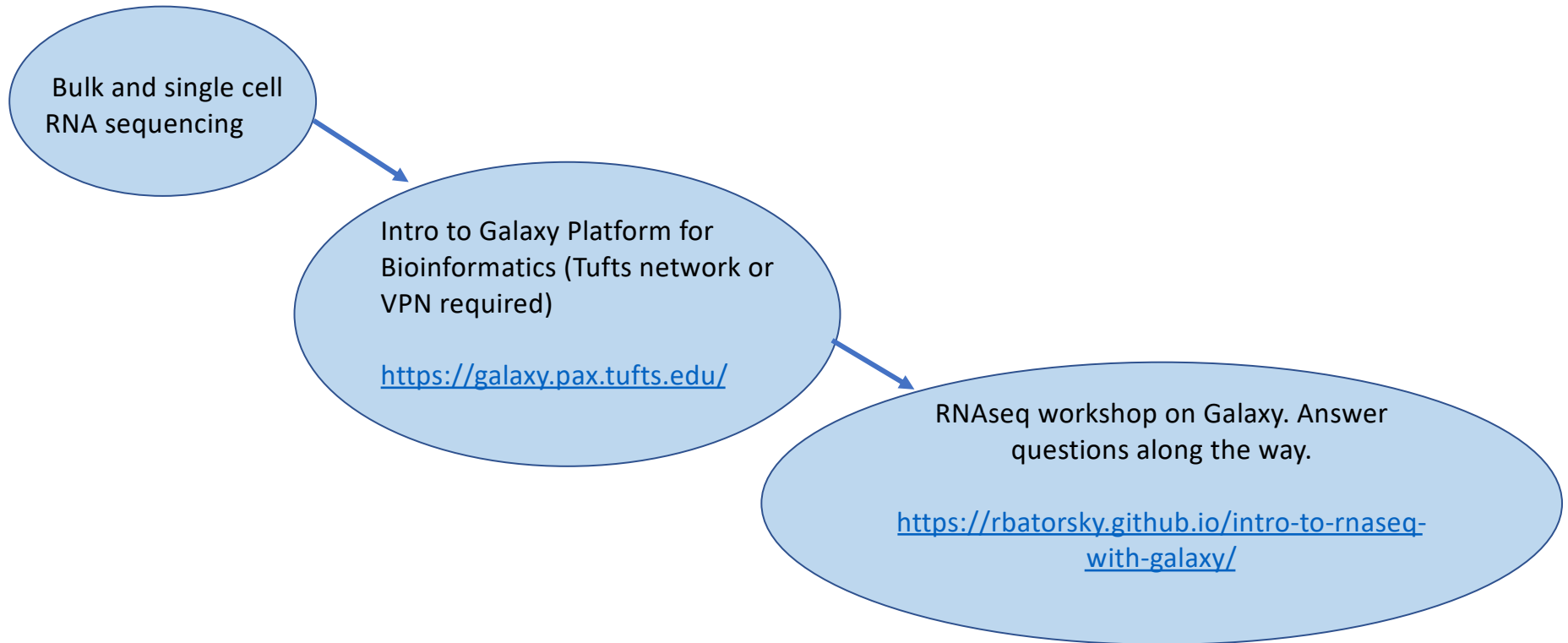
Postdoctoral



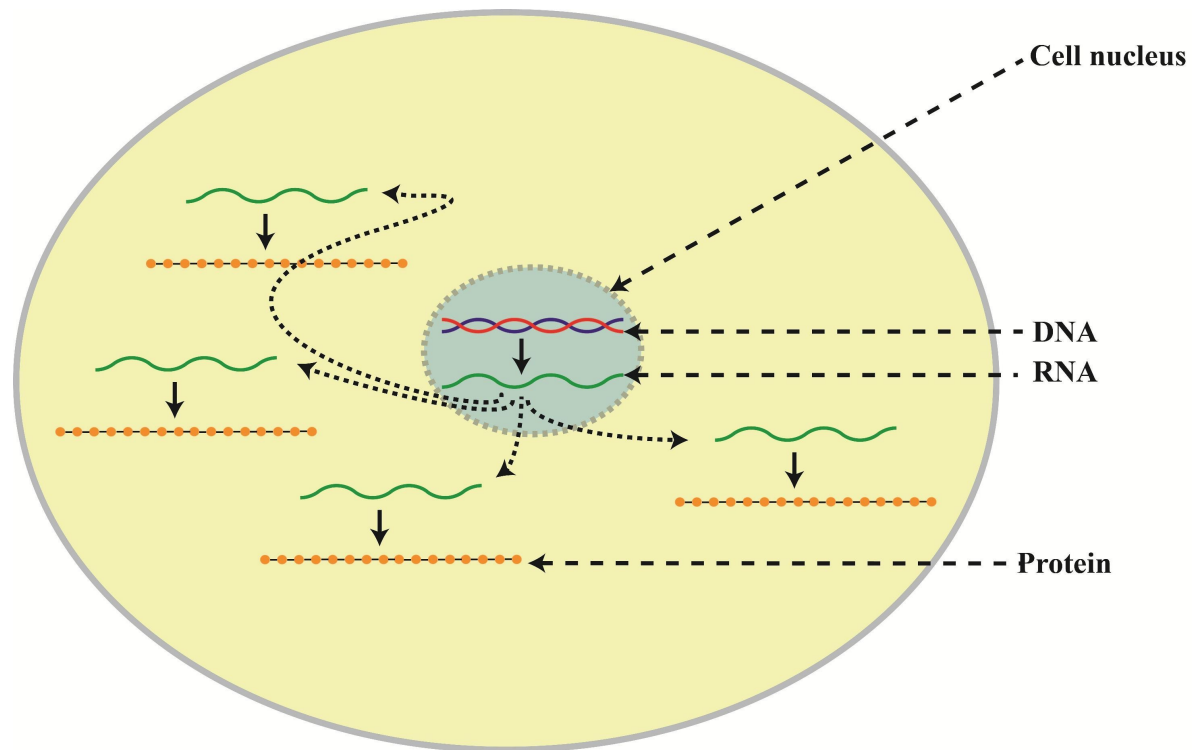
DISC Faculty Fellows



Outline



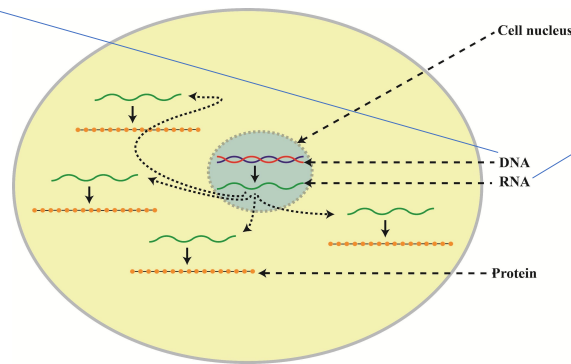
DNA and RNA in a cell



Two common analyses

DNA Sequencing

- Fixed number of copies of a gene per cell
- Analysis goal:
Variant calling and interpretation



RNA Sequencing

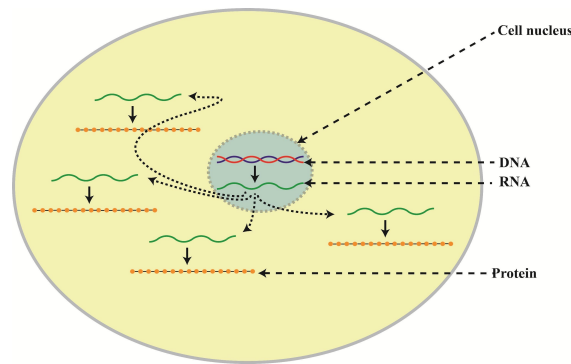
- Number of copies of a gene transcript per cell depends on gene expression
- Analysis goal:
 - Bulk : Differential expression
 - Single cell : Quantify different cell populations

<https://i0.wp.com/science-explained.com/wp-content/uploads/2013/08/Cell.jpg>

Today we will cover RNA sequencing

DNA Sequencing

- Fixed number of copies of a gene per cell
- Analysis goal: Variant calling and interpretation



RNA Sequencing

- Number of copies of a gene transcript per cell depends on gene expression
- Analysis goal:
 - Bulk : Differential expression
 - Single cell : Quantify different cell populations

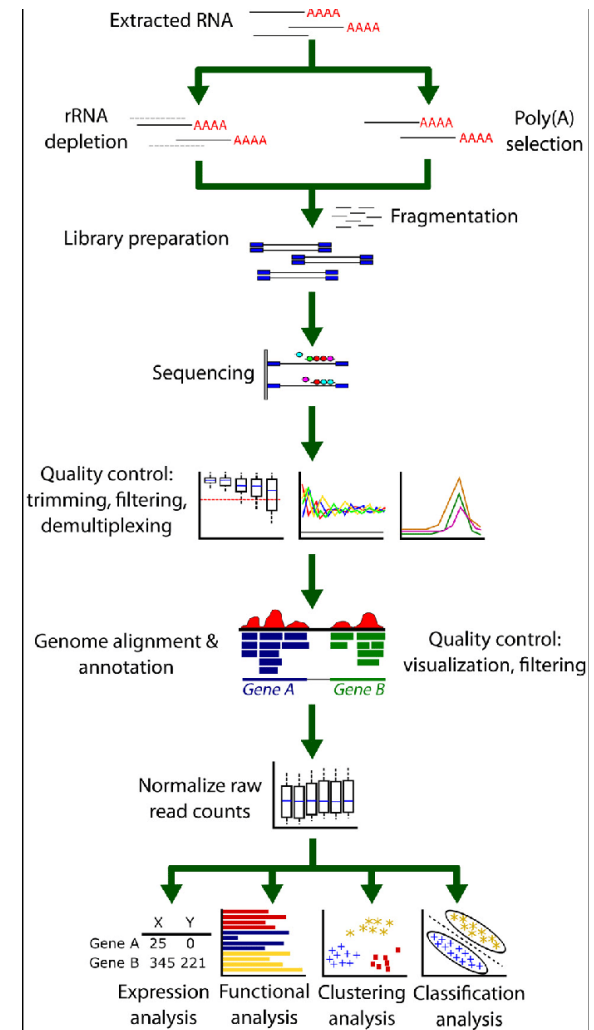
<https://i0.wp.com/science-explained.com/wp-content/uploads/2013/08/Cell.jpg>

“Bulk” RNA seq workflow

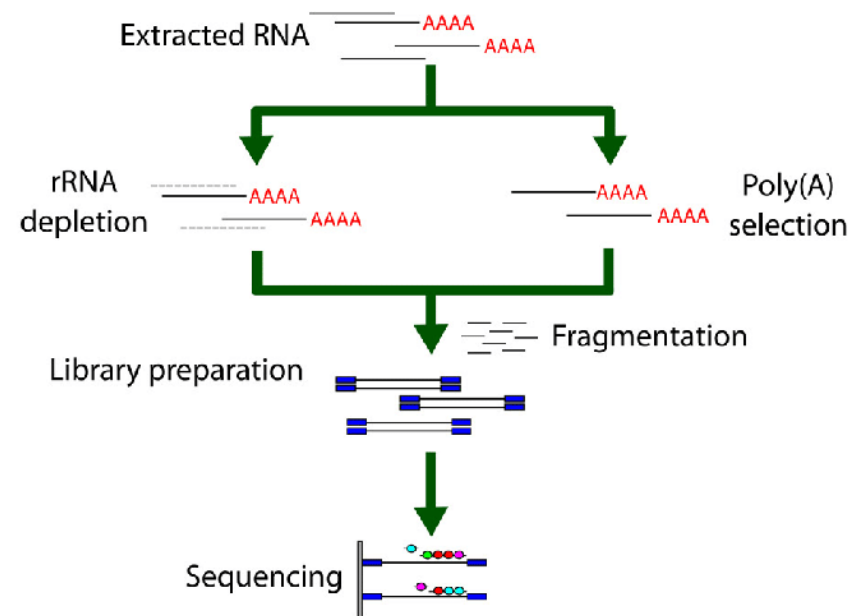
Library prep and sequencing

Bioinformatics

Good resource: [Griffiths et al Plos Comp Bio 2015](#)



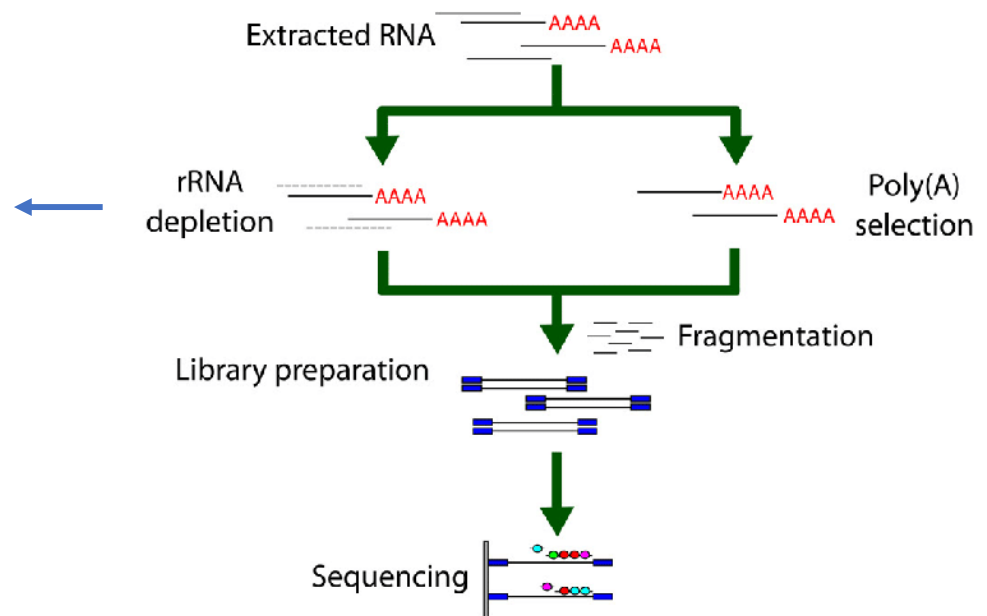
RNA seq library prep and sequencing



Good resource: [Griffiths et al Plos Comp Bio 2015](#)

RNA seq library prep and sequencing

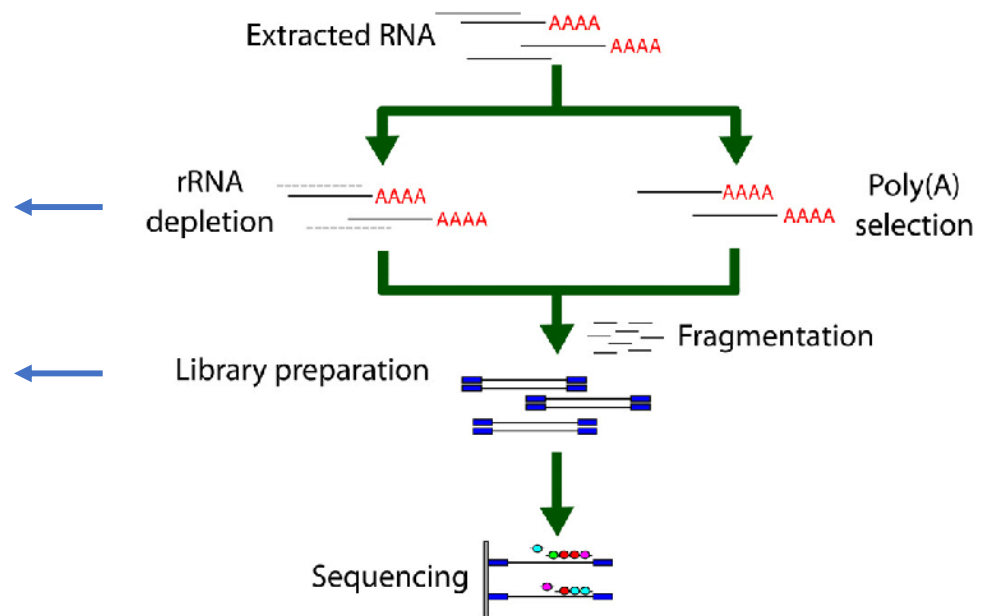
- Enrichment for mRNA, two options
- In humans, ~95%–98% of all RNA molecules are rRNAs



Good resource: [Griffiths et al Plos Comp Bio 2015](#)

RNA seq library prep and sequencing

- Enrichment for mRNA, two options
- In humans, ~95%–98% of all RNA molecules are rRNAs
- Random priming and reverse transcription
- Double stranded cDNA synthesis
- Sequencing adapter ligation



Resources:

[Illumina Sequencing by Synthesis](#)

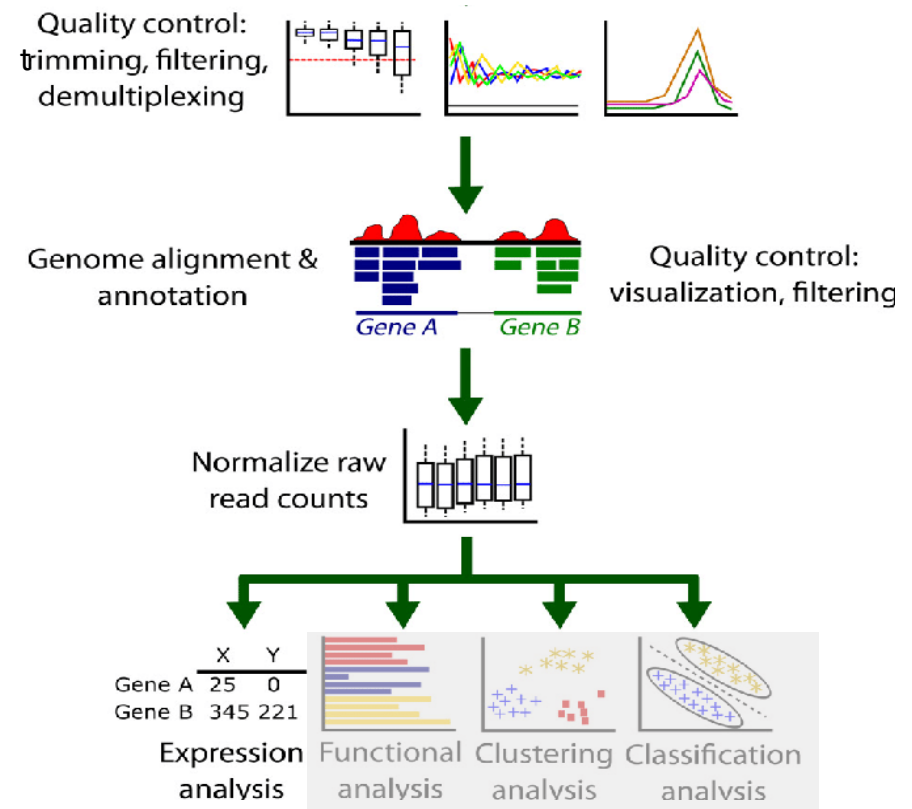
[Griffiths et al Plos Comp Bio 2015](#)

RNA seq bioinformatics

Goal of Differential Expression

“How can we detect genes for which the counts of reads change between conditions **more systematically** than as expected by chance”

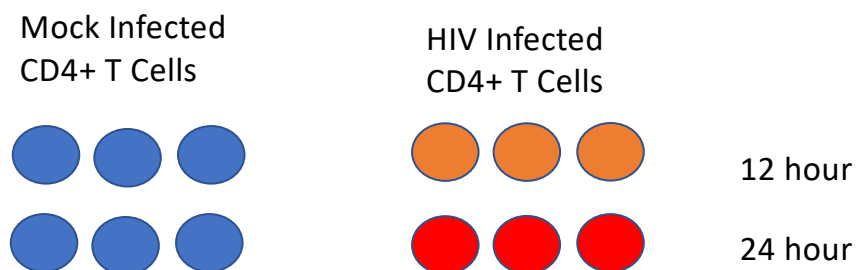
Oshlack et al. 2010. From RNA-seq reads to differential expression results. Genome Biology 2010, 11:220



Our dataset

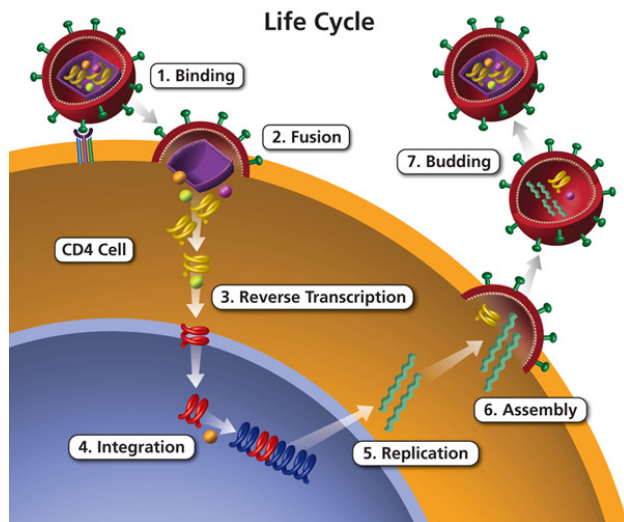
Next-Generation Sequencing Reveals HIV-1-Mediated Suppression of T Cell Activation and RNA Processing and Regulation of Noncoding RNA Expression in a CD4⁺ T Cell Line

Stewart T. Chang, Pavel Sova, Xinxia Peng, Jeffrey Weiss, G. Lynn Law, Robert E. Palermo, Michael G. Katze



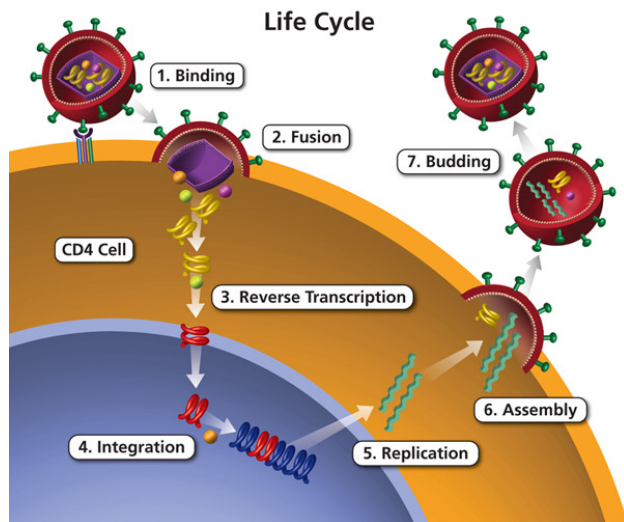
<https://www.ncbi.nlm.nih.gov/pubmed/21933919>

HIV lifecycle

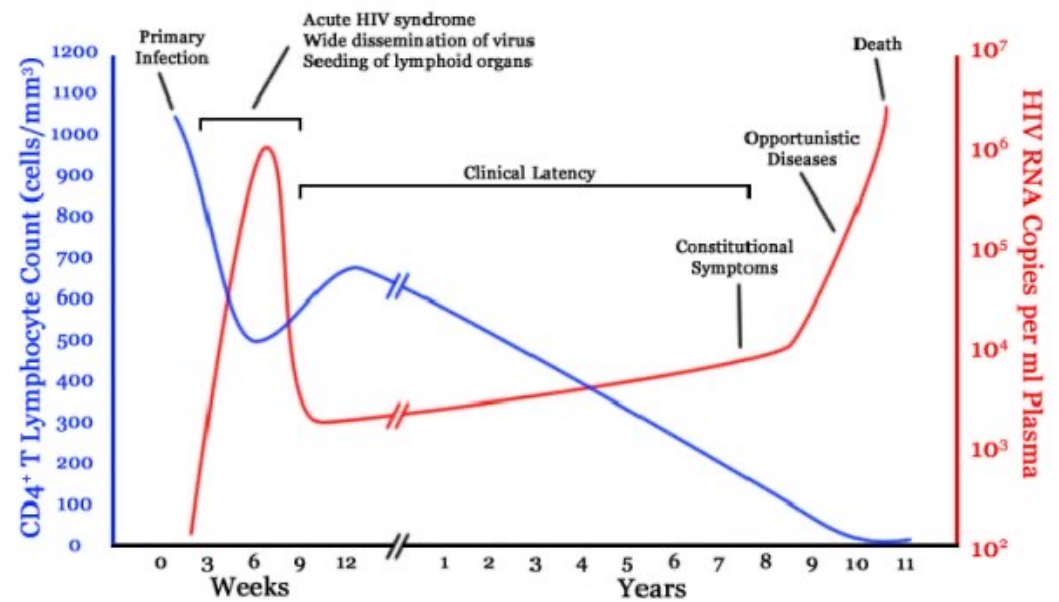


<https://aidsinfo.nih.gov/understanding-hiv-aids/glossary/1596/life-cycle>

HIV lifecycle



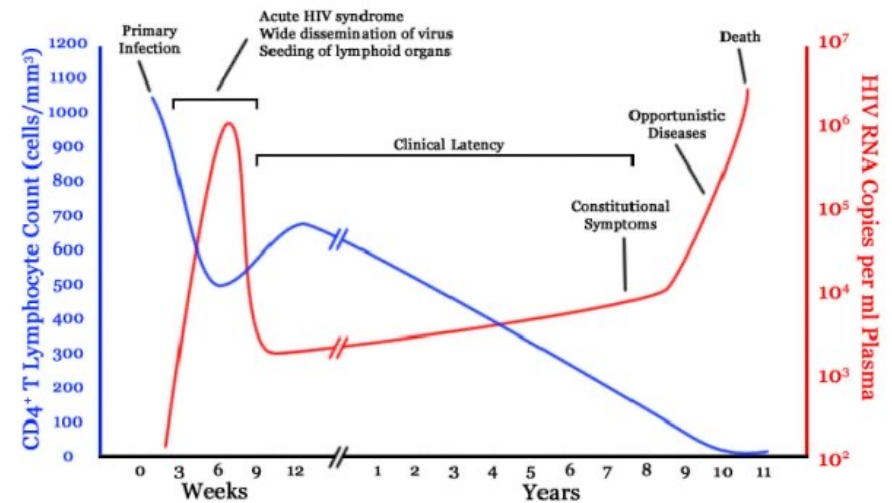
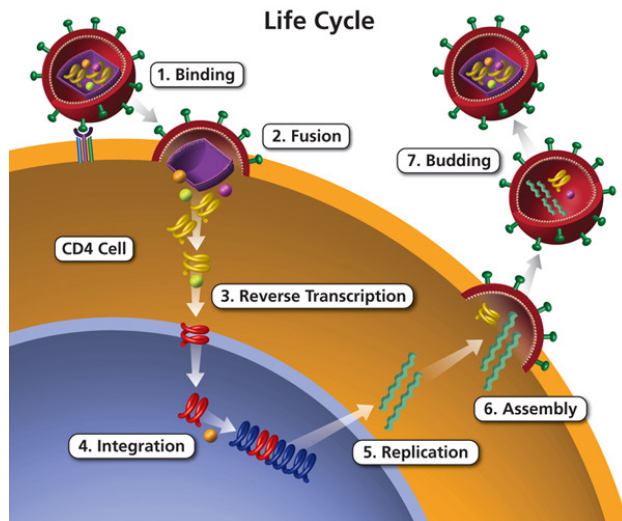
HIV infection in a human host



<https://aidsinfo.nih.gov/understanding-hiv-aids/glossary/1596/life-cycle>

The study question

What changes take place in the first 12-24 hours of HIV infection in terms of gene expression of host cell and viral replication levels?

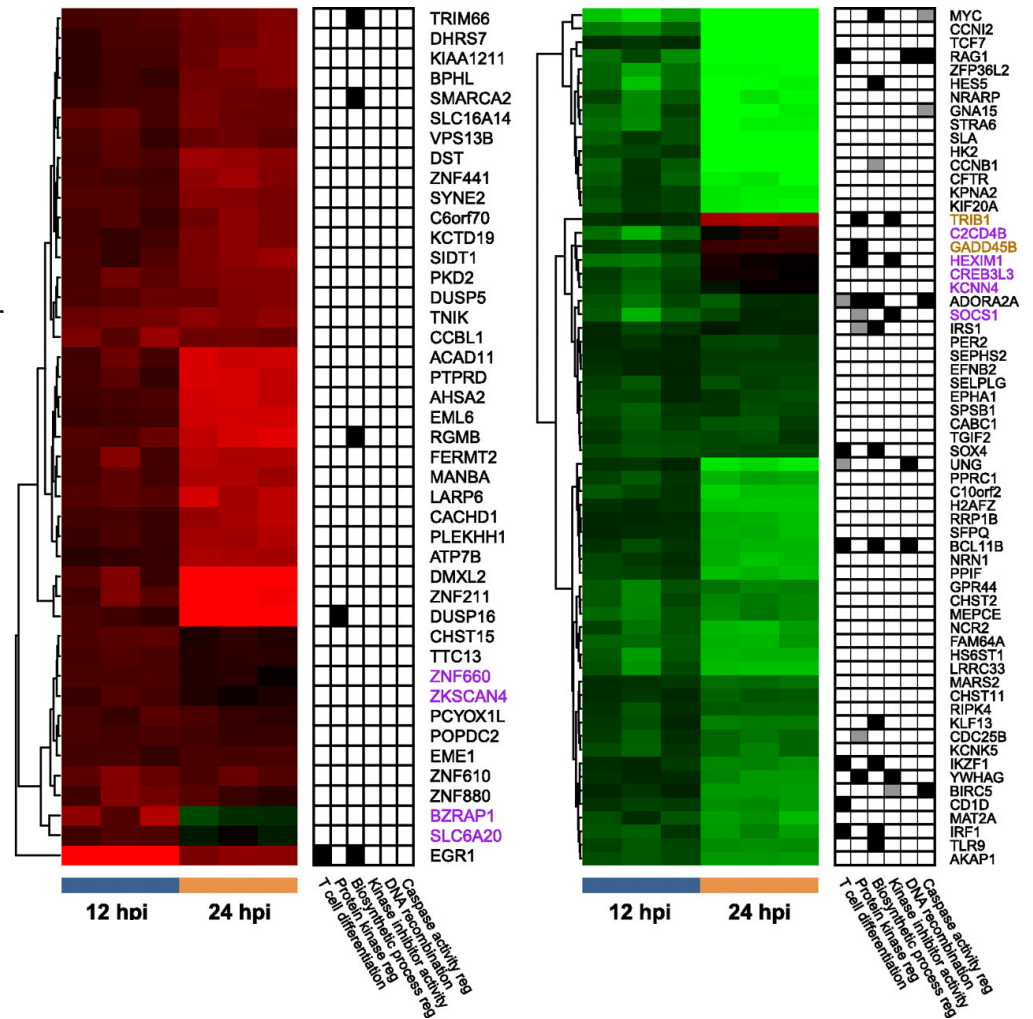
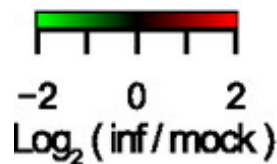


<https://aidsinfo.nih.gov/understanding-hiv-aids/glossary/1596/life-cycle>

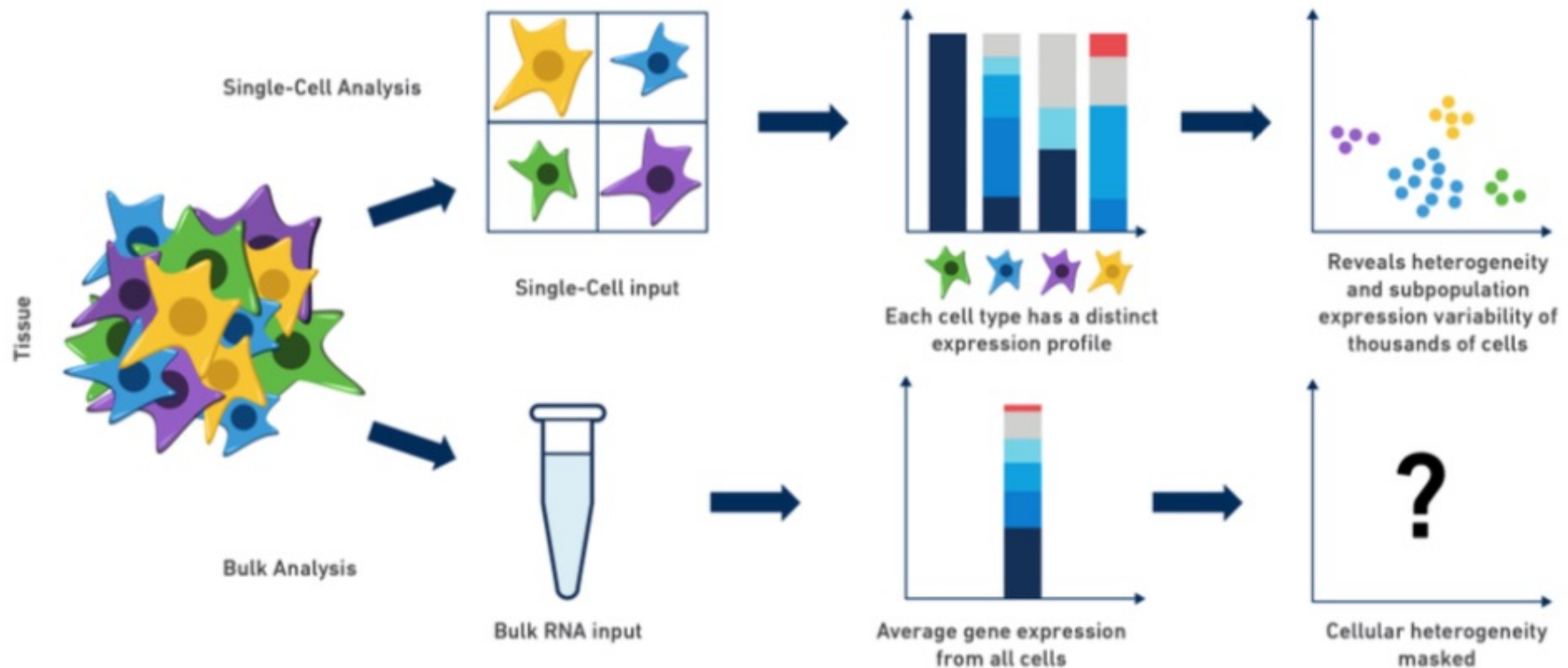
Study findings

Using RNAseq, authors demonstrate:

- 20% of reads mapped to HIV at 12 hr, 40% at 24hr
- Downregulation of T cell differentiation genes at 12hr
- 'Large-scale disruptions to host transcription' at 24hr

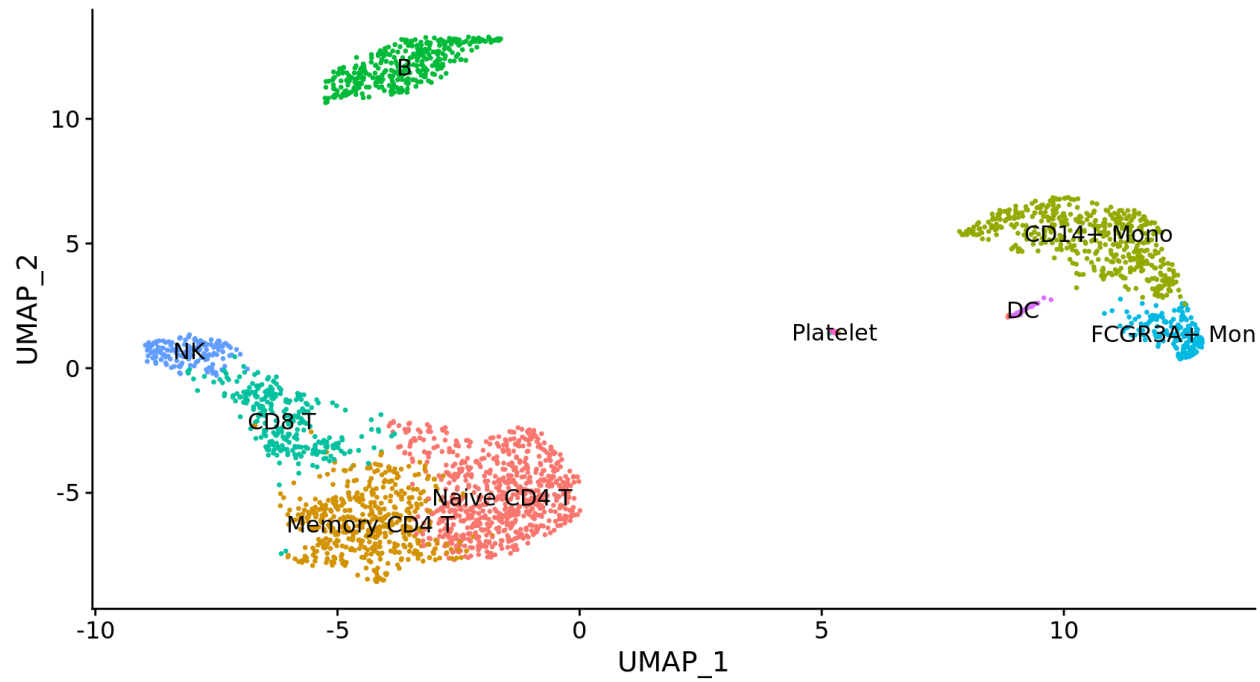


Bulk vs Single Cell RNA Sequencing



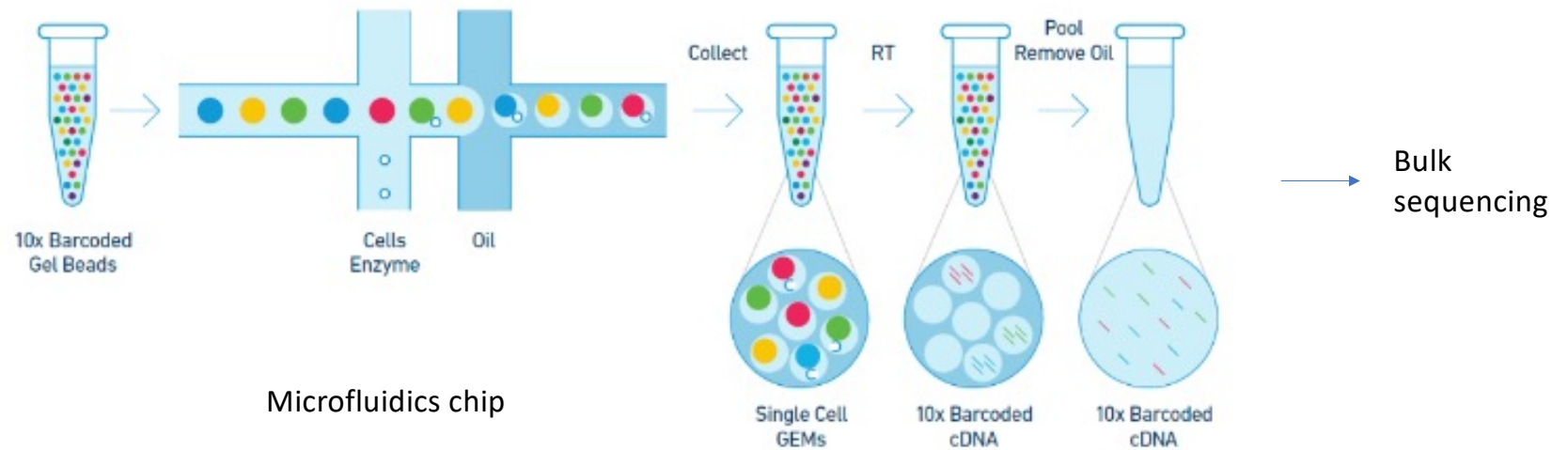
<https://www.10xgenomics.com/blog/single-cell-rna-seq-an-introductory-overview-and-tools-for-getting-started>

scRNA cell subsets in PBMC



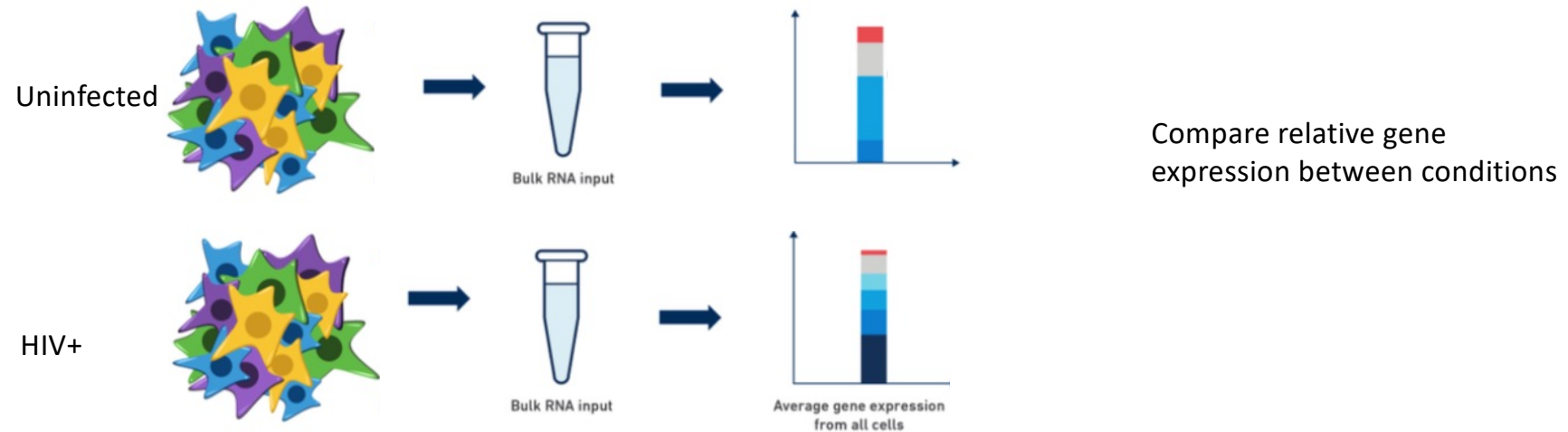
https://satijalab.org/seurat/v3.2/pbmc3k_tutorial.html

10x single cell technology

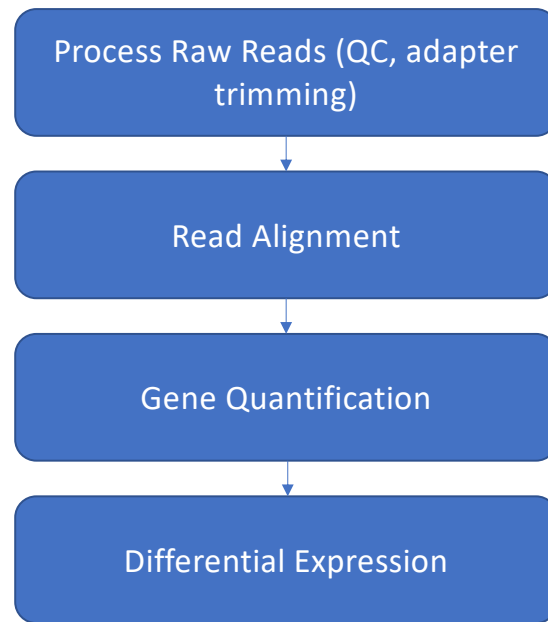


<https://github.com/hbctraining/scRNA-seq>

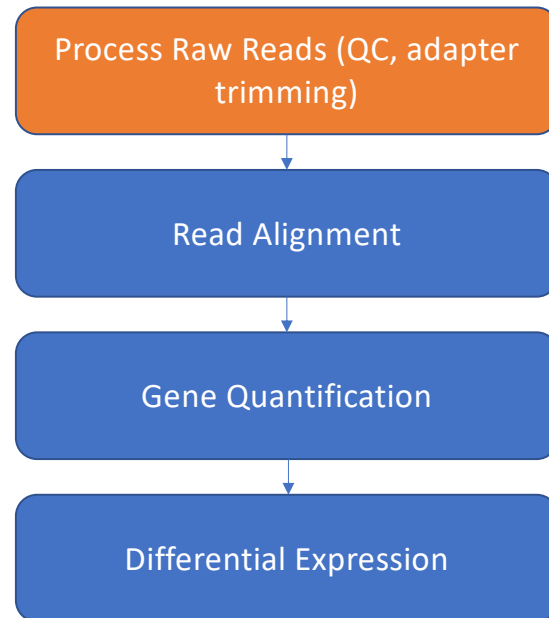
Bulk RNAseq for Differential Expression is OK!



Our (bulk) RNAseq Workflow



Quality control on Raw Reads



Raw reads in Fastq format

```
@SRR098401.109756285  
GACTCACGTAAC TTAACTCTAACAGAAATATACTA...  
+  
CAEFGDG?BCGGGEEDGGHGHGDFHEIEGGDDDD...
```

1. Sequence identifier
2. Sequence
3. + (optionally lists the sequence identifier again)
4. Quality string

Base Quality Scores

The symbols we see in the read quality string are an encoding of the quality score:

```
Quality encoding: !"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHI
                  |         |         |         |         |
Quality score: 0.....10.....20.....30.....40
```

A quality score is a prediction of the probability of an error in base calling:

Quality Score	Probability of Incorrect Base Call	Inferred Base Call Accuracy
10 (Q10)	1 in 10	90%
20 (Q20)	1 in 100	99%
30 (Q30)	1 in 1000	99.9%

Base Quality Scores

The symbols we see in the read quality string are an encoding of the quality score:

```
Quality encoding: !"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHI
                  |         |         |         |         |
Quality score: 0.....10.....20.....30.....40
```

A quality score is a prediction of the probability of an error in base calling:

Quality Score	Probability of Incorrect Base Call	Inferred Base Call Accuracy
10 (Q10)	1 in 10	90%
20 (Q20)	1 in 100	99%
30 (Q30)	1 in 1000	99.9%

Back to our read:

```
@SRR098401.109756285
GACTCACGTAAC TTAACTCTAACAGAAATATACTA...
+
CAEFGDG?BCGGGEEDGGHGHGDFHEIEGGDDDD...
```

C → Q = 34 → Probability < 1/1000 of an error

<https://www.illumina.com/science/education/sequencing-quality-scores.html>

Raw read quality control

Fastq File

@SRR497699.30343179.1 HWI-EAS39X_10175_FC61MK0_4.117_4812_10346 length=75
CAGATGGCCGACAGGAAAGCCATGAAGGCCCTGCATGGGAGATCGGAAGAGCGGTTACAGCAGGAATGCCGAGAC
+
IIIIIGIIHFIIBIIDII>IIDHIIHDIIGIFIIIEIGIBDDEFIG<EIEGEEG;<DB@A8CC7<<@C@BDDDB
@SRR497699.11626500.1 HWI-EAS39X_10175_FC61MK0_4.44_8384_16550 length=75
CGTACTGAACGTACAACGCTGATGCCATCCGCATATTAAATTCGGCAGCGTTAATTAACCTCCTGACCTCGGG
+
HHHHHHHHHHHHHHHHHHHHHHB@HHHHHHHHHHHHHHHHHHHHHHHHHHHHHGEHDHHEHHHHBHHHGHHHHHHHHG
@SRR497699.29057557.1 HWI-EAS39X_10175_FC61MK0_4.112_12508_19308 length=75
CCGAGGCTTAGCTTTTCATTATCACTGTCTCCAGGGTGTGCTTGTCAAAGAGATAAGATCGGAAGAGCGGTTACG
+
GGGBGGDGBHHHDHHEGGGHHHHHHGHGHHHHHHGBGGDGGEGDHHHHHHHHHHHH@BHHGGHGHHHHHHEEGHH
@SRR497699.1331889.1 HWI-EAS39X_10175_FC61MK0_4_5_4738_15920 length=75
CTTACTTTGTAGCCTTCATCAGGGTTTGCTGAAGATGGCGGTATATAGGCTGAGCAAGAGGTGGTGAGGTGTATC
+
HHHHHHHHHHHHGGGGHHHGHGHEEEGGEDGGGGGHHHHHGGEGBDGGGDDGBGGC<EADBEBC<GGGGBEEDG

...

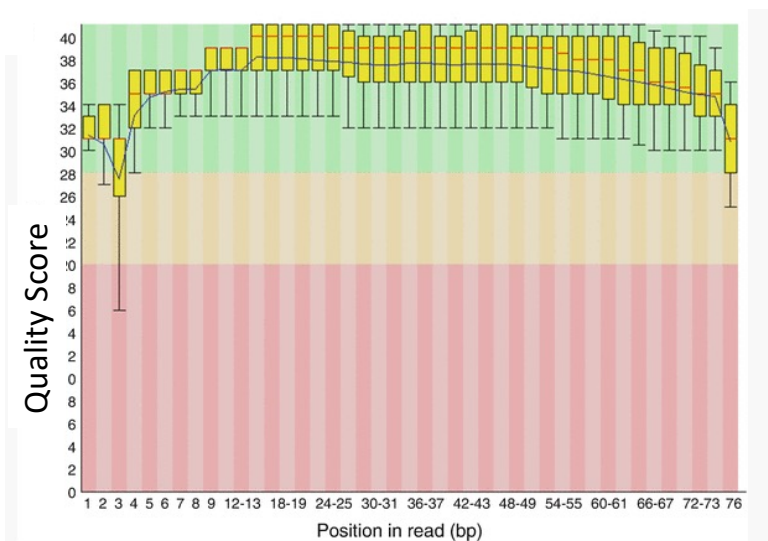
Metrics

- Sequence Quality
- GC content
- Per base sequence content
- Adapters in Sequence

FastQC Tool



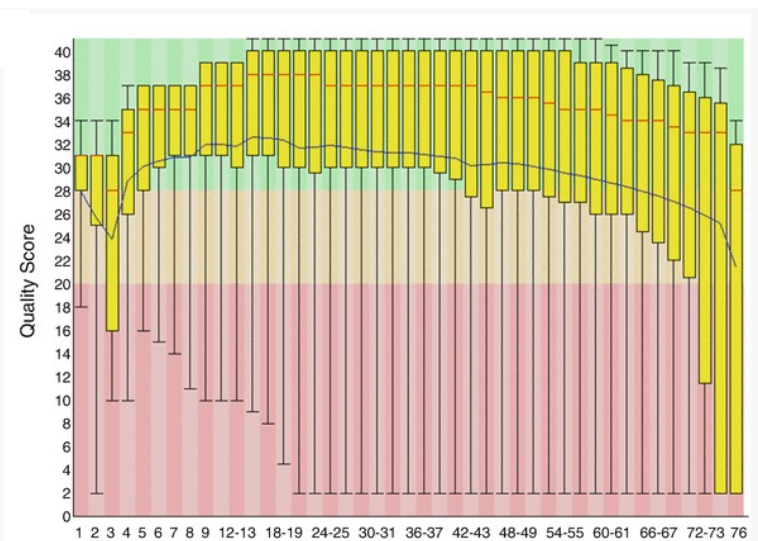
FastQC: Sequence Quality Histogram



Position in read (bp)

GOOD

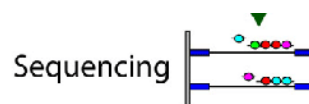
High quality over the length of the read



Position in read (bp)

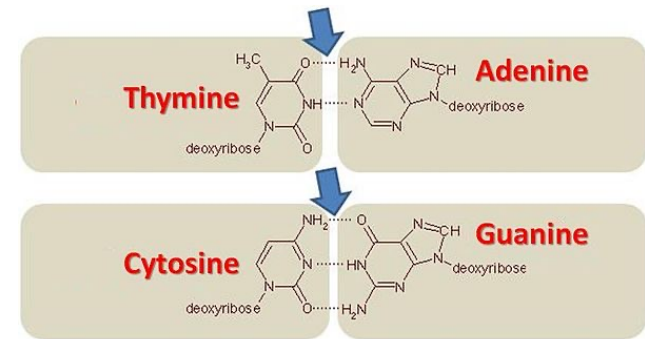
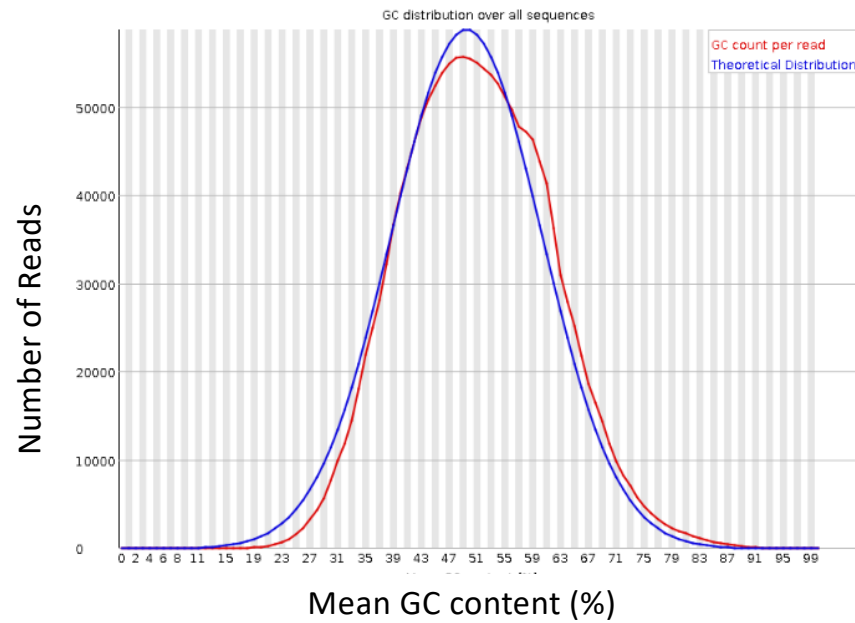
BAD

Read quality drops at the beginning and end



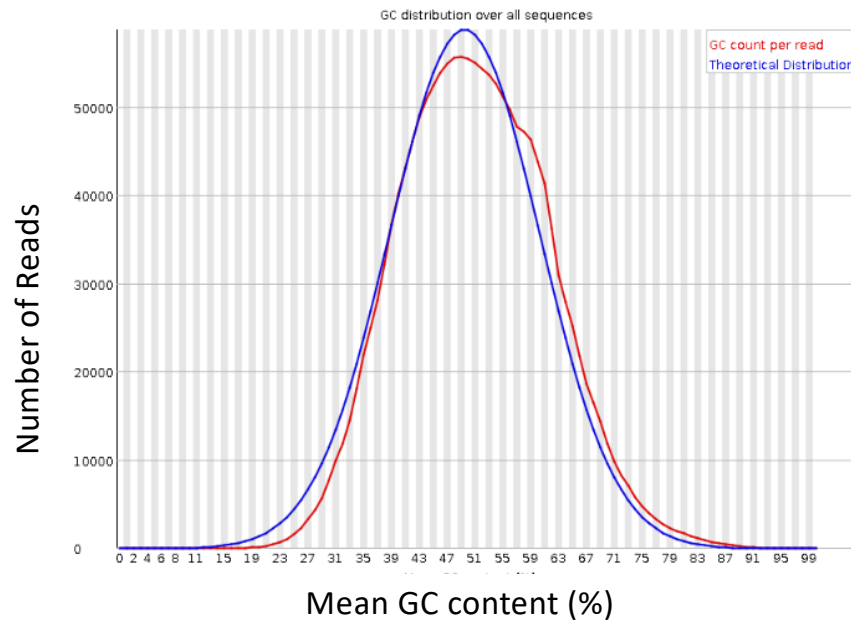
FastQC: Per sequence GC content

✓ Per sequence GC content



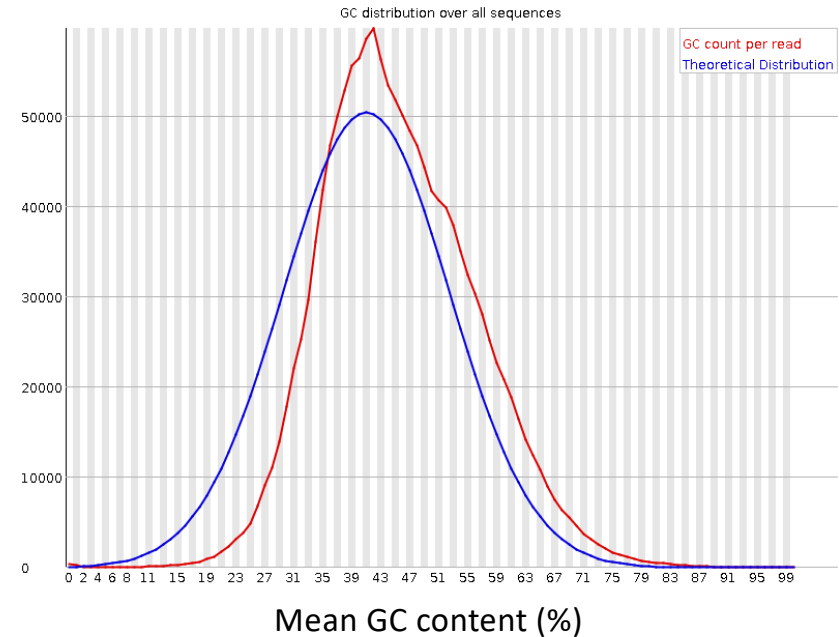
FastQC: Per sequence GC content

✓ Per sequence GC content



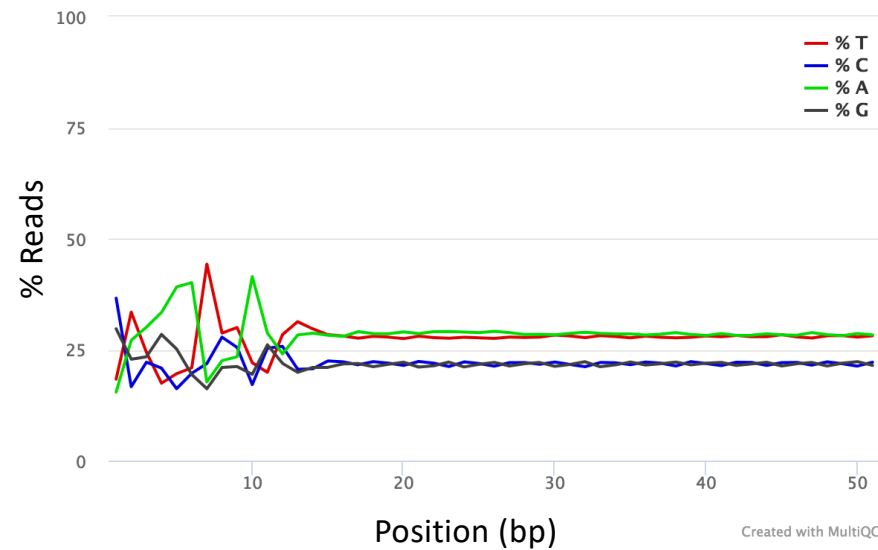
GOOD: follows normal distribution (sum of deviations is < 15% of reads)

✗ Per sequence GC content



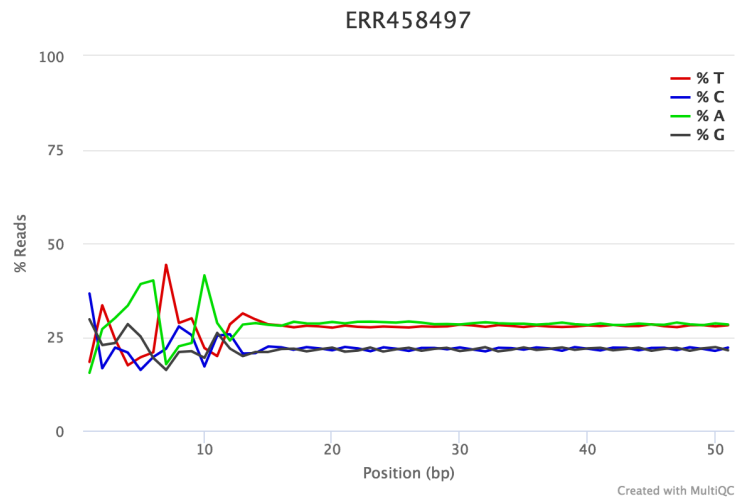
BAD: can indicate contamination with adapter dimers, or another species

FastQC: Per Base Sequence Content

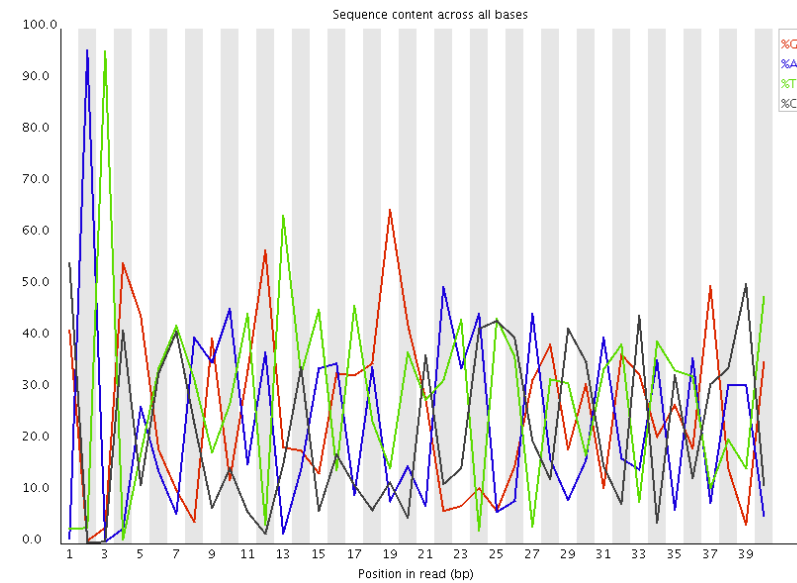


- Proportion of each position for which each DNA base has been called
- RNAseq data tends to show a positional sequence bias in the first ~12 bases
- The "random" priming step during library construction is not truly random and certain hexamers are more prevalent than others

FastQC: Per Base Sequence Content



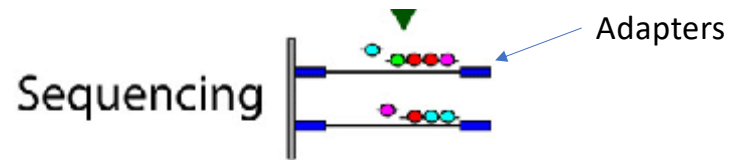
EXPECTED for RNAseq



BAD:

Shows a strong positional bias throughout the reads, which in this case is due to the library having a certain sequence that is overrepresented

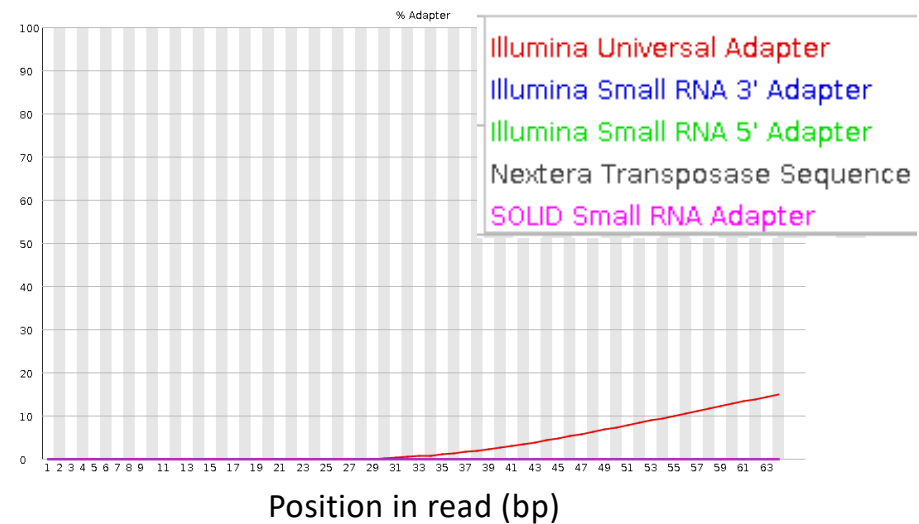
FastQC: Adapter content



FastQC will scan each read for the presence of known adapter sequences

The plot shows that the adapter content rises over the course of the read

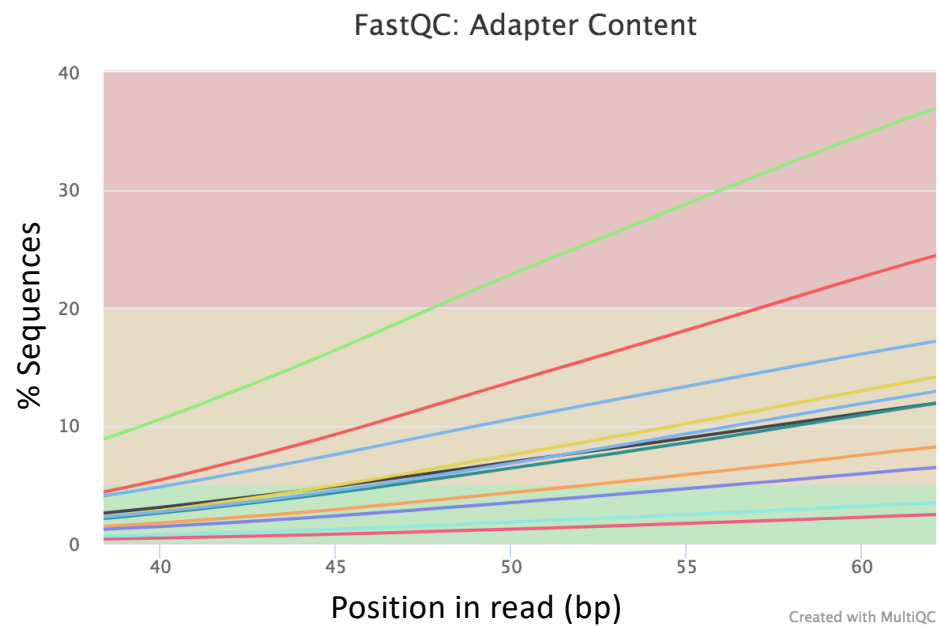
Solution – Adapter trimming!



sequencing.qcfail.com

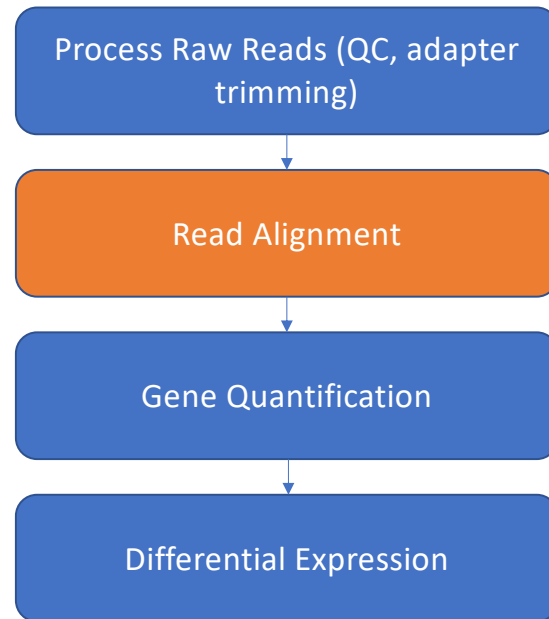
FastQC -> MultiQC

Should view all samples at once to notice abnormalities for our dataset.



We'll use a tool called "Trim Galore!" to trim adapters and remove low quality bases/reads.

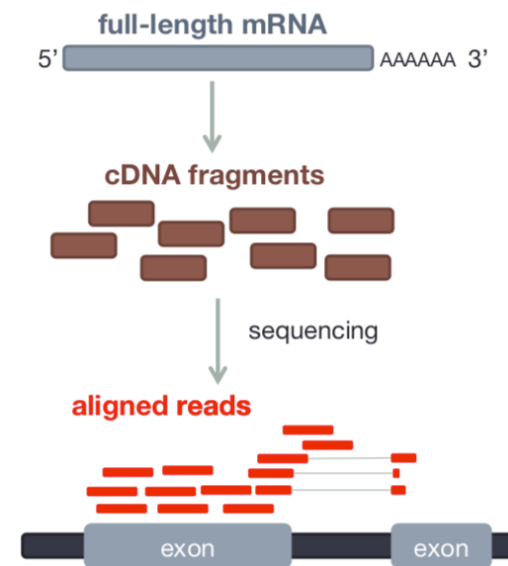
Workflow



Read Alignment

- RNAseq data originates from spliced mRNA (no introns)
- When aligning to the genome, our aligner must find a spliced alignment for reads
- We use a tool called STAR (Spliced Transcripts Alignment to a Reference) that has an exon-aware mapping algorithm.

Reference sequence



[Dobin et al Bioinformatics 2013](#)

Sequence Alignment Map (SAM)



```
@HD VN:1.5 SO:coordinate
@SQ SN:ref LN:45
r001 99 ref 7 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATACTG *
r002 0 ref 9 30 3S6M1P1I4M * 0 0 AAAAGATAAGGATA *
r003 0 ref 9 30 5S6M * 0 0 GCCTAAGCTAA * SA:Z:ref,29,-,6H5M,17,0;
r004 0 ref 16 30 6M14N5M * 0 0 ATAGCTTCAGC *
r003 2064 ref 29 17 6H5M * 0 0 TAGGC * SA:Z:ref,9,+,5S6M,30,1;
r001 147 ref 37 30 9M = 7 -39 CAGCGGCAT * NM:i:1
```

Header
section

Alignment
section

↑ CIGAR: summary of alignment, e.g. match, gap, insertion, deletion

↑ Mapping Quality

↑ Position

↑ Ref Sequence name

↑ Flag: indicates alignment information e.g. paired, aligned, etc
<https://broadinstitute.github.io/picard/explain-flags.html>

↑ Read ID

www.samformat.info

Sequence Alignment Map (SAM)



```
@HD VN:1.5 SO:coordinate
@SQ SN:ref LN:45
r001 99 ref 7 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATACTG *
r002 0 ref 9 30 3S6M1P1I4M * 0 0 AAAAGATAAGGATA *
r003 0 ref 9 30 5S6M * 0 0 GCCTAAGCTAA * SA:Z:ref,29,-,6H5M,17,0;
r004 0 ref 16 30 6M14N5M * 0 0 ATAGCTTCAGC *
r003 2064 ref 29 17 6H5M * 0 0 TAGGC * SA:Z:ref,9,+,5S6M,30,1;
r001 147 ref 37 30 9M = 7 -39 CAGCGGCAT * NM:i:1
```

Header
section

Alignment
section

Paired end info

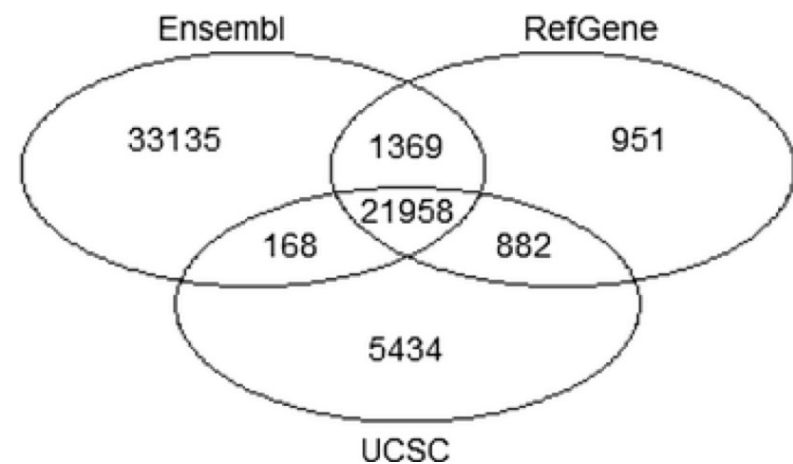
Sequence

Quality Score

Optional Fields

Genome Annotation Standards

- STAR can use an annotation file gives the location and structure of genes in order to improve alignment in known splice junctions
- Annotation is dynamic and there are at least three major sources of annotation
- The intersection among RefGene, UCSC, and Ensembl annotations shows high overlap. RefGene has the fewest unique genes, while more than 50% of genes in Ensembl are unique
- Be consistent with your choice of annotation source!



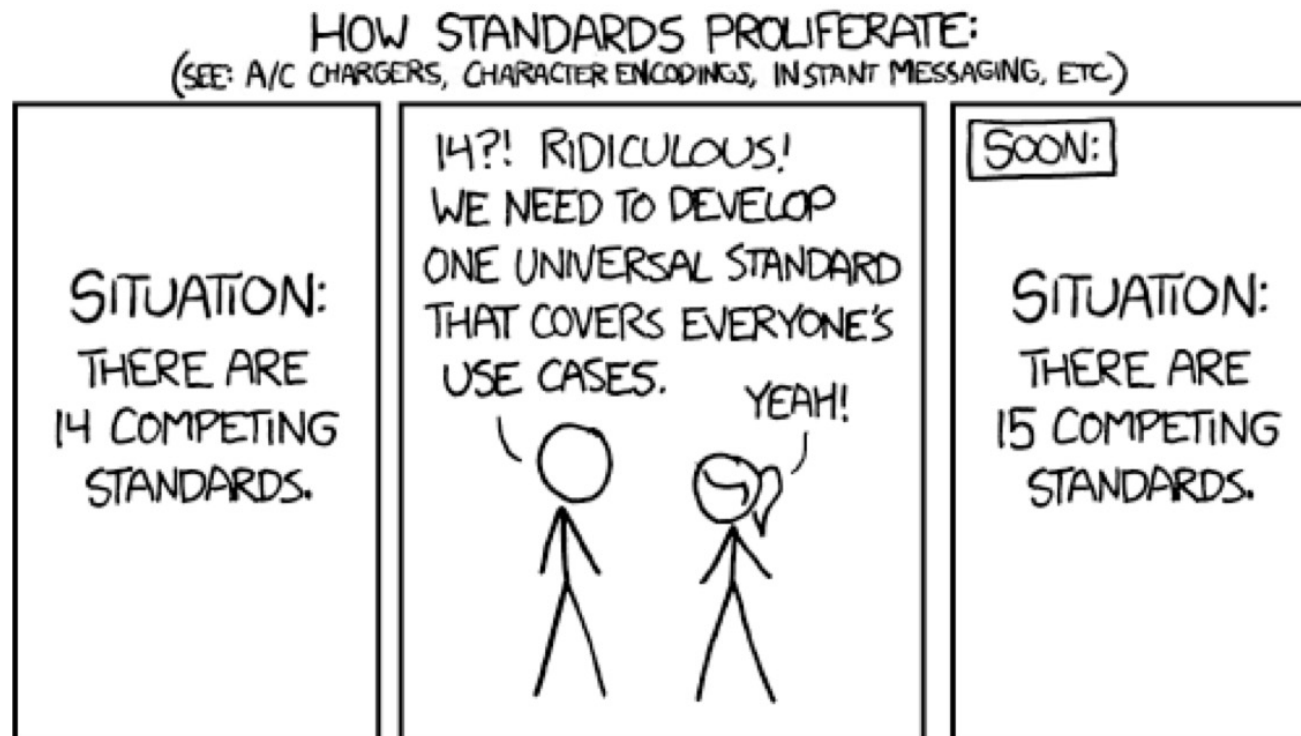
Gene Annotation Format (GTF)

In order to count genes, we need to know where they are located in the reference sequence
STAR uses a Gene Transfer Format (GTF) file for gene annotation

Chrom	Source	Feature type	Start	Stop	Strand Frame (Score)			Attribute
chr5	hg38_refGene	exon	138465492	138466068	.	+	.	gene_id "EGR1";
chr5	hg38_refGene	CDS	138465762	138466068	.	+	0	gene_id "EGR1";
chr5	hg38_refGene	start_codon	138465762	138465764	.	+	.	gene_id "EGR1";
chr5	hg38_refGene	CDS	138466757	138468078	.	+	2	gene_id "EGR1";
chr5	hg38_refGene	exon	138466757	138469315	.	+	.	gene_id "EGR1";
chr5	hg38_refGene	stop_codon	138468079	138468081	.	+	.	gene_id "EGR1";

<https://useast.ensembl.org/info/website/upload/gff.html>

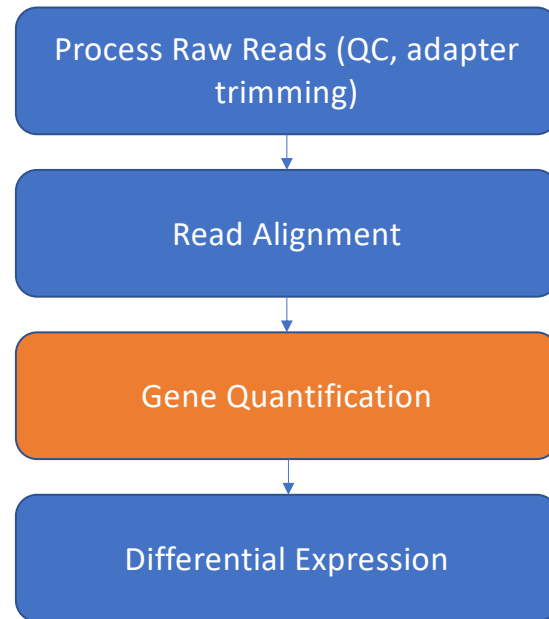
A note on standards



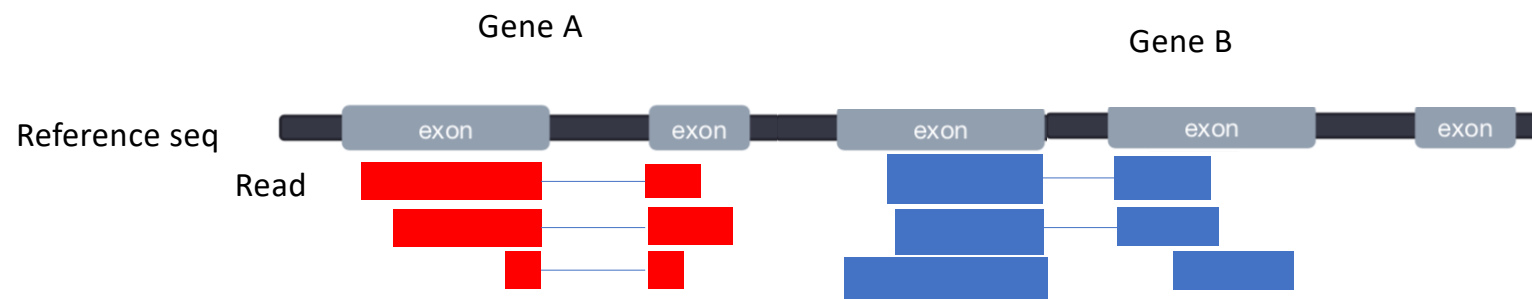
Visualizing reads with JBrowse



Workflow

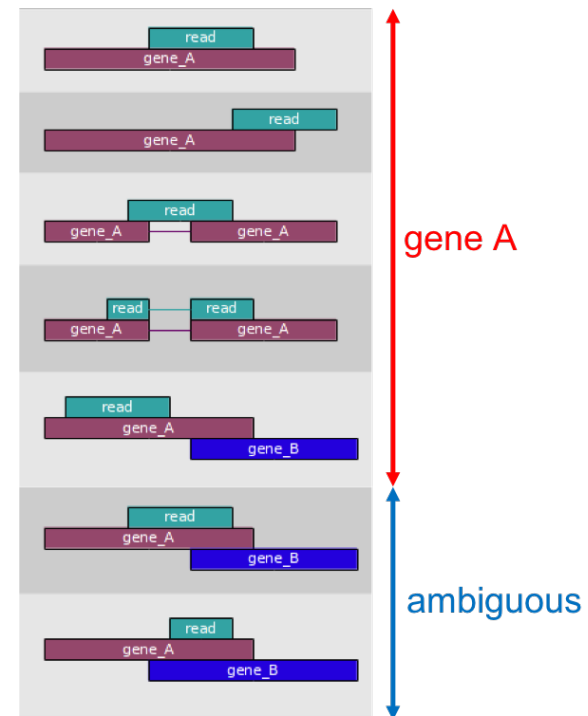


Counting reads for each gene



Counting reads: featurecounts

- The mapped coordinates of each read are compared with the features in the GTF file
- Reads that overlap with a gene by ≥ 1 bp are counted as belonging to that feature
- Ambiguous reads will be discarded

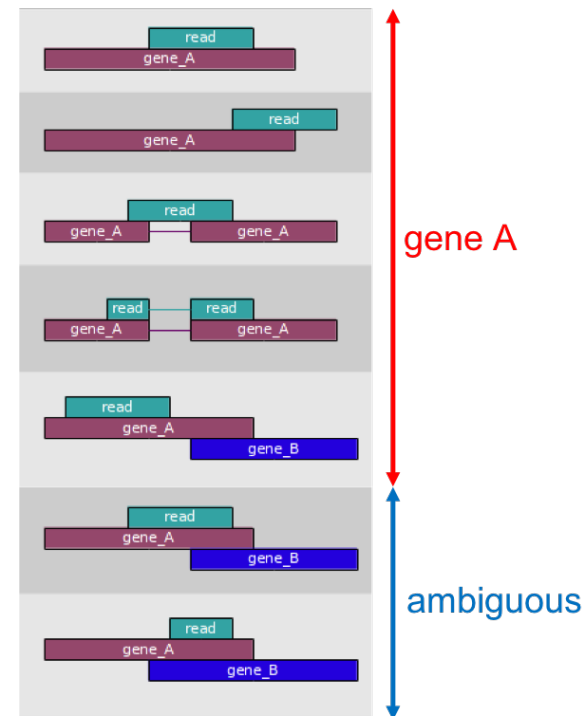


Counting reads: featurecounts

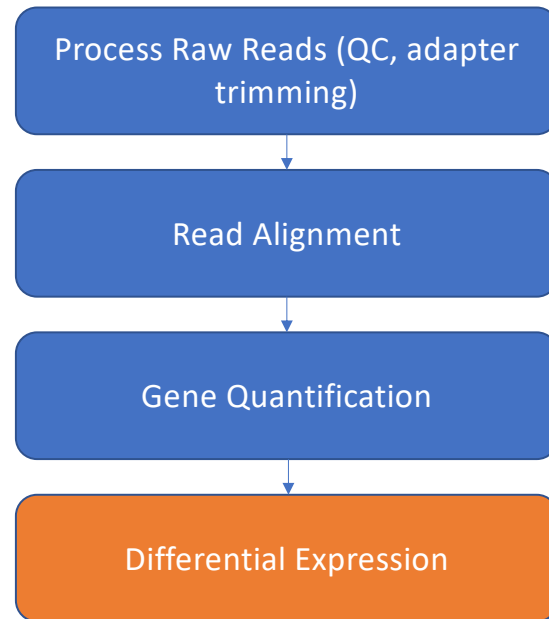
- The mapped coordinates of each read are compared with the features in the GTF file
- Reads that overlap with a gene by ≥ 1 bp are counted as belonging to that feature
- Ambiguous reads will be discarded

Result is a gene count matrix:

Gene	Sample 1	Sample 2	Sample 3	Sample 4
A	1000	1000	100	10
B	10	1	5	6
C	10	1	10	20

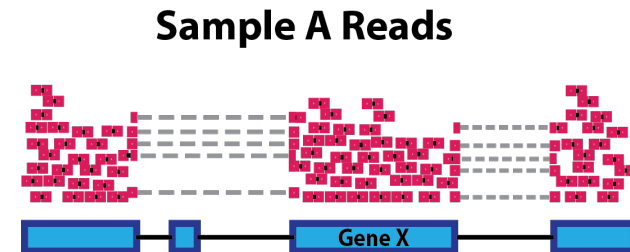


Workflow



Normalization

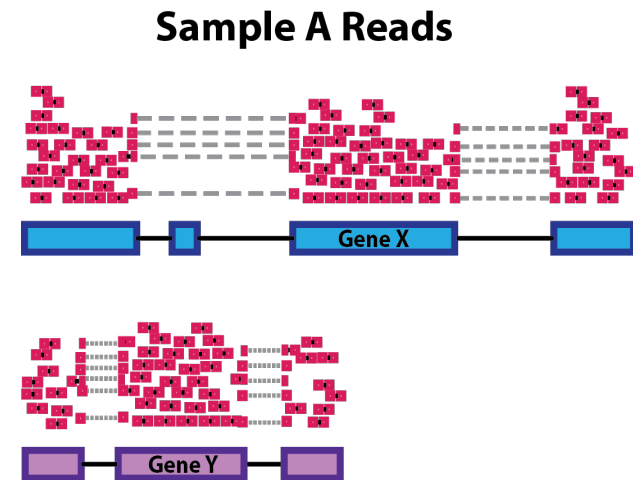
- Raw Count \neq Expression strength
- Normalization:
 - Eliminates factors that are not of interest for our experiment
 - Enables accurate comparison between samples or genes



Normalization

The number of reads mapped to a gene depends on

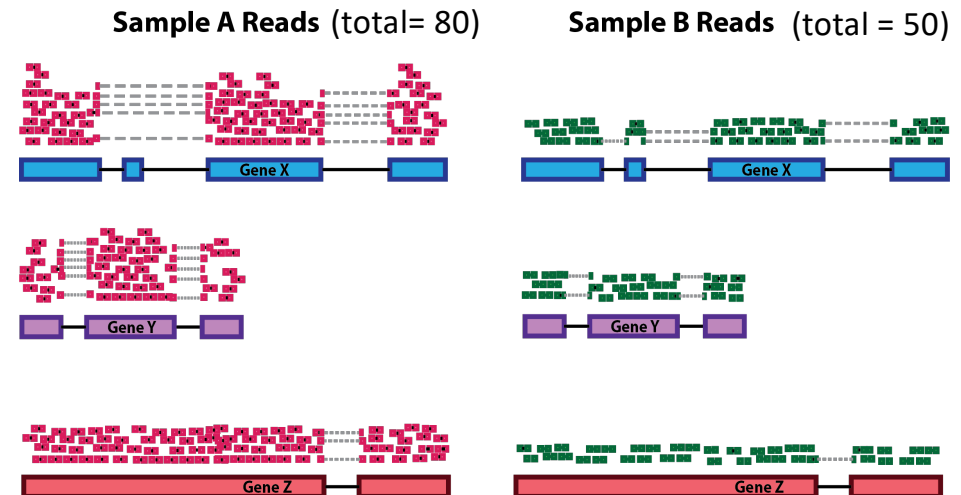
- **Gene Length**



Normalization

The number of reads mapped to a gene depends on

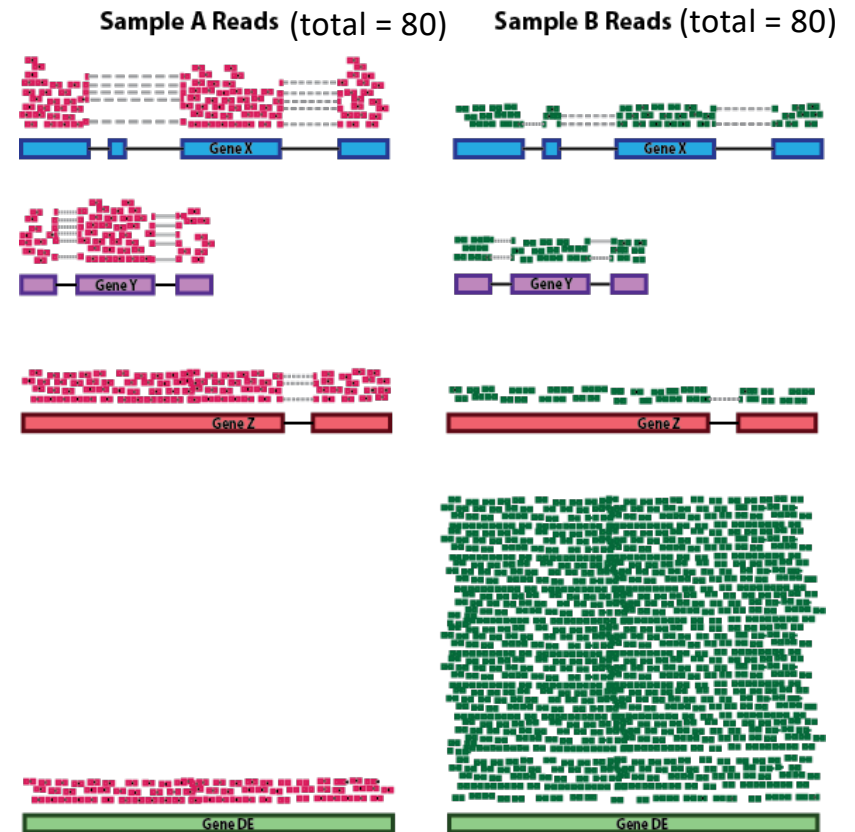
- Gene Length
- **Sequencing depth**



Normalization

The number of reads mapped to a gene depends on

- Gene Length
- Sequencing depth
- **The expression level of other genes in the sample (RNA Composition)**



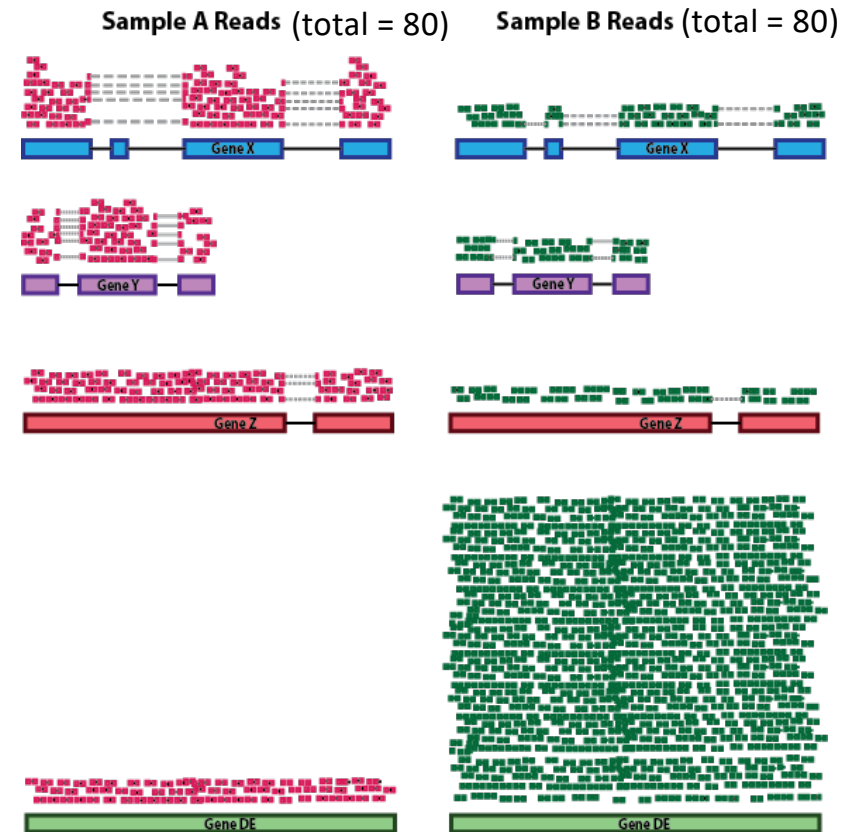
Adapted from https://hbctraining.github.io/DGE_workshop

Normalization

The number of reads mapped to a gene depends on

- Gene Length
- Sequencing depth
- The expression level of other genes in the sample (RNA Composition)

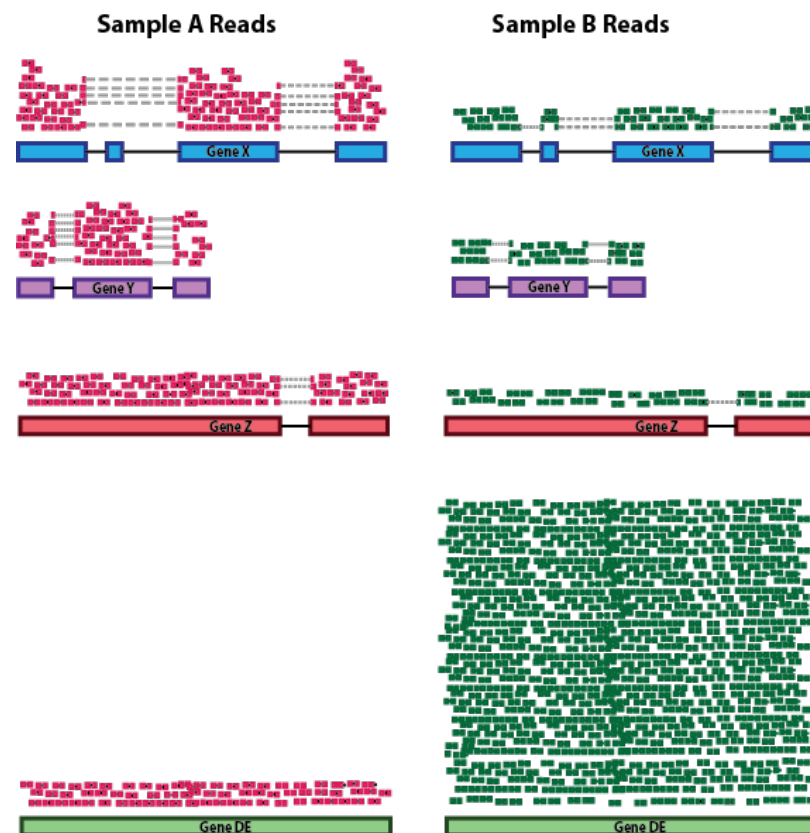
DESeq2 Median of Ratios



Adapted from https://hbctraining.github.io/DGE_workshop

Normalization: DESeq2 Median of Ratios

Gene	Sample A	Sample B
X	26	10
Y	26	10
Z	26	10
DE	2	50
Total =	80	80



Normalization: DESeq2 Median of Ratios

1. Take a row-wise average to produce an average sample (geometric mean) $\sqrt[n]{x_1 x_2 \cdots x_n}$

Gene	Sample A	Sample B	Avg. Sample
X	26	10	16
Y	26	10	16
Z	26	10	16
DE	2	50	10

Normalization: DESeq2 Median of Ratios

1. Take a row-wise average to produce an average sample (geometric mean) $\sqrt[n]{x_1 x_2 \cdots x_n}$

Gene	Sample A	Sample B	Avg. Sample
X	26	10	16
Y	26	10	16
Z	26	10	16
DE	2	50	10

2. Divide all rows by the Average Sample for that gene (**Ratio**)

Gene	Sample A/Avg.	Sample B /Avg.
X	26/16 = 1.6	10/16 = 0.6
Y	1.6	0.6
Z	1.6	0.6
DE	0.2	5

Normalization: DESeq2 Median of Ratios

1. Take a row-wise average to produce an average sample (geometric mean) $\sqrt[n]{x_1 x_2 \cdots x_n}$

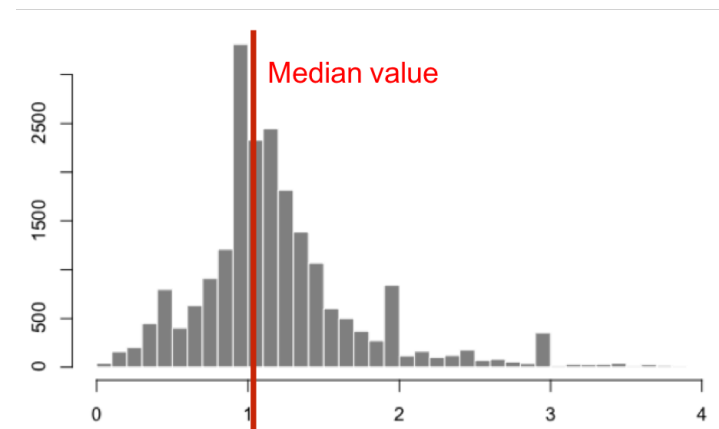
Gene	Sample A	Sample B	Avg. Sample
X	26	10	16
Y	26	10	16
Z	26	10	16
DE	2	50	16

2. Divide all rows by the Average Sample for that gene (**Ratio**)

Gene	Sample A/Avg.	Sample B /Avg.
X	26/16 = 1.6	10/16 = 0.6
Y	1.6	0.6
Z	1.6	0.6
DE	0.2	5

3. Take the **median** of each column. Should be ~1 for all

Size factor	1.6	0.6
-------------	-----	-----



Normalization: DESeq2 Median of Ratios

1. Take a row-wise average to produce an average sample (geometric mean) $\sqrt[n]{x_1 x_2 \cdots x_n}$

Gene	Sample A	Sample B	Avg. Sample
X	26	10	16
Y	26	10	16
Z	26	10	16
DE	2	50	16

2. Divide all rows by the Average Sample for that gene (**Ratio**)
4. Divide all counts by sample specific size factor

Gene	Sample A / Avg.	Sample B / Avg.
X	26/16 = 1.6	10/16 = 0.6
Y	1.6	0.6
Z	1.6	0.6
DE	0.2	5

Gene	Sample A / S_A	Sample B / S_B
X	16.3	16.7
Y	16.3	16.7
Z	16.3	16.7
DE	1.3	83.3

Normalized counts for non-DE genes are similar!

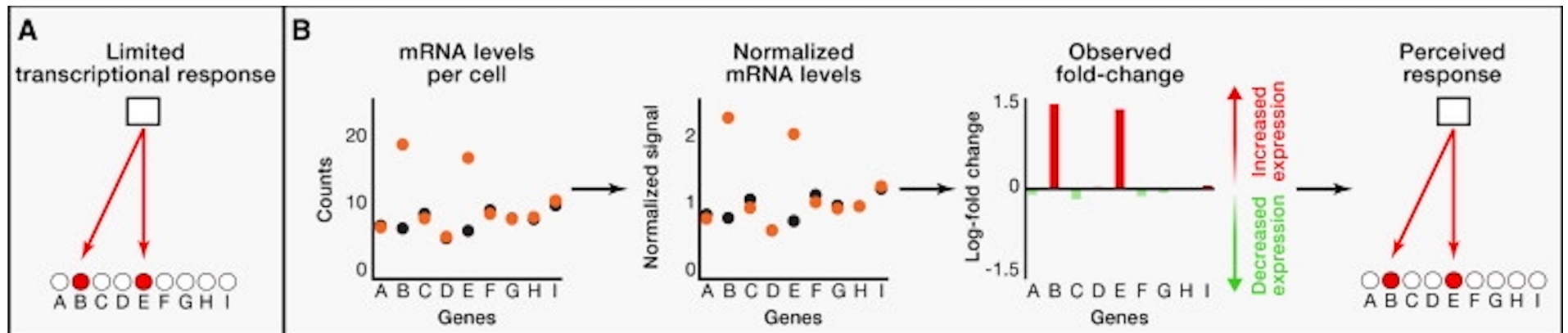
3. Take the **median** of each column. Should be ~1 for all

Size factor	1.6	0.6
-------------	-----	-----

```
estimateSizeFactors(dds)
```

Assumption of DESeq2 Median of Ratios

Median of Ratios method assumes that most genes are not Differentially Expressed between samples.

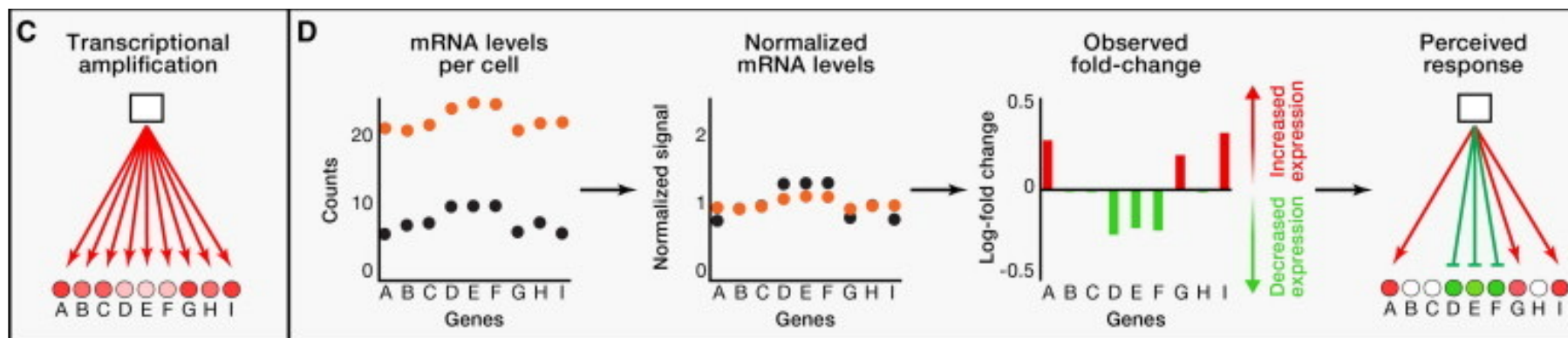


Loven et al "Revisiting Global Gene Expression Analysis" Cell 2012 <https://doi.org/10.1016/j.cell.2012.10.012>

Assumption of DESeq2 Median of Ratios

Median of Ratios method assumes that most genes are not Differentially Expressed between samples.

COUNTER EXAMPLE



NOTE: add back full picture or remove

- Late stage cell death (total RNA DOWN)
- High c-Myc cells (total RNA UP)

Known quantity spike-in transcripts (ERCC) can be used to normalize in these cases.

Loven et al "Revisiting Global Gene Expression Analysis" Cell 2012 <https://doi.org/10.1016/j.cell.2012.10.012>

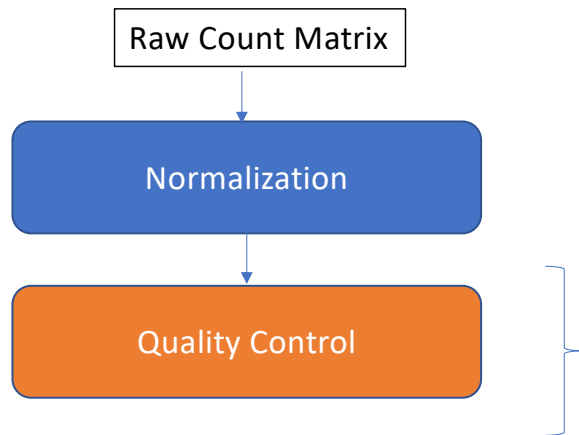
Normalization methods

Normalization method	Description	Accounted factors	Recommended use
CPM (counts per million)	$\frac{K_i}{Total\ Reads\ per\ Sample/10^6}$	sequencing depth	Comparison between replicates of the sample group
R/FPKM (reads/fragments per kilobase of exon per million reads/fragments mapped)	$\frac{K_i}{Gene\ Length/10^3 * Total\ Reads\ per\ Sample/10^6}$	sequencing depth and gene length	Comparison between genes in a sample
DESeq2's median of ratios [1]	K_i divided by sample-specific size factors	sequencing depth and RNA composition	Differential Expression between samples

Similar to DESeq2: EdgeR, limma-voom

Adapted from https://hbctraining.github.io/DGE_workshop

Quality Control Visualizations



Examine sources of variation in the data

- Principal Component Analysis
- Hierarchical Clustering

(Log2 + 1) Transformed, Normalized Count Table

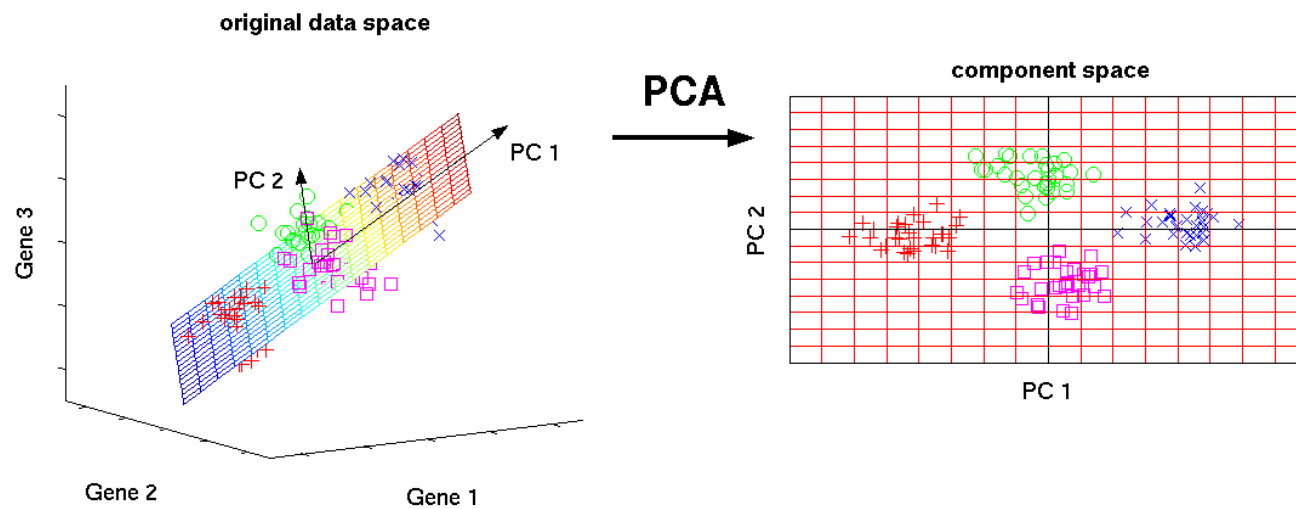
Gene	Sample A	Sample B	Sample C
1	1	1.6	0.5
2	2.2	-0.2	1
3	-1	1	3.1

Principle Component Analysis

Dimension reduction technique

Example: 3 gene dimensions -> 2 PC

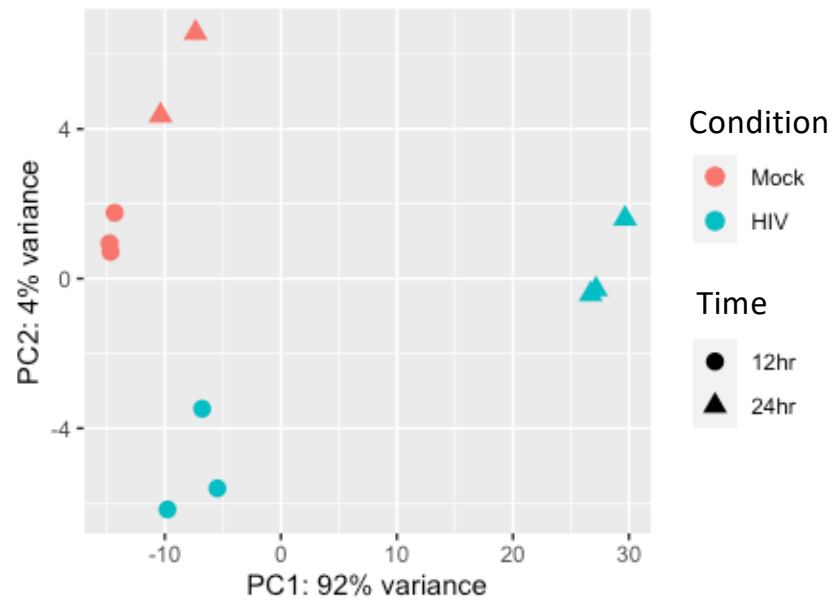
Gene	Mock_12h	Mock_12h	Mock_24h	Mock_24h	HIV_12h	HIV_12h	HIV_24h	HIV_24h
Gene 1	8.9	8.9	8.9	9.0	8.9	8.9	9.0	6.8
Gene 2	0.6	-1.0	0.6	-1.0	0.6	-1.0	0.6	3.8
Gene 3	4.1	11.9	4.1	-0.5	4.1	8.7	4.0	4.4



Do your samples cluster as expected?

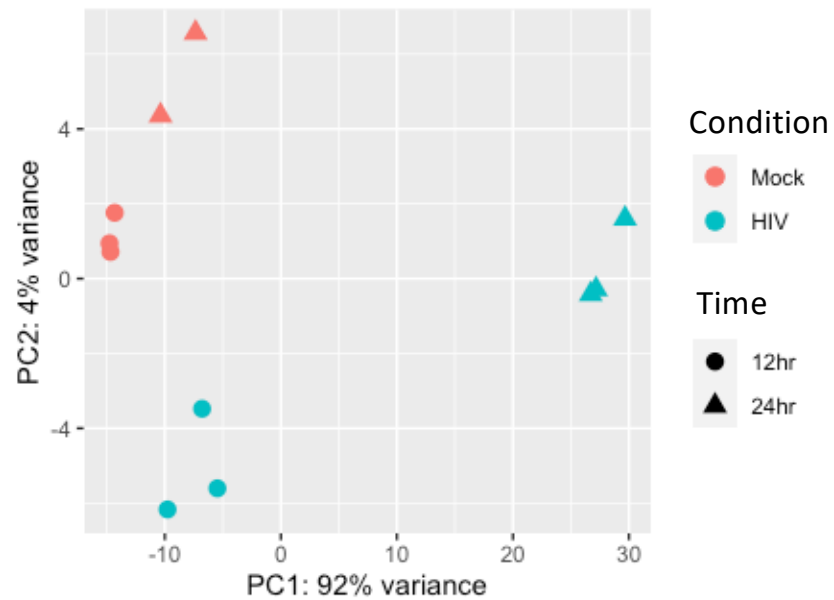
What are the major sources of variation in the data?

Principle Component Analysis



- ✓ Do your samples cluster as expected?
- ✓ What are the major sources of variation in the data?

Principle Component Analysis



- ✓ Do your samples cluster as expected?
- ✓ What are the major sources of variation in the data?
- ✓ Is there a batch effect?

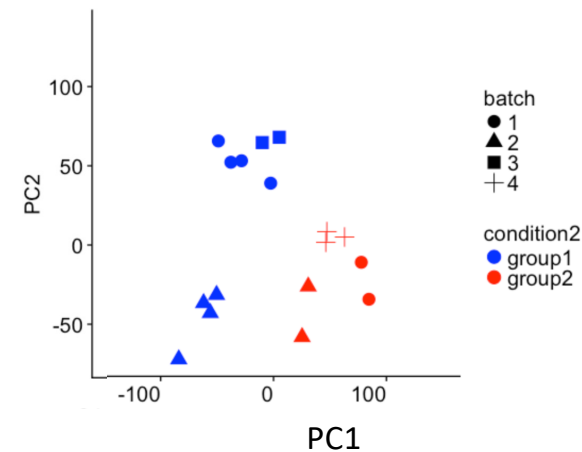
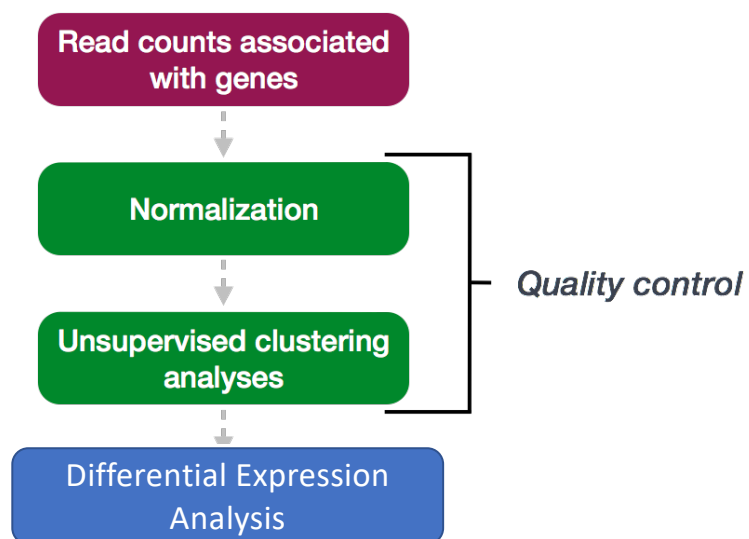


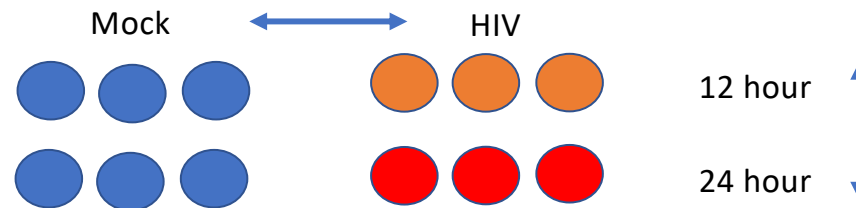
Image <https://support.bioconductor.org/p/111491/>

Differential Expression with DESeq2



https://hbctraining.github.io/DGE_workshop

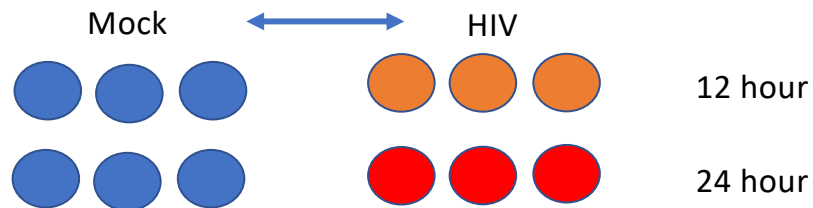
Multi-factor experiment design



Factor 1:
Infection status (Mock or HIV)

Factor 2:
Time (12 or 24 hr)

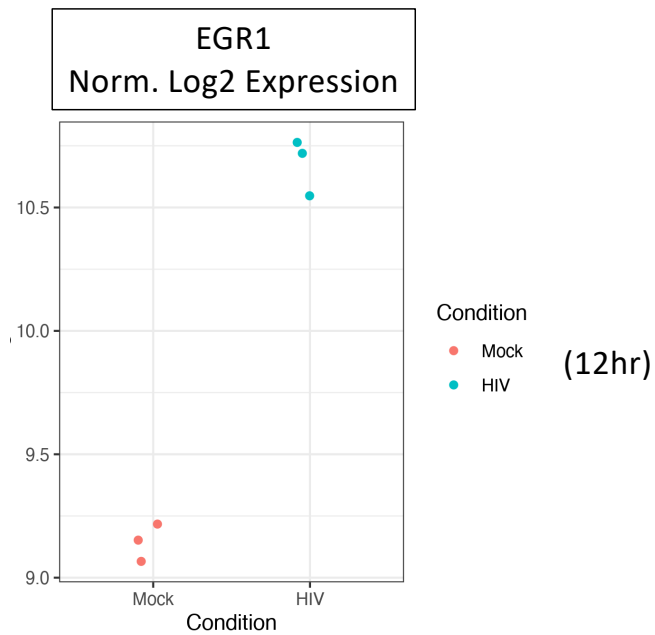
Multi-factor experiment design



- Differential Expression compares two conditions
- We'll choose Infection status at 12 hr (Mock or HIV) for comparison
- We could also choose time, or a combination of multiple factors

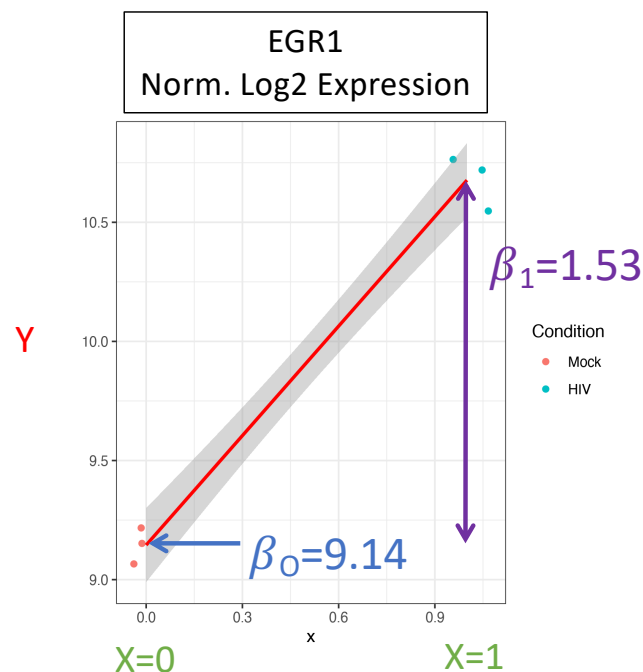
Step 1: Modeling gene expression values

All leading DE tools use **regression models** to estimate the fold change between conditions for **each gene**



Step 1: Modeling gene expression values

All leading DE tools use **regression models** to estimate the fold change between conditions for **each gene**
Example, simple linear regression:



$$Y = \beta_0 + \beta_1 X + e$$

Log2 Expression Values

Intercept

Condition (0-Mock, 1-HIV)

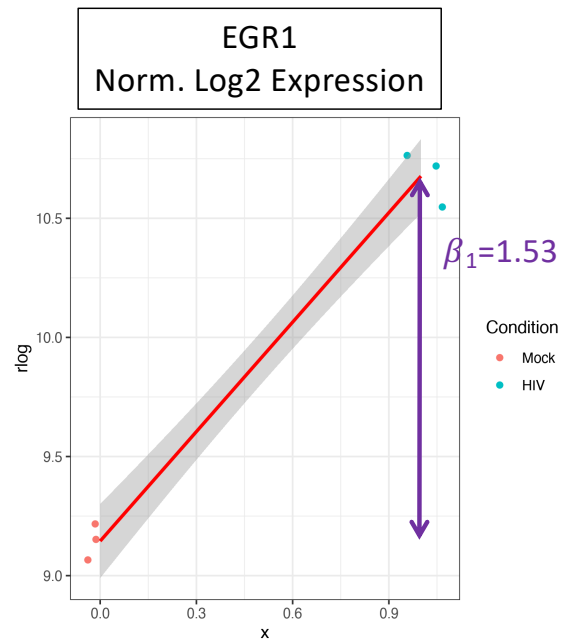
Slope: difference between Mock /HIV

Error

DESeq2 uses a Generalized Linear Model with a Negative Binomial error Distribution, which has been shown to be best fit for RNAseq data.

Love, M.I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol* **15**, 550 (2014).

Step 2: Hypothesis Testing



$$H_0 : \beta_1 = 0 \quad \text{vs.} \quad H_A : \beta_1 \neq 0$$

H_0 : there is no systematic difference between the average read count values for Mock vs. HIV

- Statistical test – Wald test (similar to t-test) on β_1
- $Z = \beta_1 / SE_{\beta_1}$
- Z-statistic is compared to the normal distribution and probability of getting a statistic at least as extreme is computed

Is EGR1 differentially expressed?

Yes! $p \ll 0.05$

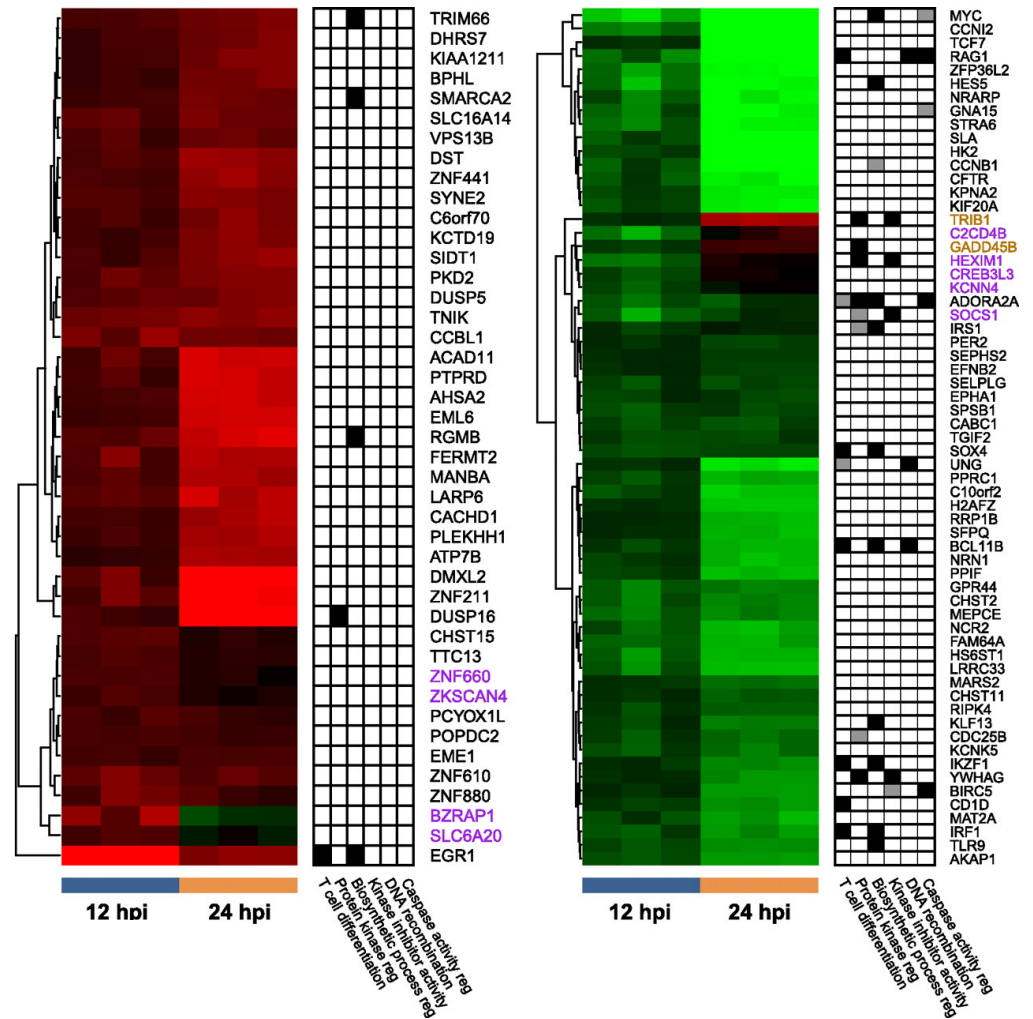
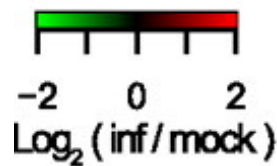
DESeq2 Results table

GeneID	Base mean	log2FoldChange	StdErr	P-value	P-adj
EGR1	1273	1.55	0.13	1.19e-77	1.52e-73
MYC	5226	-1.53	0.14	1.63e-36	1.03e-32

- Mean of normalized counts – averaged over all samples from two conditions (HIV, Mock)
- Log of the fold change between two conditions
- StdErr – Standard error of coefficient (e.g. b_1)
- P-value – the probability that the Wald statistic is as extreme as observed if H_0 were true
- P-adj – accounting for multiple testing correction

Study findings

- T cell differentiation-related genes were overrepresented in the DEG at 24hr
- 'Large-scale disruptions to host transcription' at 24hr



References

DESeq2 vignette (R/Rstudio):

<http://www.bioconductor.org/packages/release/bioc/vignettes/DESeq2/inst/doc/DESeq2.html#differential-expression-analysis>

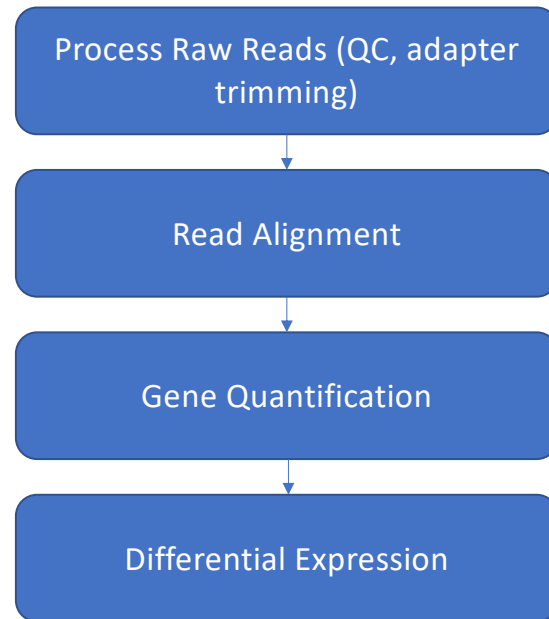
HBC Training (Command line/R):

https://hbctraining.github.io/DGE_workshop

Galaxy Training:

https://galaxyproject.org/tutorials/rb_rnaseq/

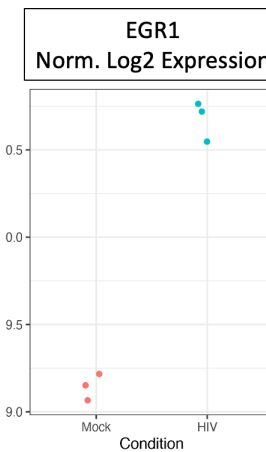
Review



```
@SRR497699.30343179.1 HWI-EAS39X_10175_FC61MK0_4_117_4812_10346 length=75
CAGATGGCCGACAGAGGAGCCATGAAGGCCCTGCATGGGGAGATCGGAAGAGCGGTTTCAGCAGGAATGCCGAGAC
+
IIIIIGIIHFIIIIIBIIDII>IIDHIIHDIIGIFIIIEIGIBDDEFIG<EIEGEEG;<DB@A8CC7<><C@BBD0B
```

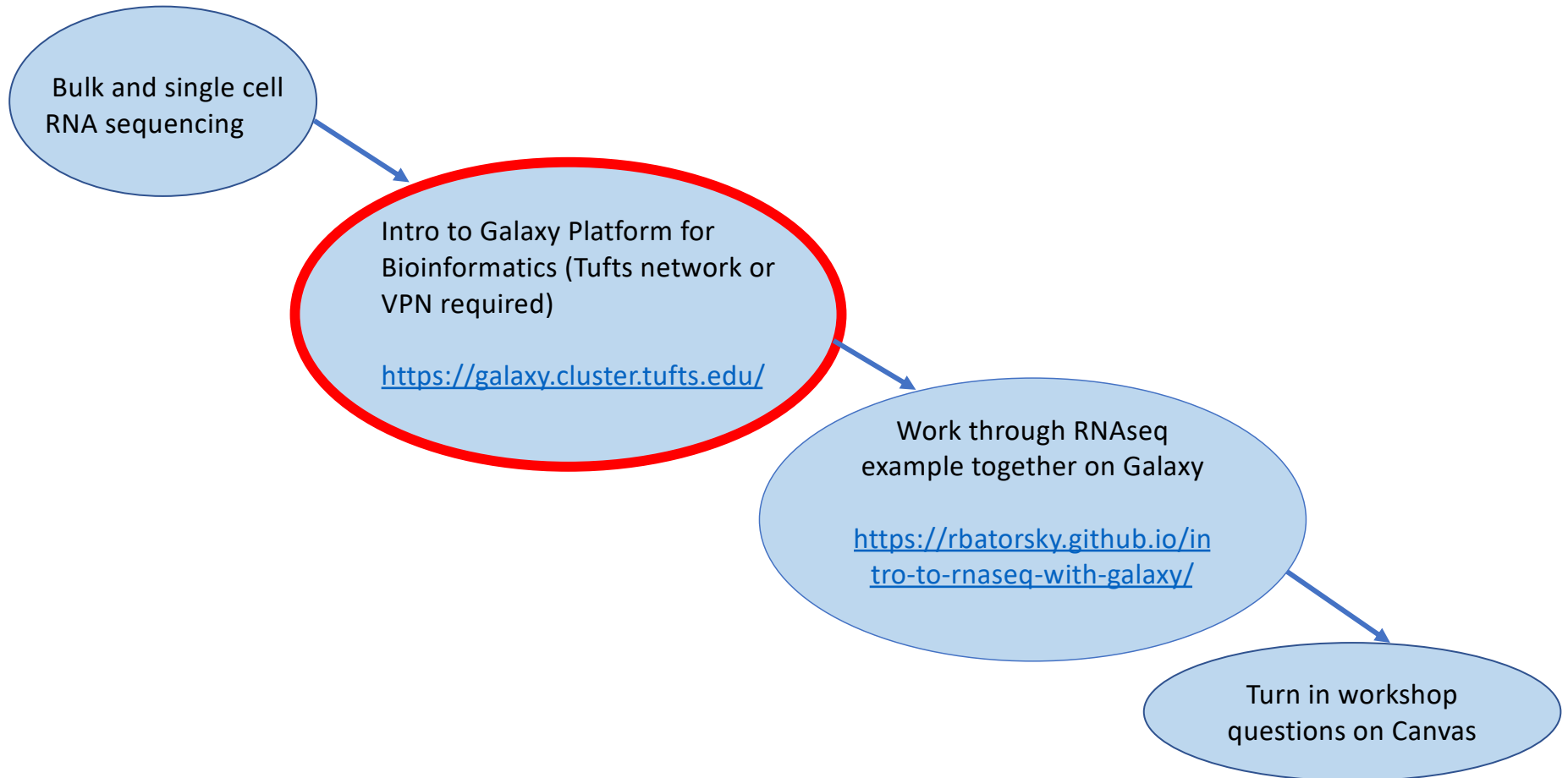


Gene	Sample 1	Sample 2	Sample 3	Sample 4
A	1000	1000	100	10
B	10	1	5	6
C	10	1	10	20



$\log_2\text{FoldChange} = 1.55$
 $\text{Adjusted } p\text{-val} < 0.05$

Outline



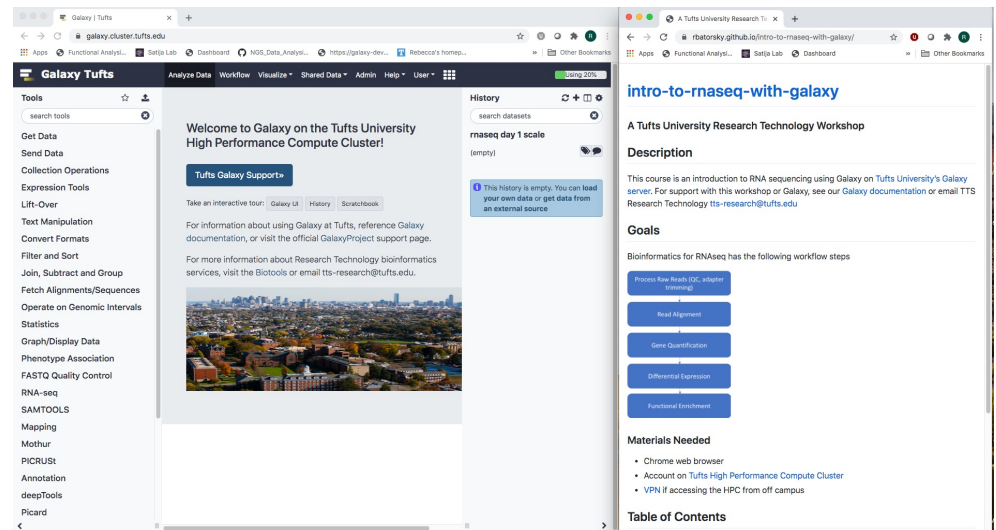


- ❖ **Web-based** platform for running data analysis and integration, geared towards bioinformatics
 - Open-source
 - Developed at Penn State, Johns Hopkins, OHSU and Cleveland Clinic with many more outside contributions
 - Large and extremely responsive community

Access Galaxy

1. Connect to Tufts Network, either on campus or via [VPN](#)
2. Visit <https://galaxy.cluster.tufts.edu/>
3. Log in with you cluster username and password
4. In another browser window go to course workflow:
<https://rbatorsky.github.io/intro-to-rnaseq-with-galaxy/>

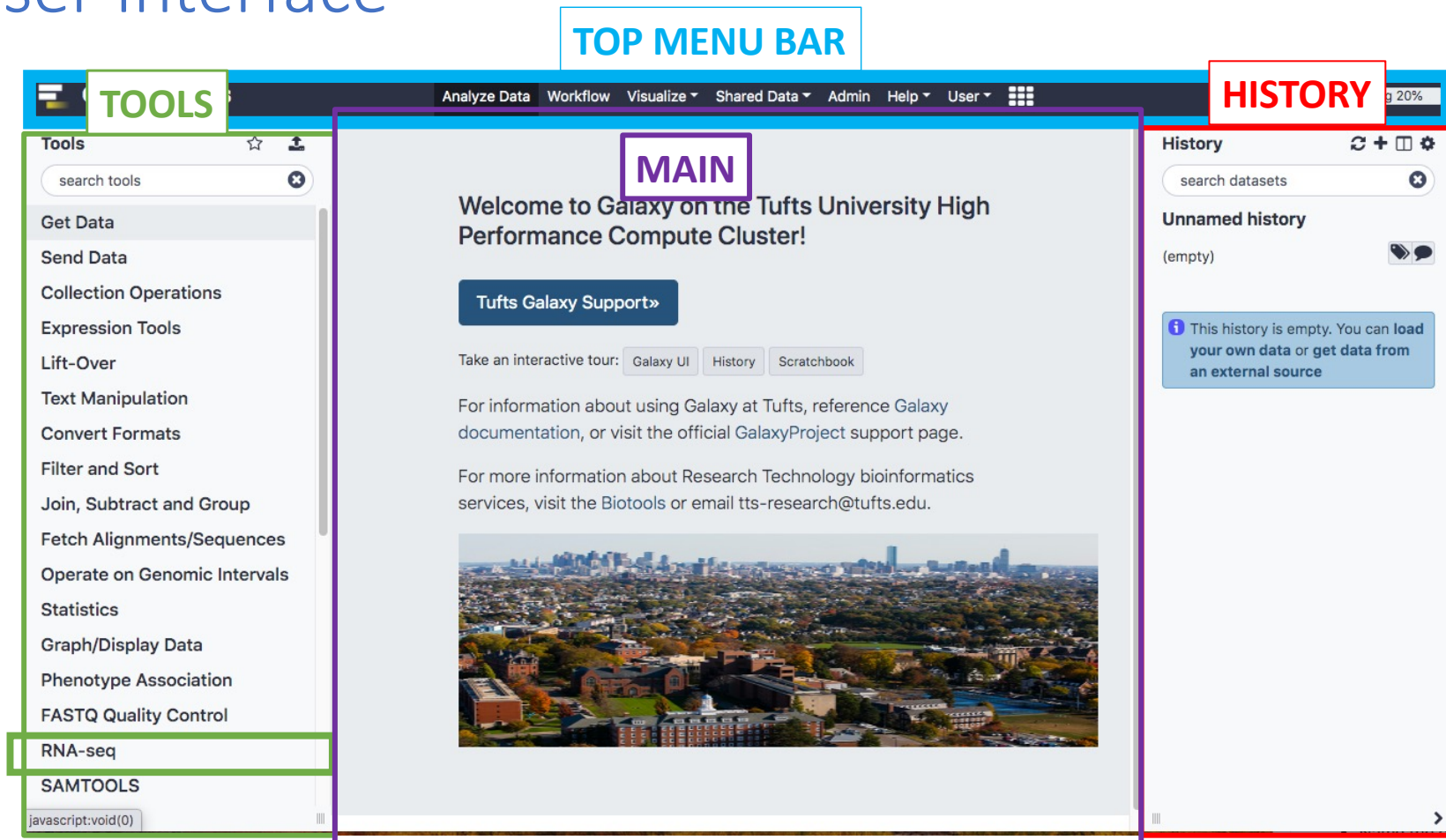
Suggested screen layout



User Interface

The screenshot displays the Galaxy Tufts web interface. At the top is a dark navigation bar with the 'Galaxy Tufts' logo and a menu containing 'Analyze Data', 'Workflow', 'Visualize', 'Shared Data', 'Admin', 'Help', 'User', and a grid icon. A status indicator on the right shows 'Using 20%'. On the left, a 'Tools' sidebar lists various categories: 'Get Data', 'Send Data', 'Collection Operations', 'Expression Tools', 'Lift-Over', 'Text Manipulation', 'Convert Formats', 'Filter and Sort', 'Join, Subtract and Group', 'Fetch Alignments/Sequences', 'Operate on Genomic Intervals', 'Statistics', 'Graph/Display Data', 'Phenotype Association', 'FASTQ Quality Control', 'RNA-seq', and 'SAMTOOLS'. The main content area features a 'Welcome to Galaxy on the Tufts University High Performance Compute Cluster!' message, a 'Tufts Galaxy Support»' button, and links for an interactive tour ('Galaxy UI', 'History', 'Scratchbook'). It also provides information on using Galaxy at Tufts and contacting the Research Technology bioinformatics services. Below the text is a wide aerial photograph of the Tufts University campus. On the right, a 'History' panel shows an empty 'Unnamed history' list with a message: 'This history is empty. You can load your own data or get data from an external source'.

User Interface



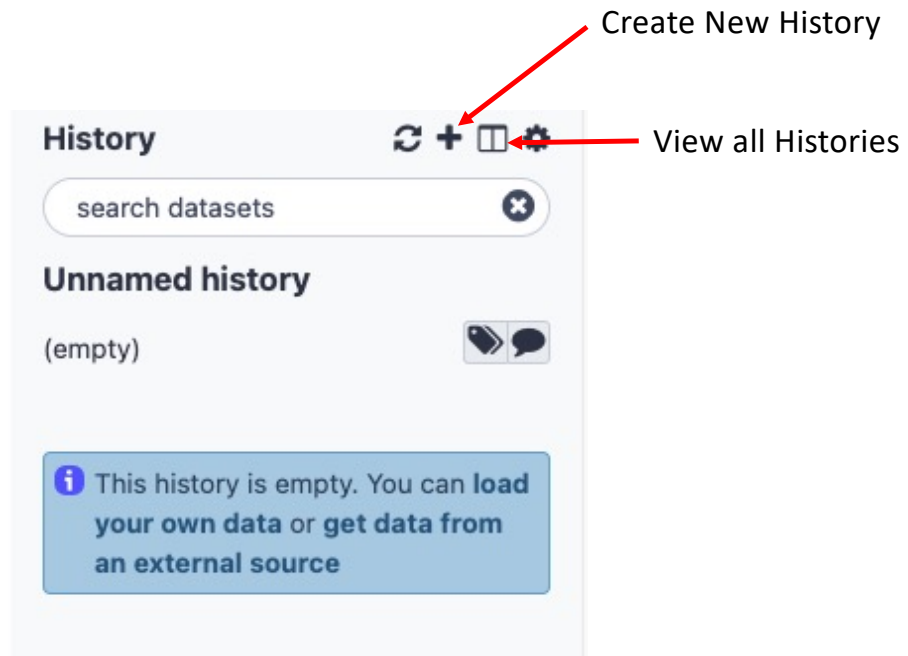
Galaxy User Interface

To return to
home screen

The screenshot displays the Galaxy Tufts web interface. The top navigation bar includes links for 'Analyze Data', 'Workflow', 'Visualize', 'Shared Data', 'Admin', 'Help', and 'User'. A red circle highlights the 'Galaxy Tufts' logo in the top left corner. On the left, a 'Tools' sidebar lists various functions such as 'Get Data', 'Send Data', 'Collection Operations', 'Expression Tools', 'Lift-Over', 'Text Manipulation', 'Convert Formats', 'Filter and Sort', 'Join, Subtract and Group', 'Fetch Alignments/Sequences', 'Operate on Genomic Intervals', 'Statistics', 'Graph/Display Data', 'Phenotype Association', 'FASTQ Quality Control', 'RNA-seq', 'SAMTOOLS', 'Mapping', 'Mothur', and 'PICRUST'. A red circle highlights the left toolbar at the bottom of this sidebar. The main content area features a 'Welcome to Galaxy on the Tufts University High Performance Compute Cluster!' message, a 'Tufts Galaxy Support»' button, and links for an interactive tour ('Galaxy UI', 'History', 'Scratchbook'). It also provides information about using Galaxy at Tufts and mentions Research Technology bioinformatics services. A cityscape image is shown below the text. On the right, a 'History' panel shows an 'Unnamed history' (empty) and a message: 'This history is empty. You can load your own data or get data from an external source'. A red circle highlights the right toolbar at the bottom of this panel. At the very bottom, four red circles highlight the minimize/adjust toolbars for the Tools sidebar, the main content area, and the History panel.

Minimize/Adjust
toolbars

History



History

Create New History

View all Histories

The image shows the Galaxy web interface's History panel. On the left, a panel titled 'History' contains a search bar, a button to 'Create New History' (indicated by a red arrow), and a button to 'View all Histories' (indicated by a red arrow). Below these is an 'Unnamed history' section, which is currently empty. A blue information box at the bottom of this section states: 'This history is empty. You can load your own data or get data from an external source'. On the right, a larger panel shows a list of recent jobs, including 'RNA-seq', 'WT_3_collection', 'WT_1_collection', 'SNF2_3_collection', 'SNF2_1_collection', and 'Concatenate datasets on data 1 and data 2'. The 'Concatenate datasets on data 1 and data 2' job is highlighted in red, and a red arrow points from the 'View all Histories' button to it.

Tools

The screenshot displays the Galaxy web interface. On the left, a sidebar lists various tool categories. The 'Tools' category is circled in red, and a green arrow points to the 'RNA-seq' category, which is also circled. Below 'RNA-seq', several tools are listed, including 'DESeq2', 'featureCounts', and 'RNA STAR'. The main area of the interface shows a 'Welcome to Galaxy on the Tufts cluster' message, a 'Bioinformatics @ Tufts' button, and a 'Take an interactive tour' section with links to 'Galaxy UI', 'History', and 'Scratchbook'. On the right, a 'History' panel shows an 'Unnamed history' which is empty. The top navigation bar includes links for 'Analyze Data', 'Workflow', 'Visualize', 'Shared Data', 'Admin', 'Help', and 'User', along with a 'Using 14.7 GB' status indicator.

Galaxy

Analyze Data Workflow Visualize Shared Data Admin Help User Using 14.7 GB

Tools

search tools

Send Data

Collection Operations

Lift-Over

Text Manipulation

Convert Formats

Filter and Sort

Join, Subtract and Group

Fetch Alignments/Sequences

Operate on Genomic Intervals

Statistics

Graph/Display Data

Phenotype Association

FASTQ Quality Control

RNA-seq

DESeq2 Determines differentially expressed features from count tables

featureCounts Measure gene expression in RNA-Seq experiments from SAM or BAM files.

RNA STAR Gapped-read mapper for RNA-seq data

SAMTOOLS

Mapping

Workflows

All workflows

Welcome to Galaxy on the Tufts cluster

Bioinformatics @ Tufts

Take an interactive tour: Galaxy UI History Scratchbook

Galaxy is an open platform for supporting data intensive research. Galaxy is developed by The Galaxy Team with the support of many contributors.

The Galaxy Project is supported in part by NHGRI, NSF, The Huck Institutes of the Life Sciences, The Institute for CyberScience at Penn State, and Johns Hopkins University.

History

search datasets

Unnamed history

(empty)

This history is empty. You can load your own data or get data from an external source

Tools

Click on the name of the tool to open it in the main panel

The screenshot displays the Galaxy web interface. On the left is a sidebar with a 'Tools' section containing a search bar and a list of tool categories: Get Data, Send Data, Collection Operations, Lift-Over, Text Manipulation, Convert Formats, Filter and Sort, Join, Subtract and Group, Fetch Alignments/Sequences, Operate on Genomic Intervals, Statistics, Graph/Display Data, Phenotype Association, FASTQ Quality Control, RNA-seq, RNA STAR, SAMTOOLS, Mapping, and Workflows. The 'featureCounts' tool is highlighted in the 'RNA-seq' category. An arrow points from the text 'Click on the name of the tool to open it in the main panel' to the 'featureCounts' tool name. The main panel shows the 'featureCounts' tool configuration page, version 1.6.4. It includes sections for 'Alignment file' (with a file selection button and a message 'No bam or sam dataset available.'), 'Specify strand information' (set to 'Unstranded'), 'Gene annotation file' (set to 'locally cached'), 'Using locally cached annotation' (set to 'No options available'), 'Output format' (set to 'Gene-ID "\t" read-count (MultiQC/DESeq2/edgeR/limma-voom compatible)'), and 'Create gene-length file' (with 'Yes' and 'No' buttons). There are also links for 'Options for paired-end reads' and 'Advanced options'. At the bottom is an 'Execute' button. The right sidebar shows the 'History' section, which is currently empty, with a message: 'This history is empty. You can load your own data or get data from an external source'.

Importing data

The screenshot displays the Galaxy web interface. The top navigation bar includes links for Analyze Data, Workflow, Visualize, Shared Data, Admin, Help, and User. The 'Shared Data' link is highlighted with a red box and an arrow pointing to it with the text 'Import shared data libraries'. On the left sidebar, under the 'Tools' section, there is a search bar and a list of tool categories. A red box highlights the 'Upload data from local storage or from the cluster' option, with an arrow pointing to it from the text 'Upload data from local storage or from the cluster'. The main content area shows a welcome message for 'Galaxy on the Tufts cluster' and a button for 'Bioinformatics @ Tufts'. The right sidebar shows the 'History' section, which is currently empty, with a message indicating that the history is empty and suggesting ways to load data.

Galaxy

Analyze Data Workflow Visualize Shared Data Admin Help User

Using 14.7 GB

Tools

search tools

Get Data

Send Data

Collection Operations

Lift-Over

Text Manipulation

Convert Formats

Filter and Sort

Join, Subtract and Group

Fetch Alignments/Sequences

Operate on Genomic Intervals

Statistics

Graph/Display Data

Phenotype Association

FASTQ Quality Control

RNA-seq

SAMTOOLS

Mapping

Workflows

All workflows

Welcome to Galaxy on the Tufts cluster

Bioinformatics @ Tufts

Take an interactive tour: Galaxy UI History Scratchbook

Galaxy is an open platform for supporting data intensive research. Galaxy is developed by The Galaxy Team with the support of many contributors.

The Galaxy Project is supported in part by NHGRI, NSF, The Huck Institutes of the Life Sciences, The Institute for CyberScience at Penn State, and Johns Hopkins University.

History

search datasets

Unnamed history

(empty)

This history is empty. You can load your own data or get data from an external source

17

Access Galaxy

1. Connect to Tufts Network, either on campus or via [VPN](#)

2. Visit <https://galaxy.cluster.tufts.edu/>

3. Log in with you cluster username and password

4. In another browser window go to course workflow:

<https://rbatorsky.github.io/intro-to-rnaseq-with-galaxy/>

5. Under Table of Contents click on “**Process Raw Reads**”

Suggested screen layout

The image shows two side-by-side browser windows. The left window displays the Galaxy Tufts interface at <https://galaxy.cluster.tufts.edu/>. It features a sidebar with various tool categories like 'Tools', 'Get Data', 'Collection Operations', etc. The main content area has a 'Welcome to Galaxy on the Tufts University High Performance Compute Cluster!' message and a 'Tufts Galaxy Support' button. The right window shows a course page titled 'intro-to-rnaseq-with-galaxy' from the 'A Tufts University Research Technology Workshop'. It includes a 'Description' section, 'Goals', a list of 'Bioinformatics for RNAseq' workflow steps (Process Raw Reads, Read Alignment, Gene Quantification, Differential Expression, Functional Enrichment), 'Materials Needed' (Chrome web browser, Account on Tufts High Performance Compute Cluster, VPN), and a 'Table of Contents' section.