



# RNA-seq to study HIV Infection in cells

Rebecca Batorsky  
Sr Bioinformatics Scientist  
March 2021

# Research Technology Team



**Delilah Maloney**

High Performance Computing Specialist



**Kyle Monahan**

Senior Data Science Specialist



**Shawn Doughty**

Manager, Research Computing



**Rebecca Batorsky**

Senior Bioinformatics Scientist



**Meg Farley**

Bioinformatics Intern



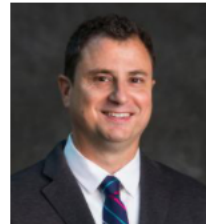
**Chris Barnett**

Senior Geospatial Analyst



**Tom Phimmasen**

Senior Data Consultant



**Patrick Florance**

Director, Academic Data Services



**Jake Perl**

Digital Humanities NLP Specialist



**Carolyn Talmadge**

Senior GIS Specialist



**Uku-Kaspar Uustalu**

Data Science Specialist

- ✓ Consultation on Projects and Grants
- ✓ High Performance Compute Cluster
- ✓ Workshops

<https://it.tufts.edu/research-technology>

# Outline

Bulk and single cell  
RNA sequencing

```
graph TD; A([Bulk and single cell RNA sequencing]) --> B([Intro to Galaxy Platform for Bioinformatics (Tufts network or VPN required)  
https://galaxy.cluster.tufts.edu/]); B --> C([Work through RNAseq example together on Galaxy  
https://rbatorsky.github.io/intro-to-rnaseq-with-galaxy/]);
```

Intro to Galaxy Platform for  
Bioinformatics (Tufts network or  
VPN required)

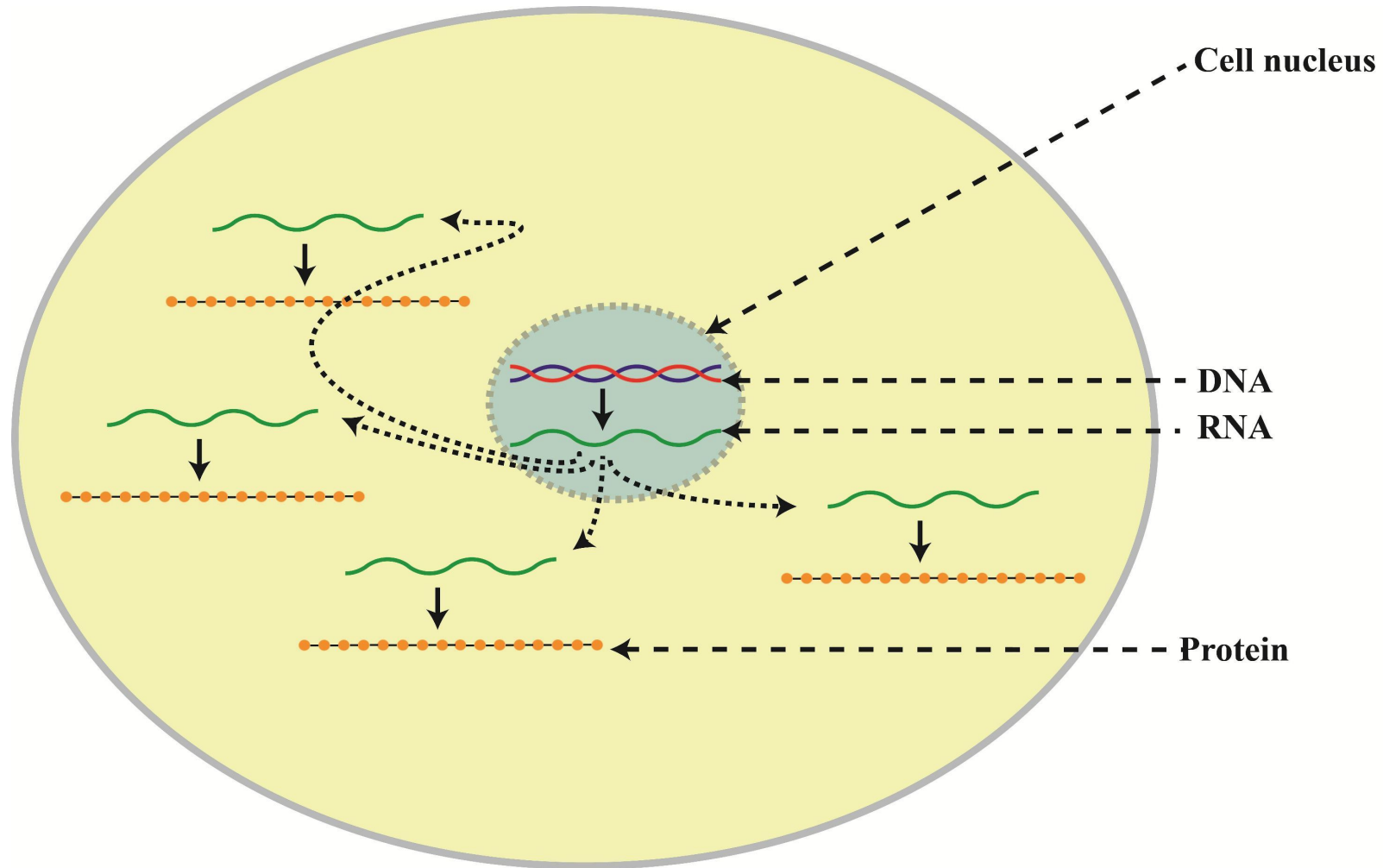
<https://galaxy.cluster.tufts.edu/>

Work through RNAseq  
example together on Galaxy

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

2 days!

# DNA and RNA in a cell

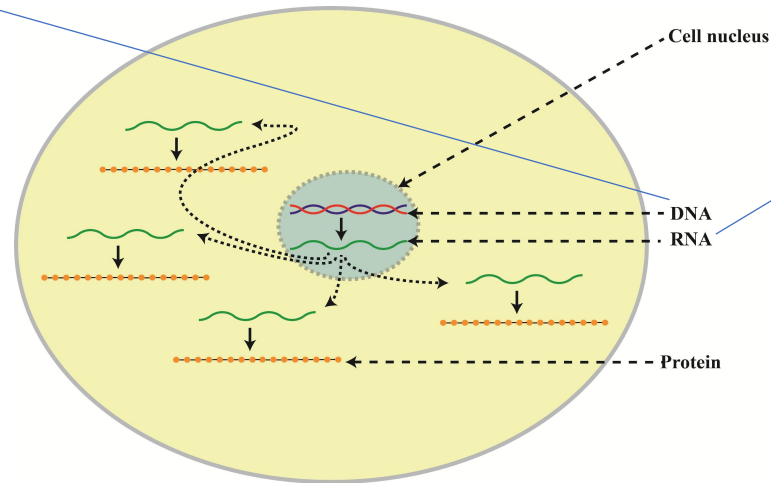




# Two common analyses

## DNA Sequencing

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



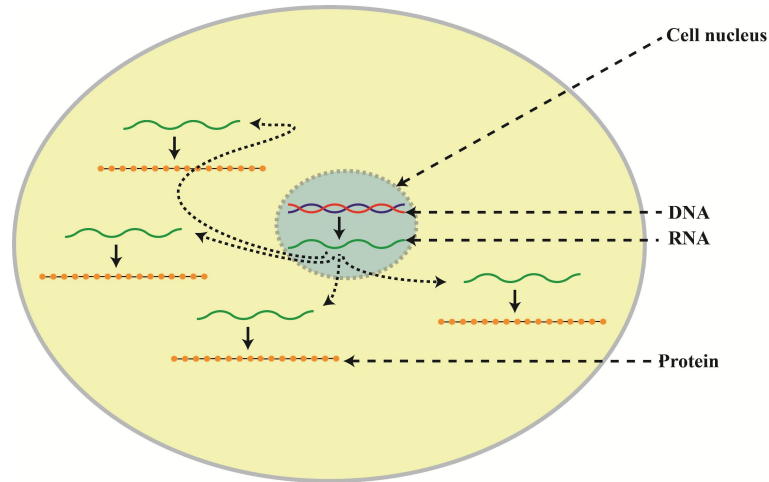
## RNA Sequencing

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

# Today we will cover RNA sequencing

## DNA Sequencing

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



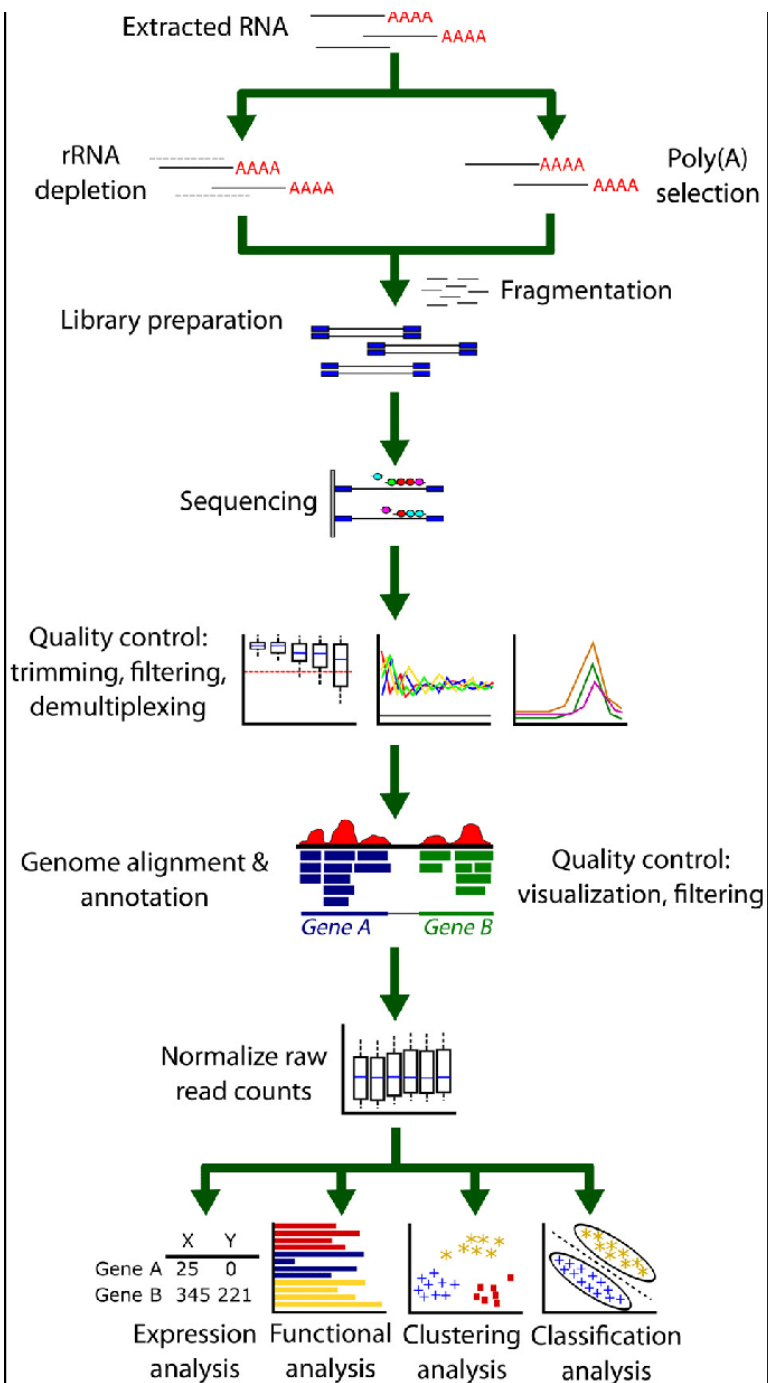
## RNA Sequencing

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

# “Bulk” RNA seq workflow

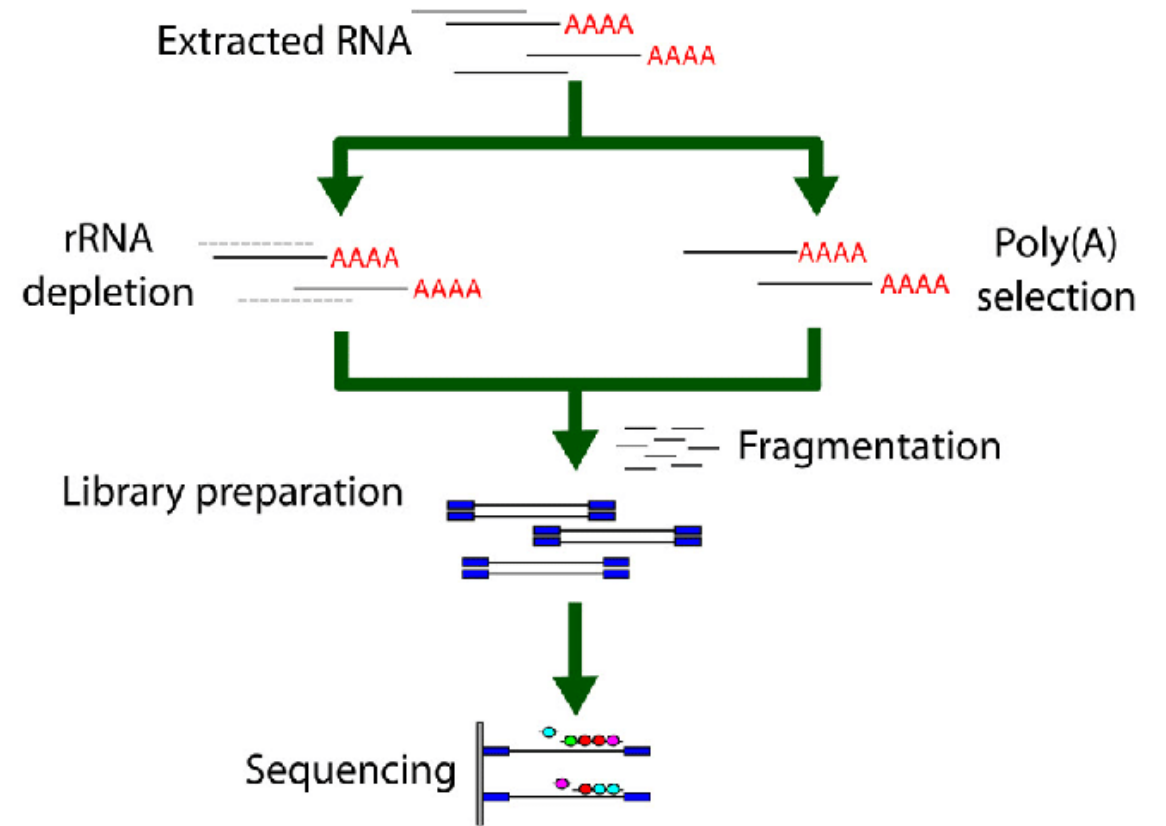
Library prep and sequencing

Bioinformatics



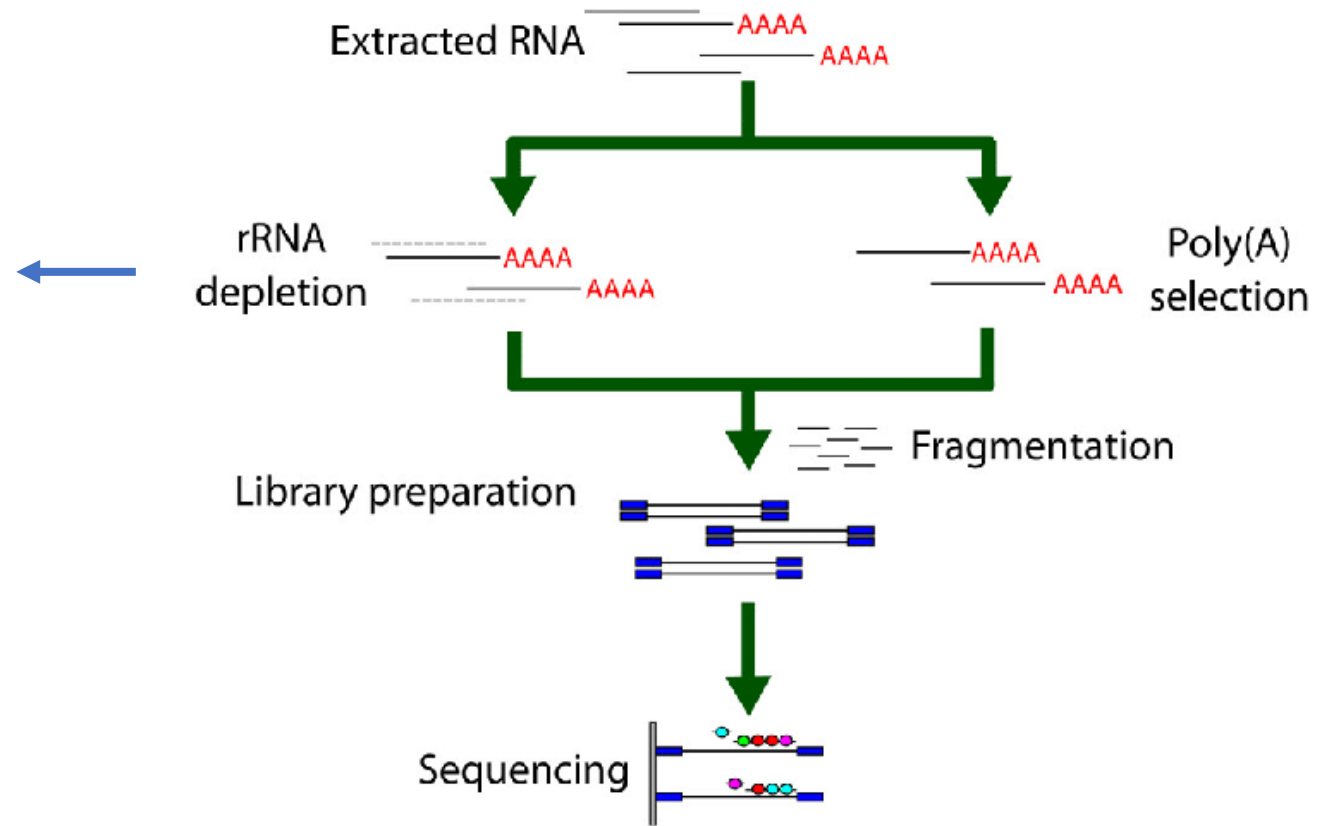
Good resource: [Griffiths et al Plos Comp Bio 2015](#)

# RNA seq library prep and sequencing



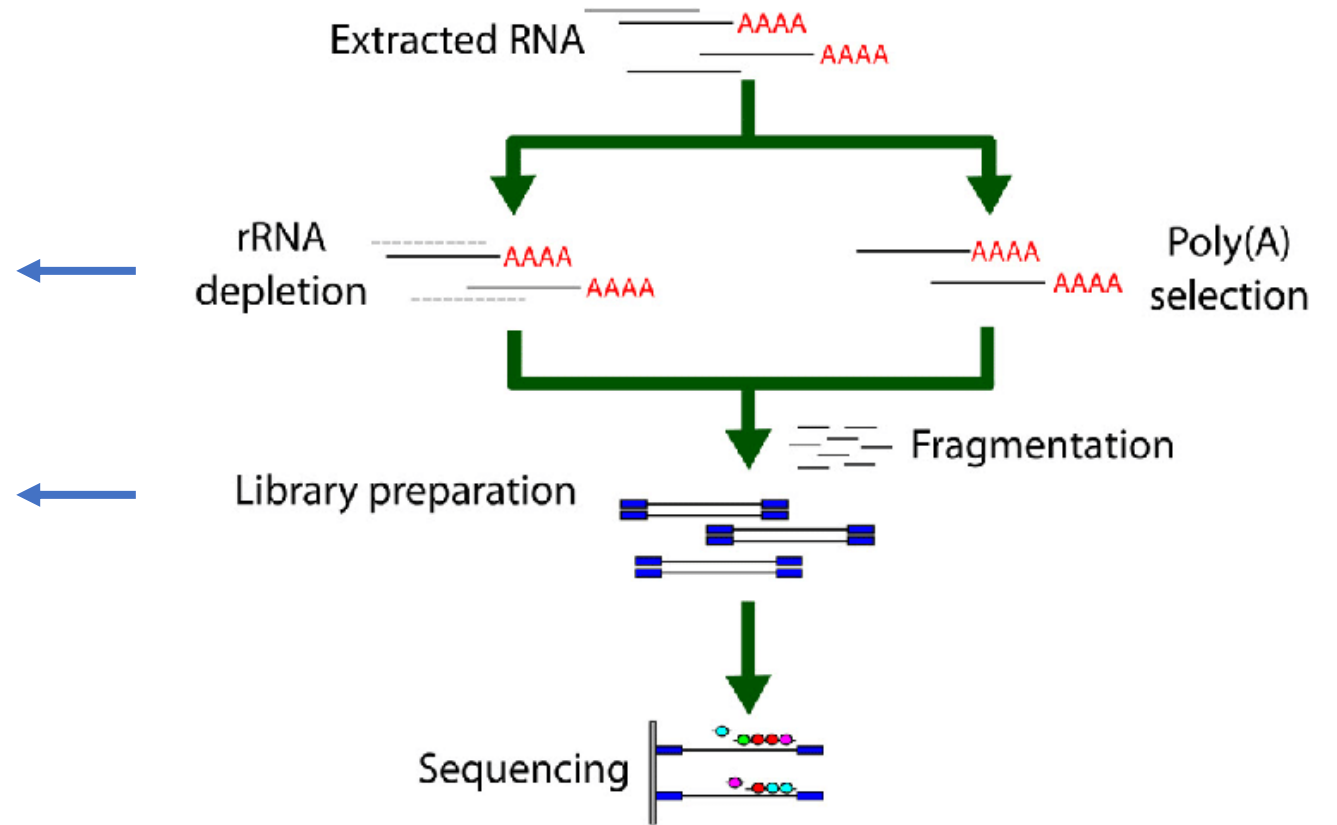
# RNA seq library prep and sequencing

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



# RNA seq library prep and sequencing

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

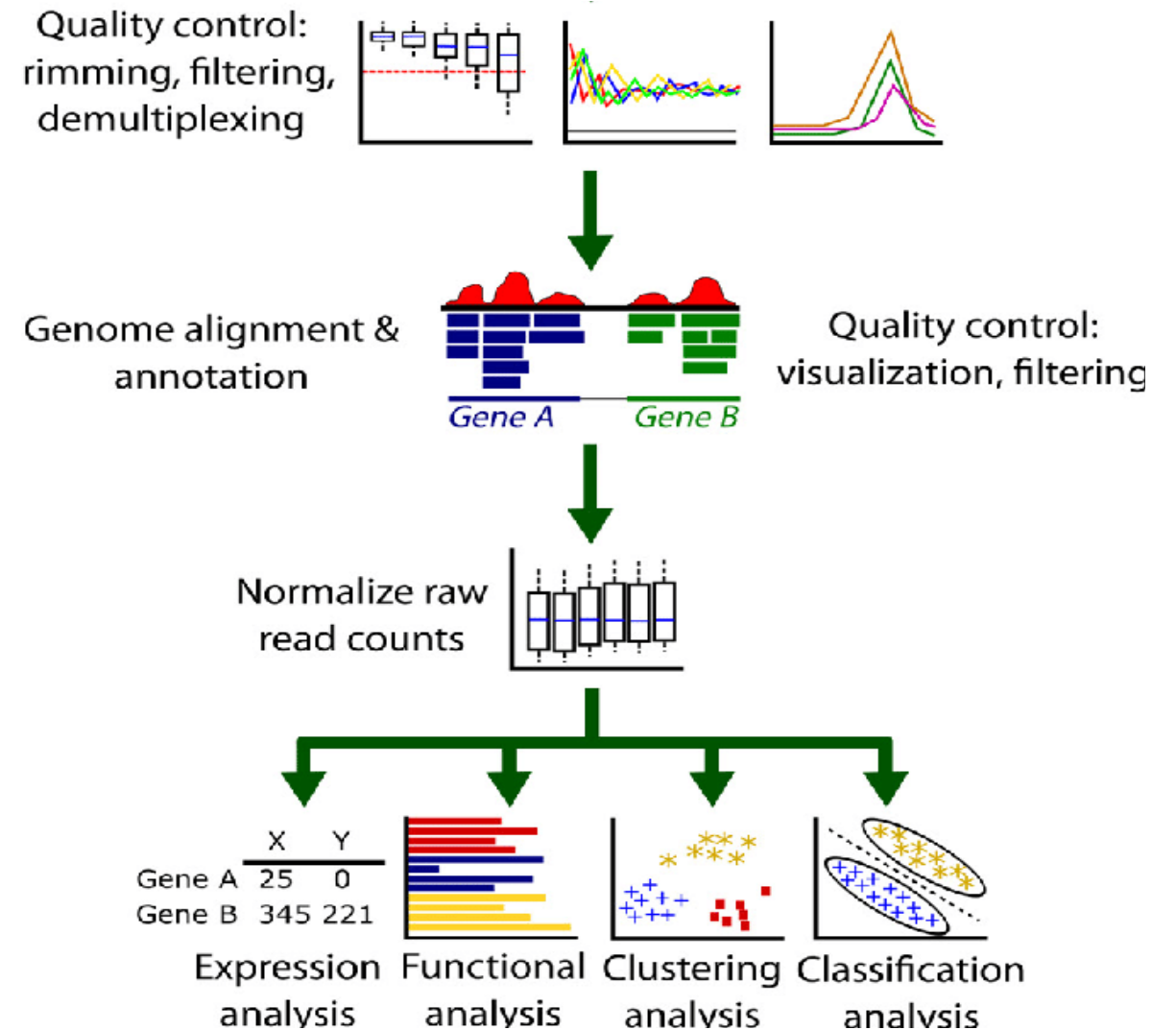


Resources:

[Illumina Sequencing by Synthesis](#)

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

# RNA seq bioinformatics



# Goal of Differential Expression in RNAseq

“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

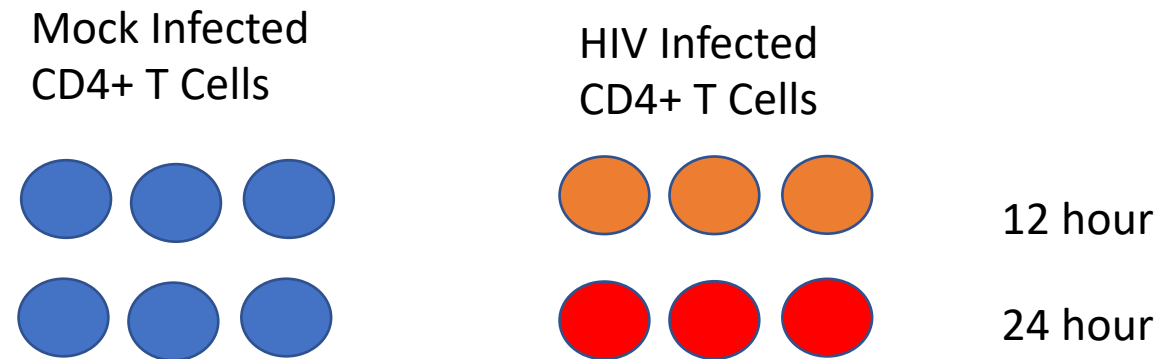
<http://genomebiology.com/2010/11/12/220>



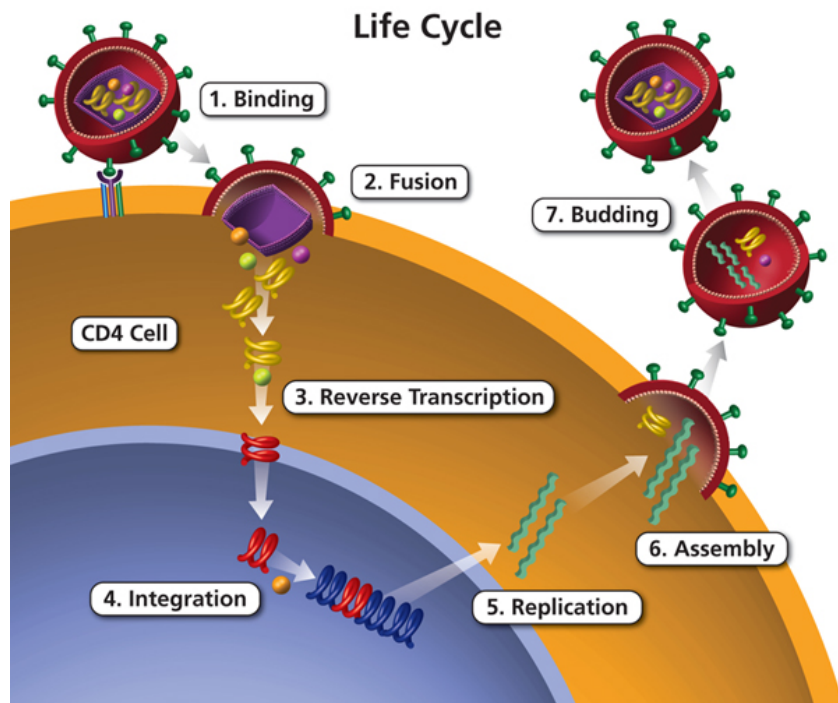
# Our dataset

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

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

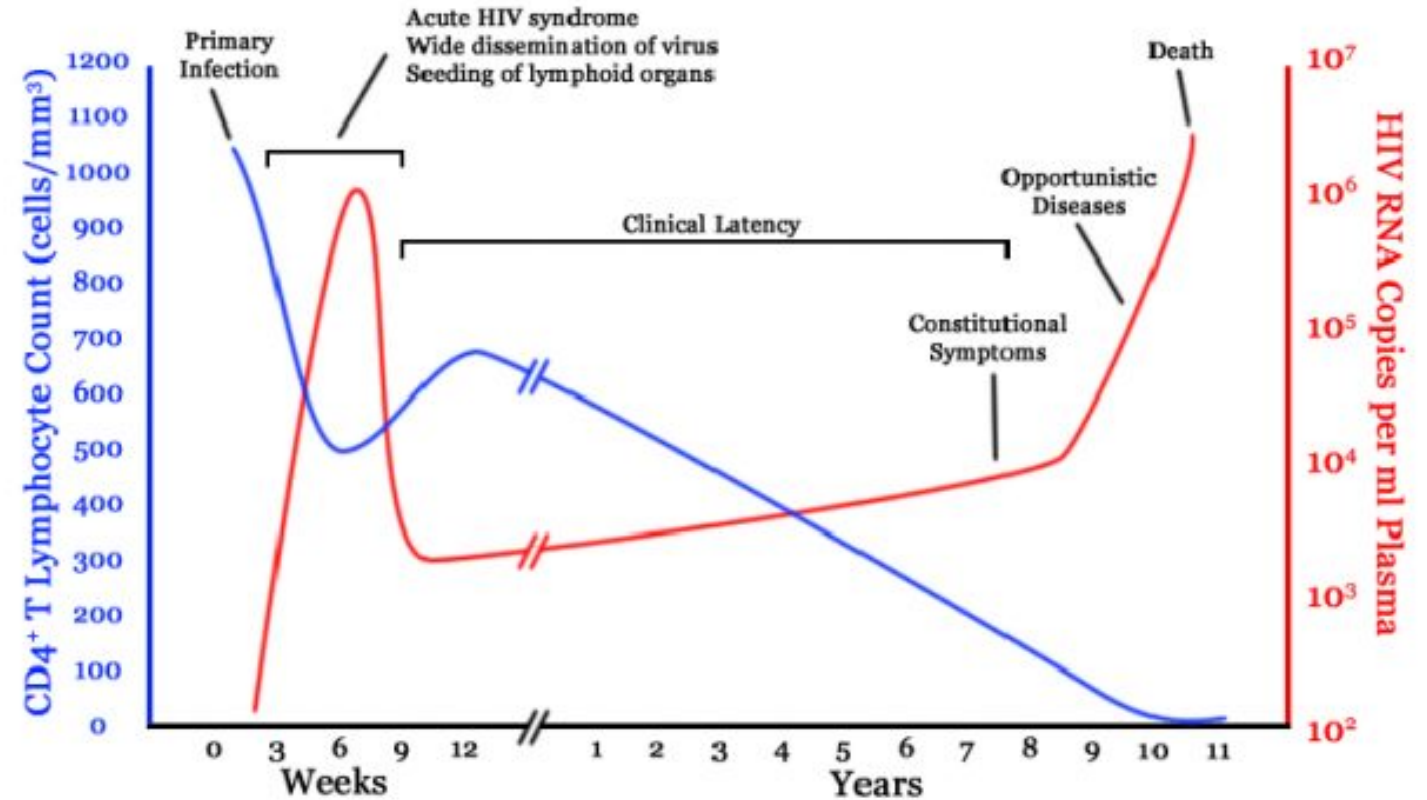
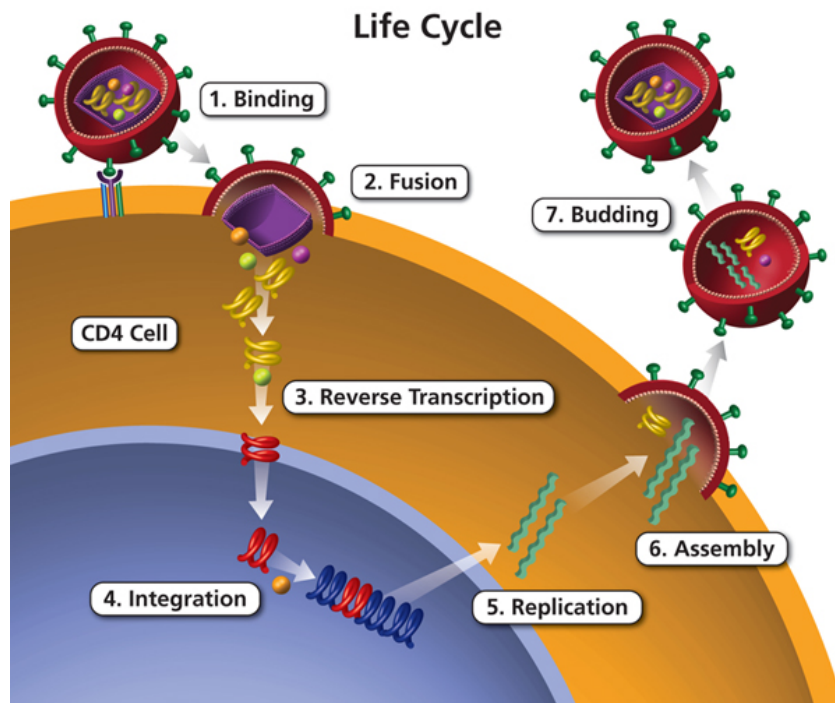


# HIV lifecycle



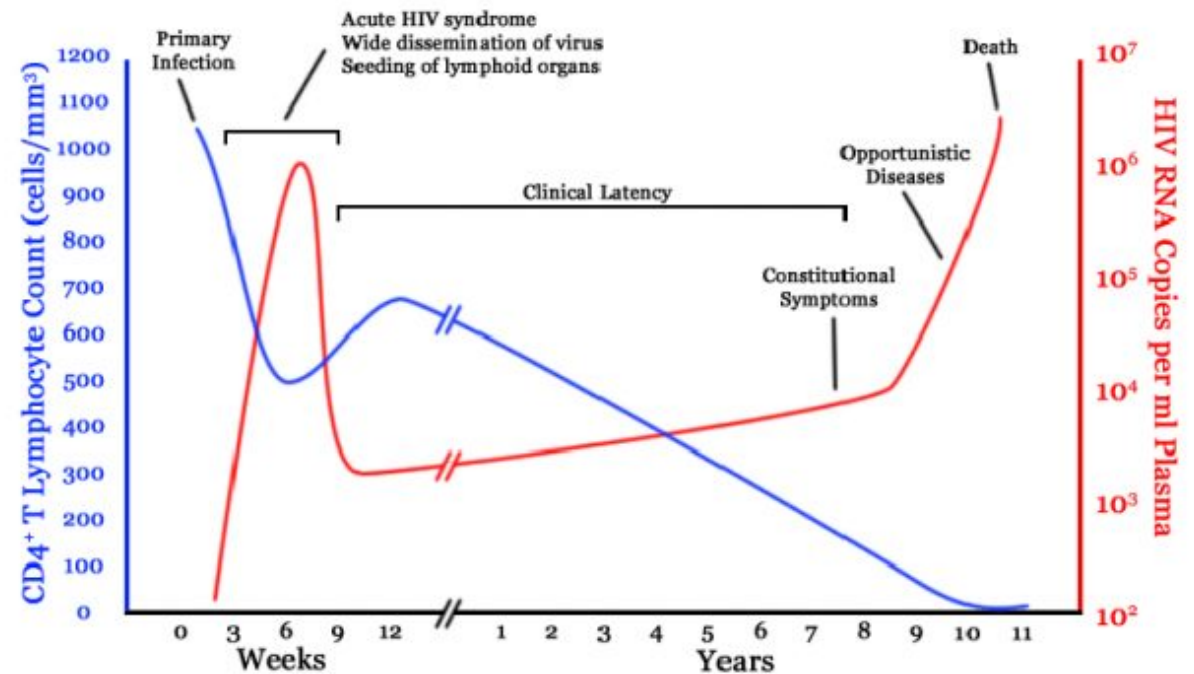
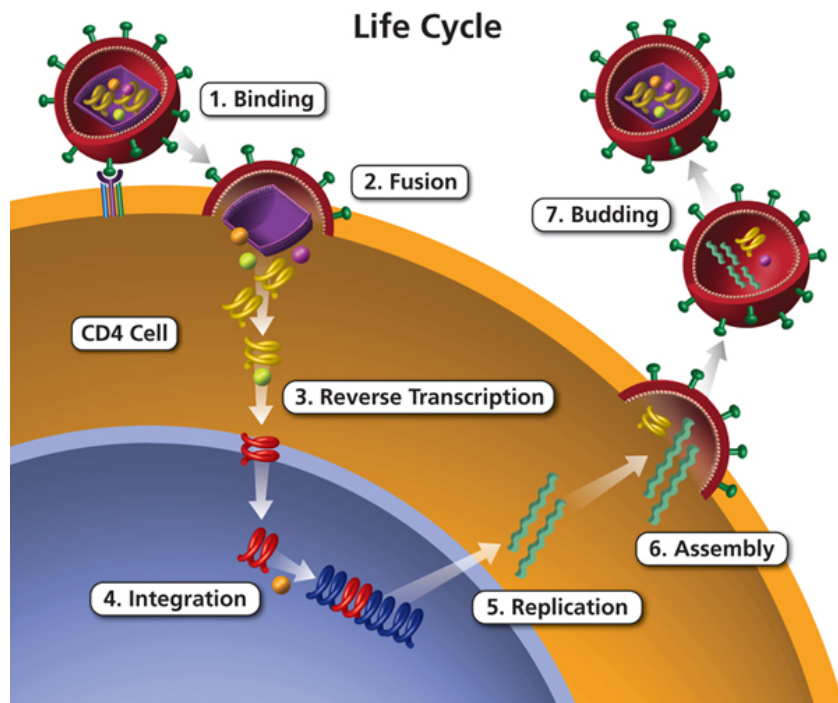
# HIV lifecycle

## HIV infection in a human host



# The study question

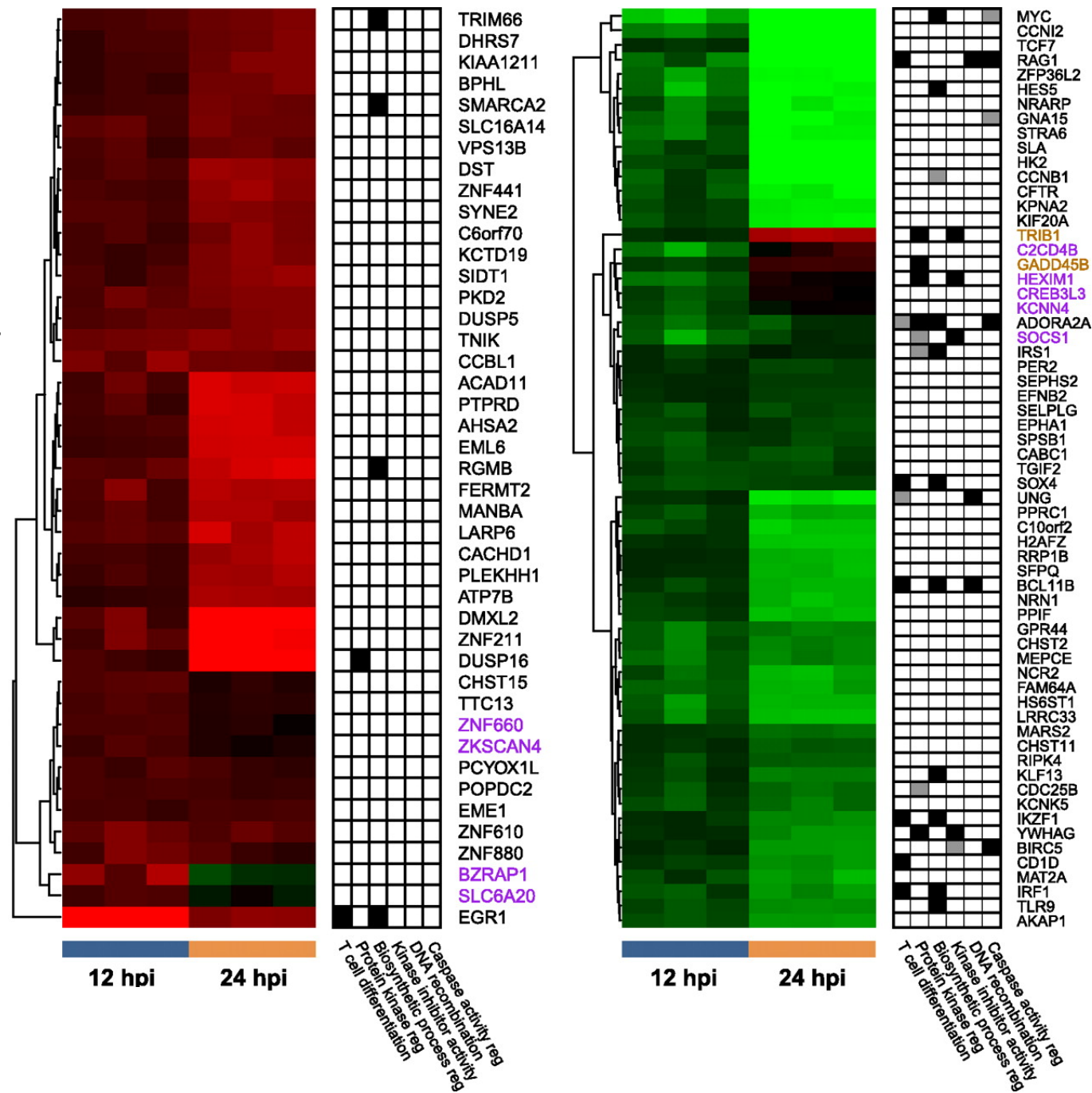
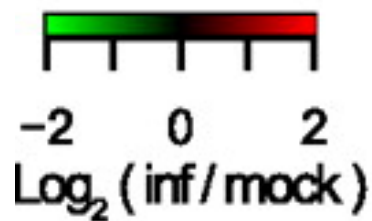
What changes take place in the first 12-24 hours of HIV infection in terms of gene expression of host cell and viral replication levels?



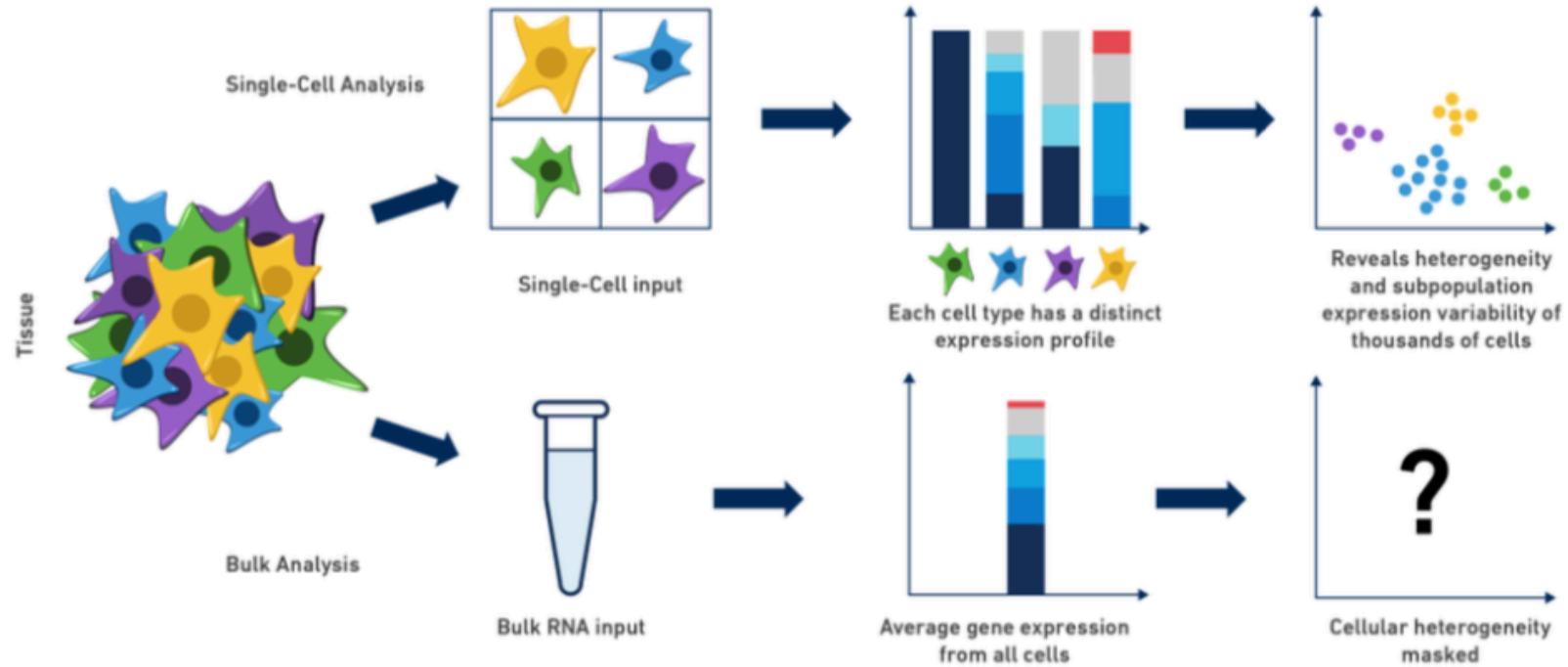
# Study findings

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

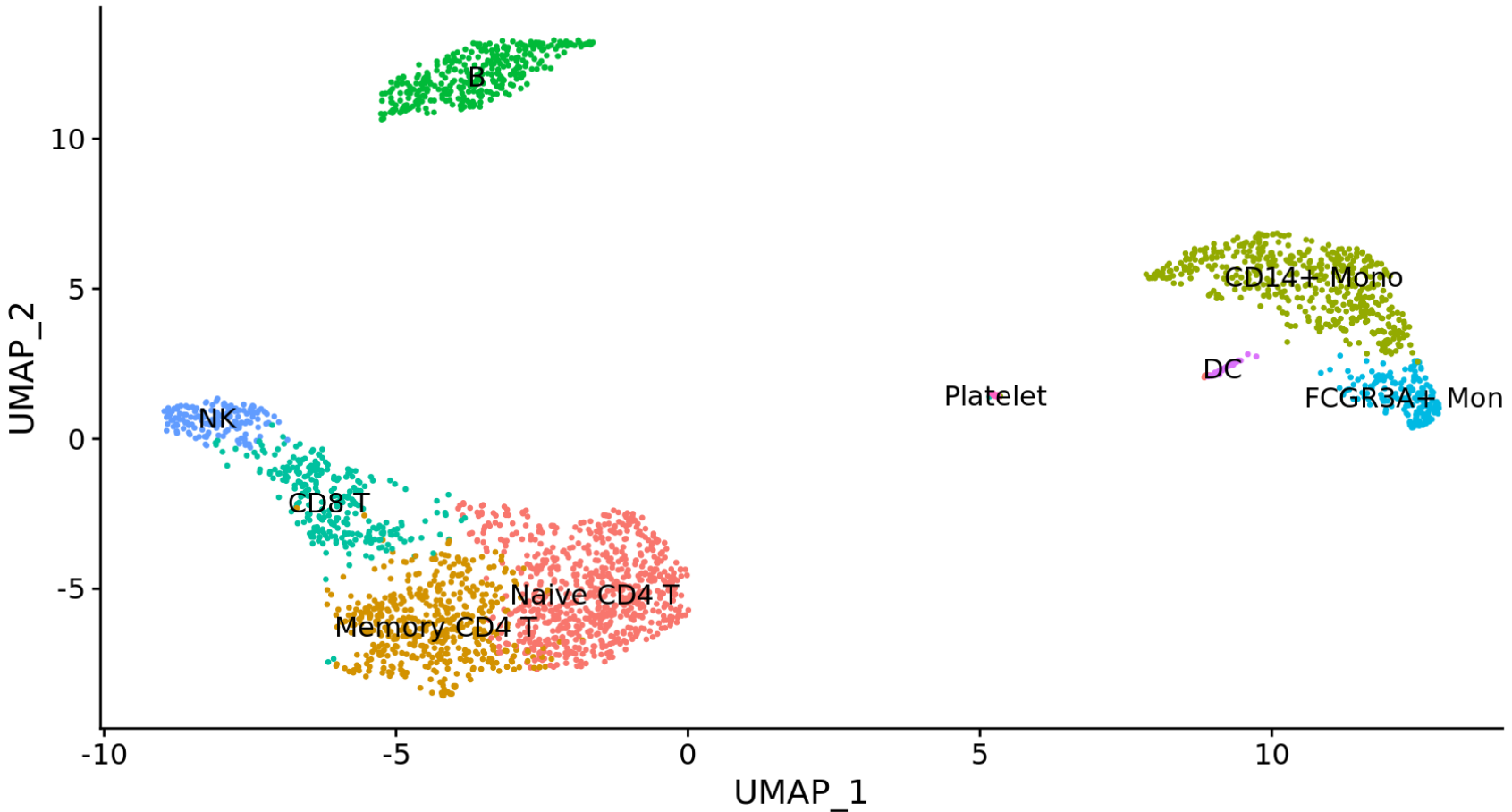


# Bulk vs Single Cell RNA Sequencing

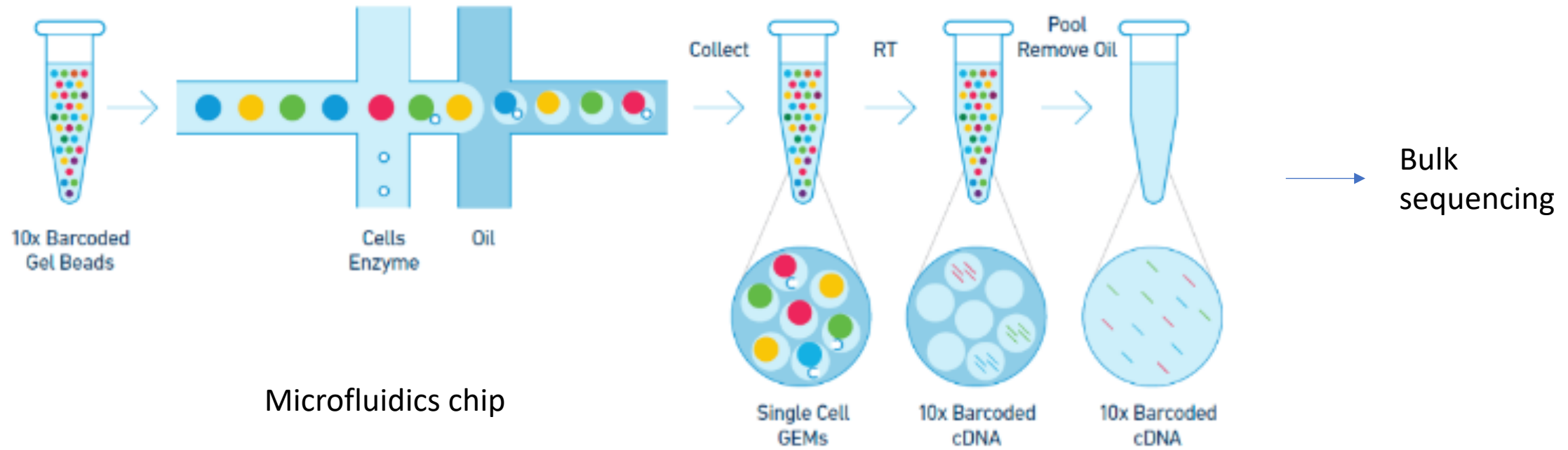




# scRNA cell subsets in PBMC

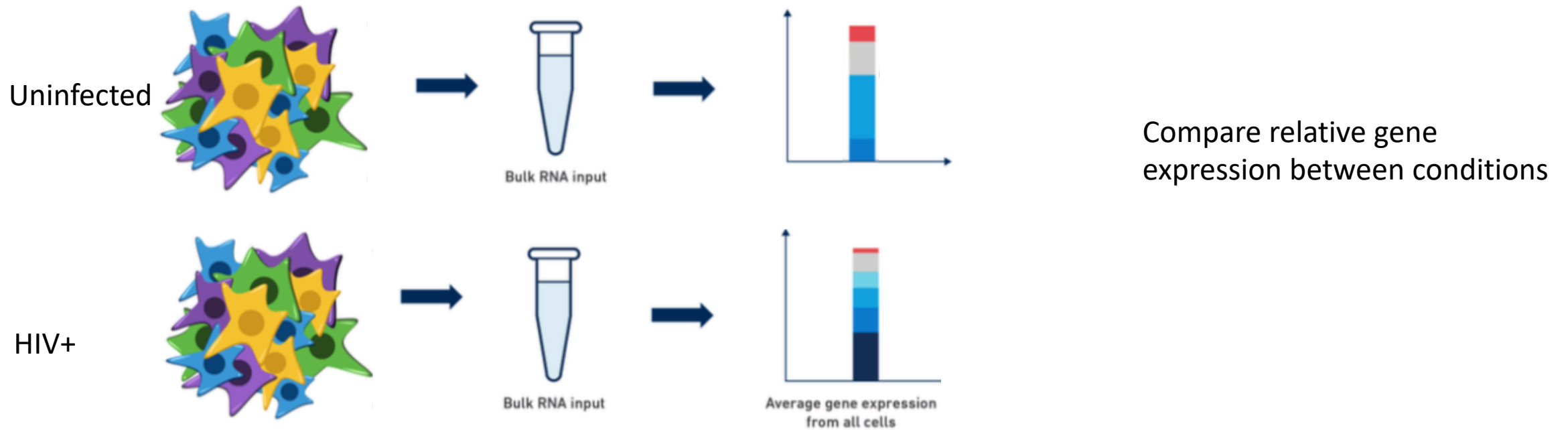


# 10x single cell technology

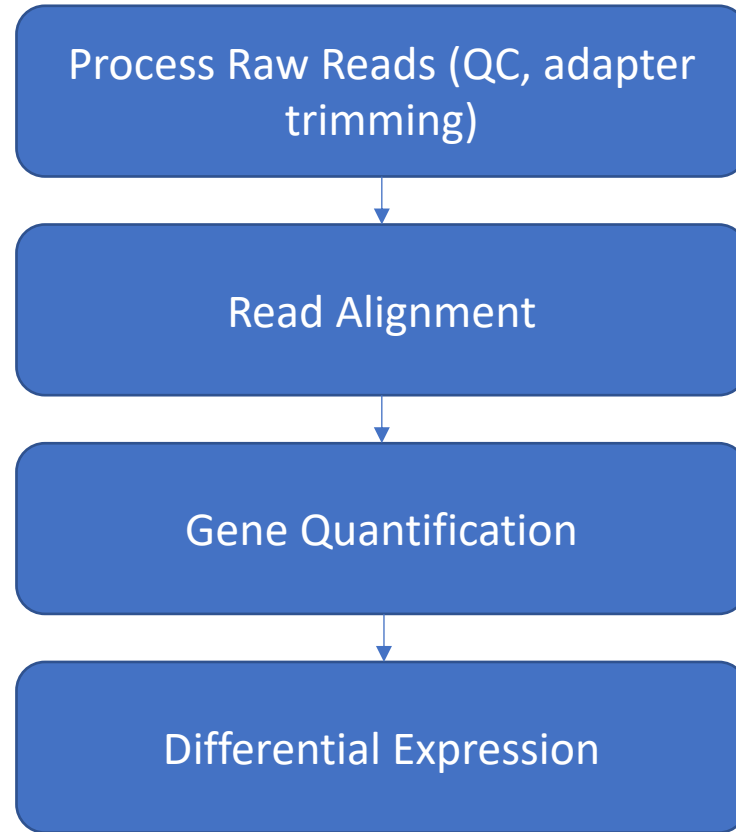




# Bulk RNAseq for Differential Expression is OK!



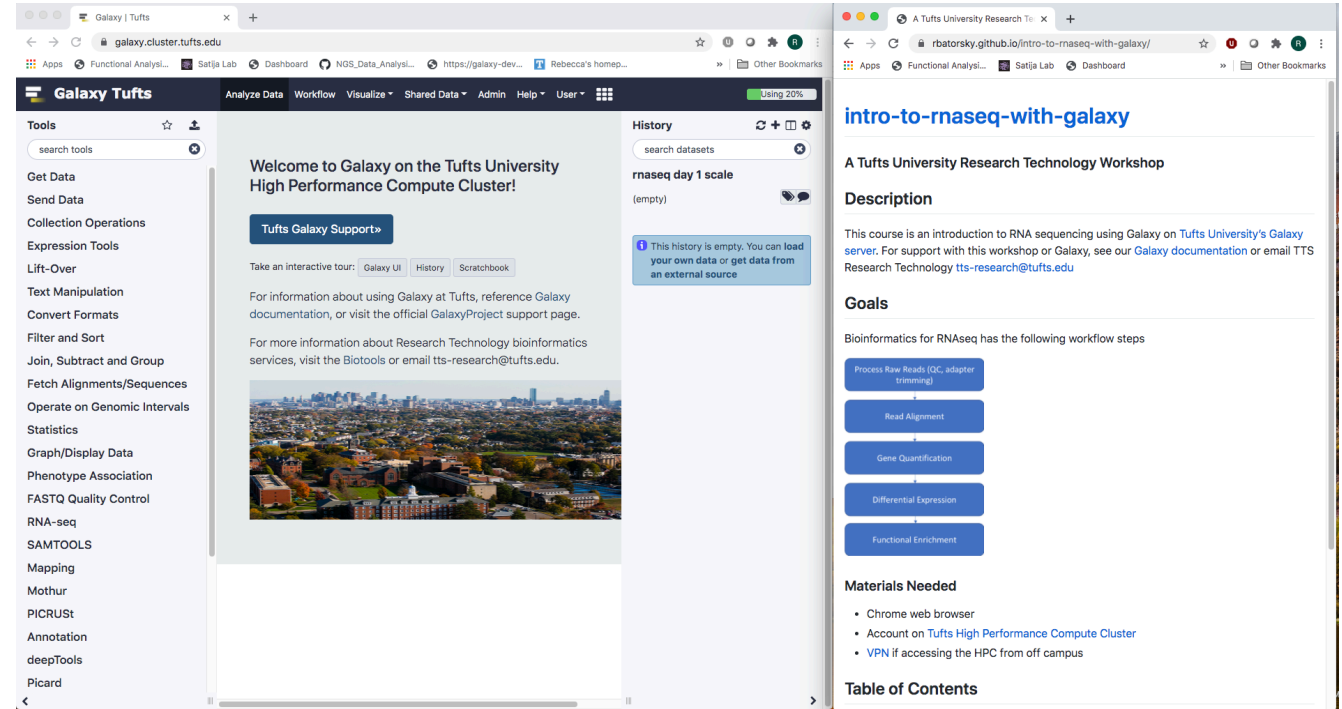
# Our (bulk) RNAseq Workflow



# Access Galaxy

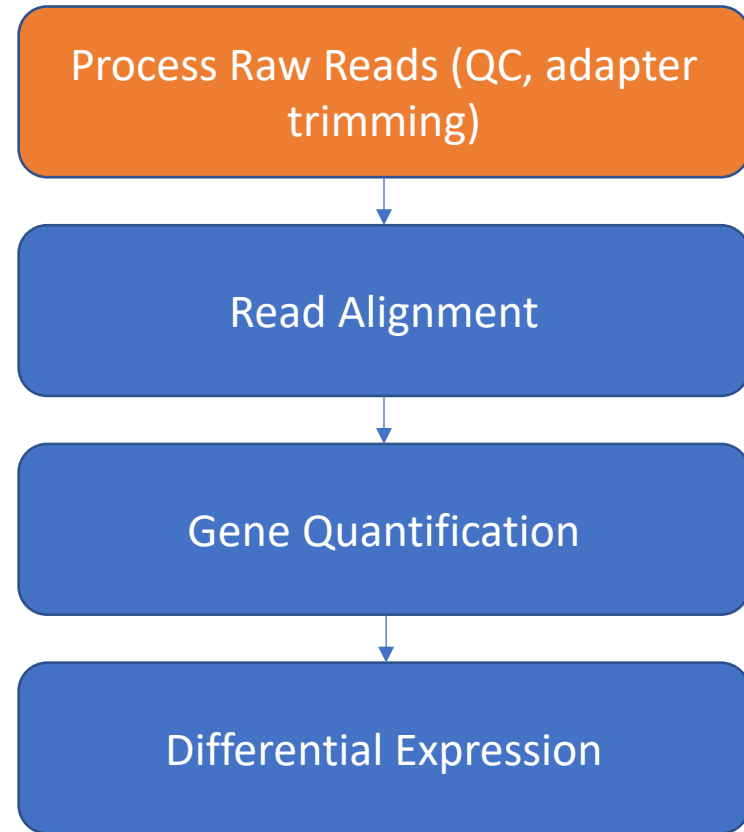
1. Connect to Tufts Network, either on campus or via [VPN](#)
2. Visit <https://galaxy.cluster.tufts.edu/>
3. Log in with you cluster username and password
4. In another browser window go to course workflow:  
<https://rbatorsky.github.io/intro-to-rnaseq-with-galaxy/>

## Suggested screen layout



The image shows two browser windows side-by-side. The left window displays the Galaxy Tufts interface. The top navigation bar includes 'Galaxy Tufts' and menu items like 'Analyze Data', 'Workflow', 'Visualize', 'Shared Data', 'Admin', 'Help', and 'User'. A left sidebar lists various tool categories such as 'Tools', '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', 'PICRUST', 'Annotation', 'deepTools', and 'Picard'. 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). Below this is a cityscape image. The right window shows a course workflow page titled 'intro-to-rnaseq-with-galaxy'. It includes a description of the course, a 'Goals' section with a workflow diagram, and a 'Materials Needed' section listing requirements like a Chrome browser, Tufts HPC account, and VPN access. The workflow diagram consists of five blue buttons: 'Process Raw Reads (QC, adapter trimming)', 'Read Alignment', 'Gene Quantification', 'Differential Expression', and 'Functional Enrichment', connected by downward arrows.

# Quality control on Raw Reads



# Raw reads in Fastq format

```
@SRR098401.109756285  
GACTCACGTAAC TTTAAACTCTAACAGAAATATACTA...  
+  
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
GACTCACGTA ACTTTAACTCTAACAGAAATATACTA...
+
CAEFGDG?BCGGGEEDGGHGHGDFHEIEGGDDDD...
```

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

# Base Quality Scores



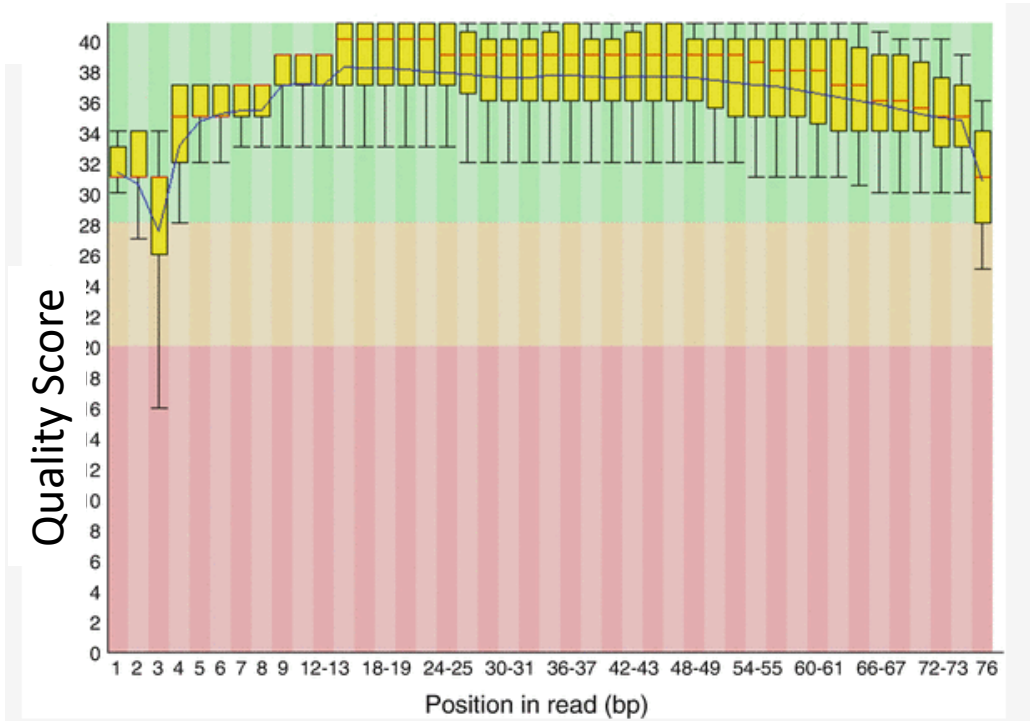
S - Sanger Phred+33, raw reads typically (0, 40)  
X - Solexa Solexa+64, raw reads typically (-5, 40)  
I - Illumina 1.3+ Phred+64, raw reads typically (0, 40)  
J - Illumina 1.5+ Phred+64, raw reads typically (3, 41)  
with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (bold)  
(Note: See discussion above).  
L - Illumina 1.8+ Phred+33, raw reads typically (0, 41)



# Raw read quality control

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

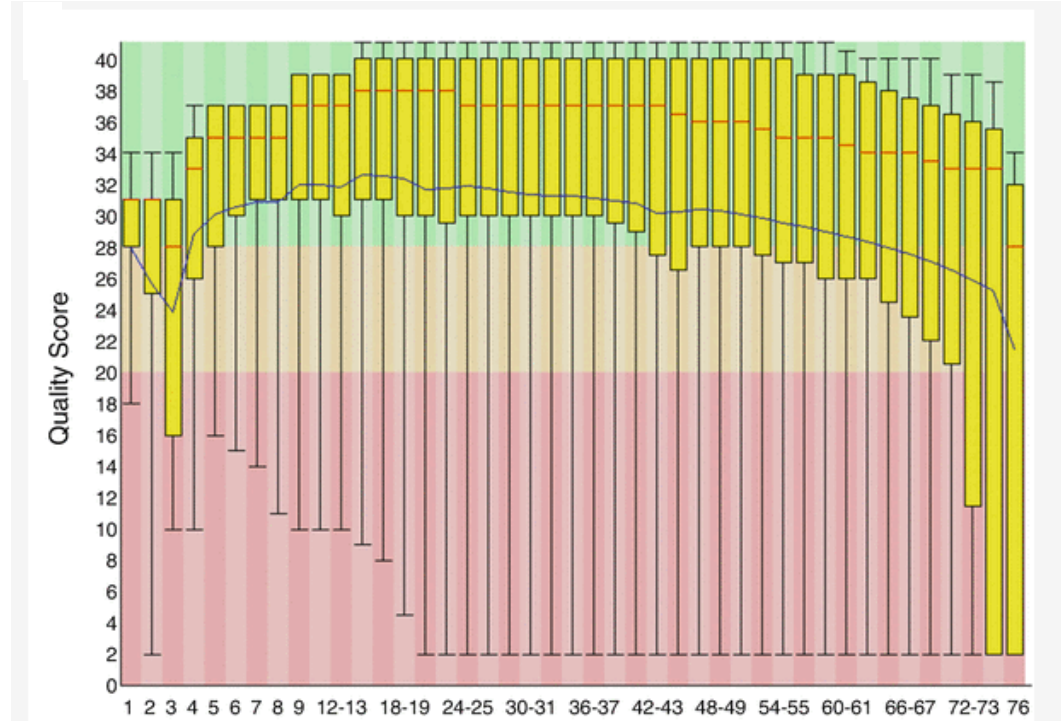
# FastQC: Sequence Quality Histogram



Position in read (bp)

GOOD

High quality over the length of the read



Position in read (bp)

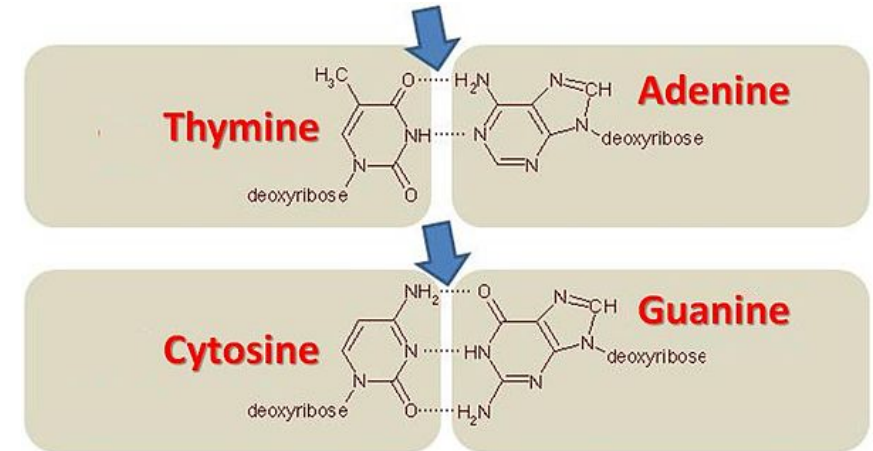
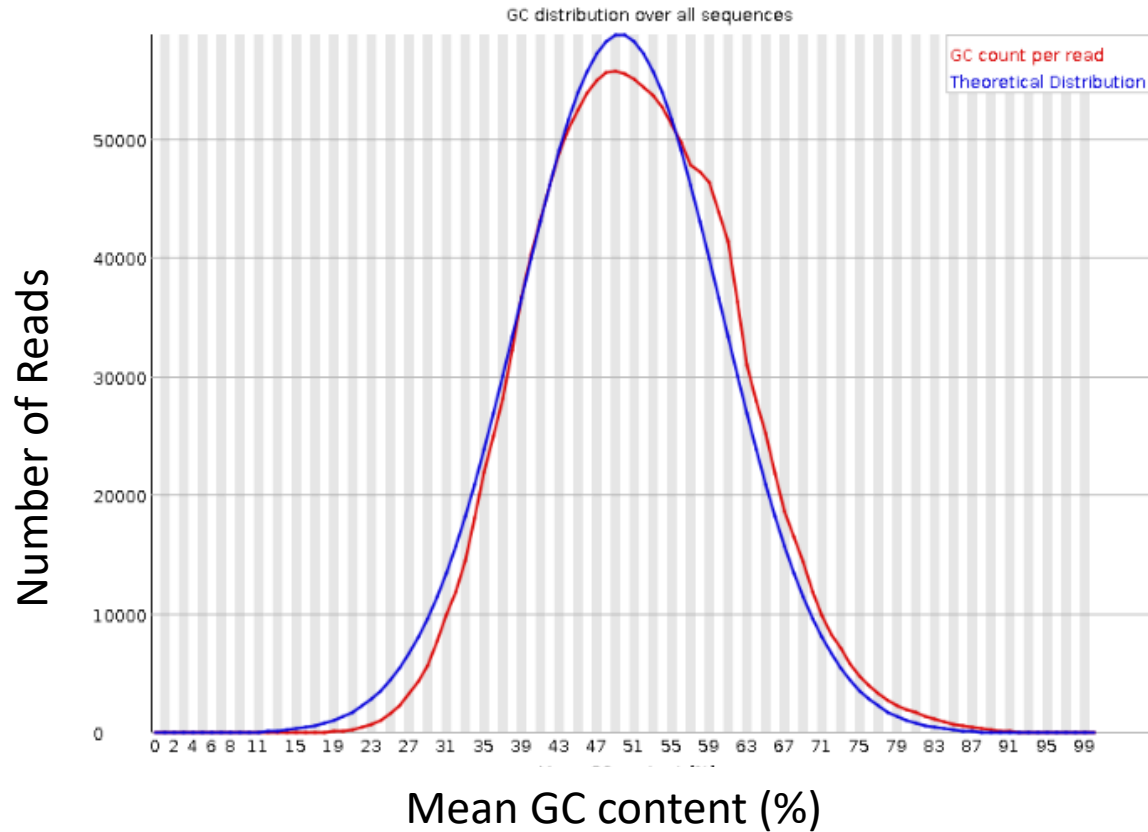
BAD

Read quality drops at the beginning and end



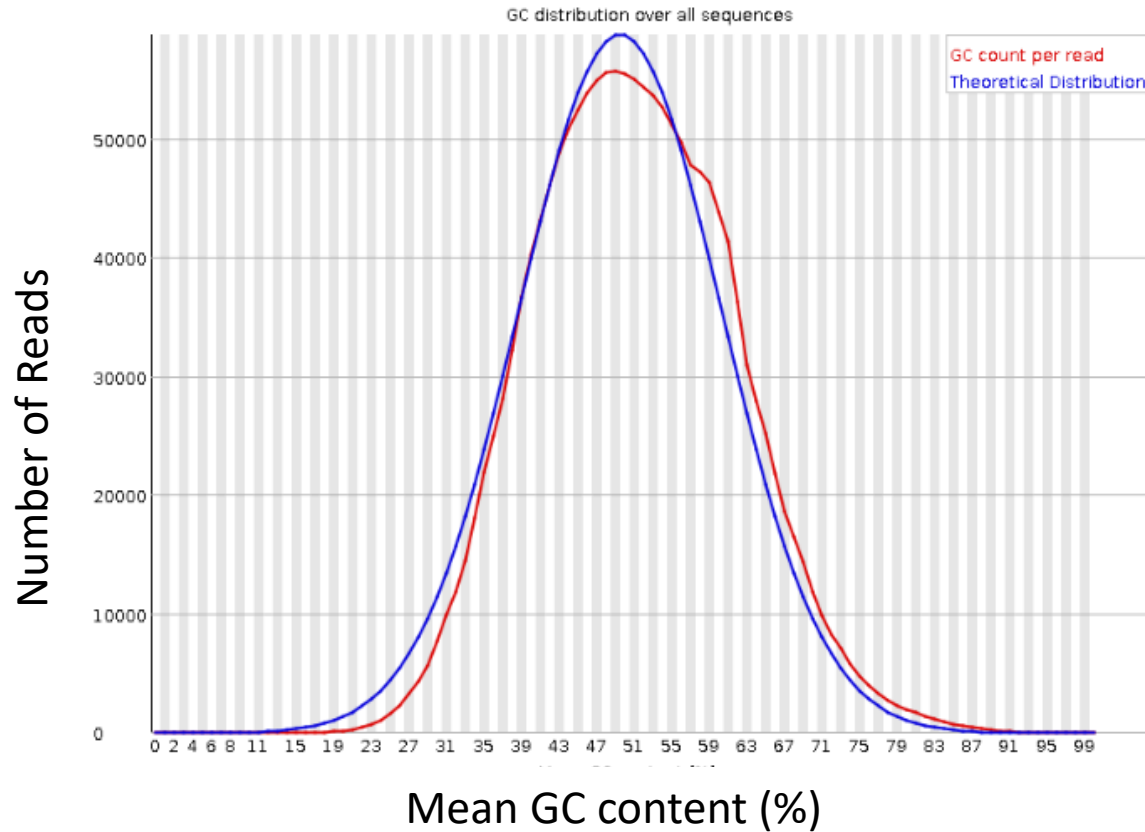
# FastQC: Per sequence GC content

## ✔ Per sequence GC content



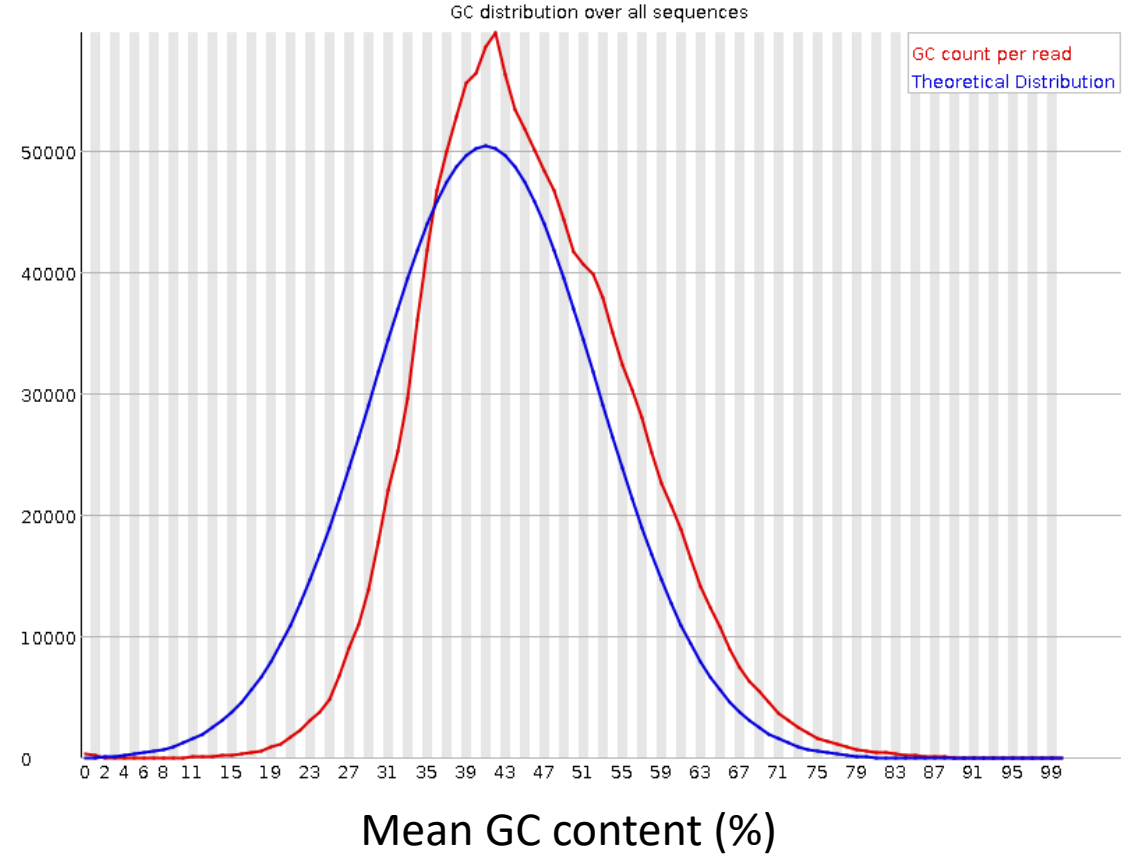
# FastQC: Per sequence GC content

## ✔ Per sequence GC content



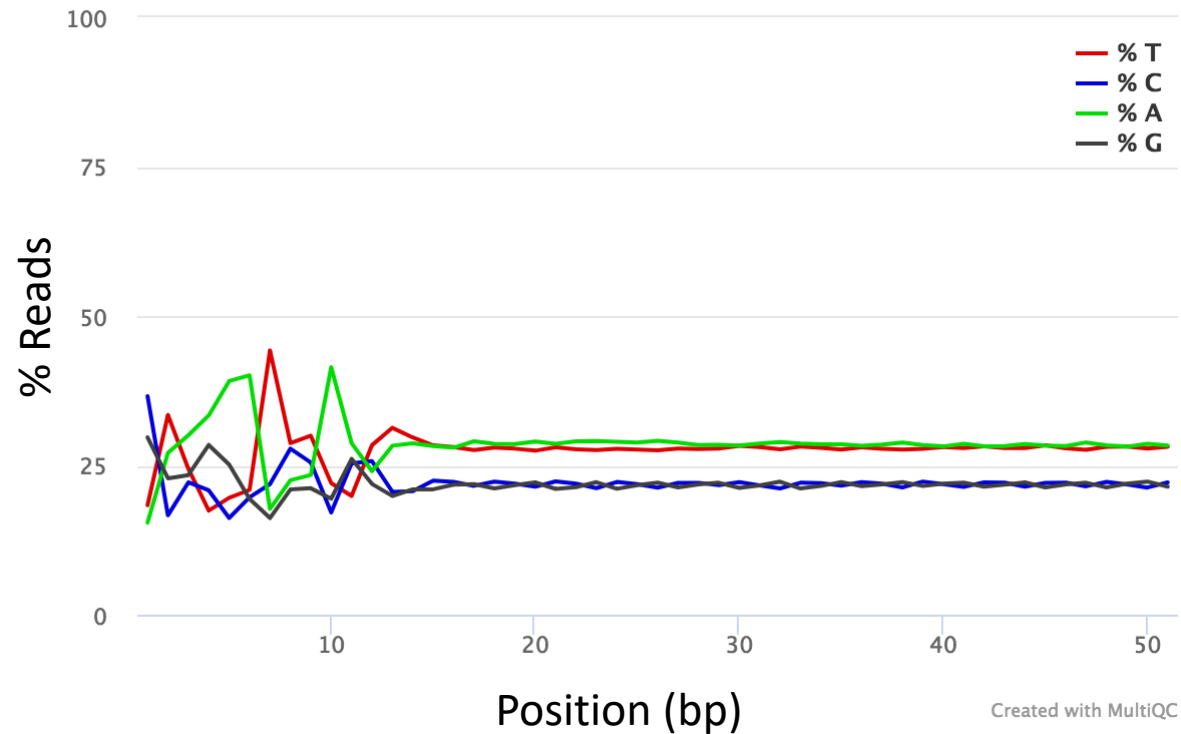
GOOD: follows normal distribution (sum of deviations is < 15% of reads)

## ✘ Per sequence GC content



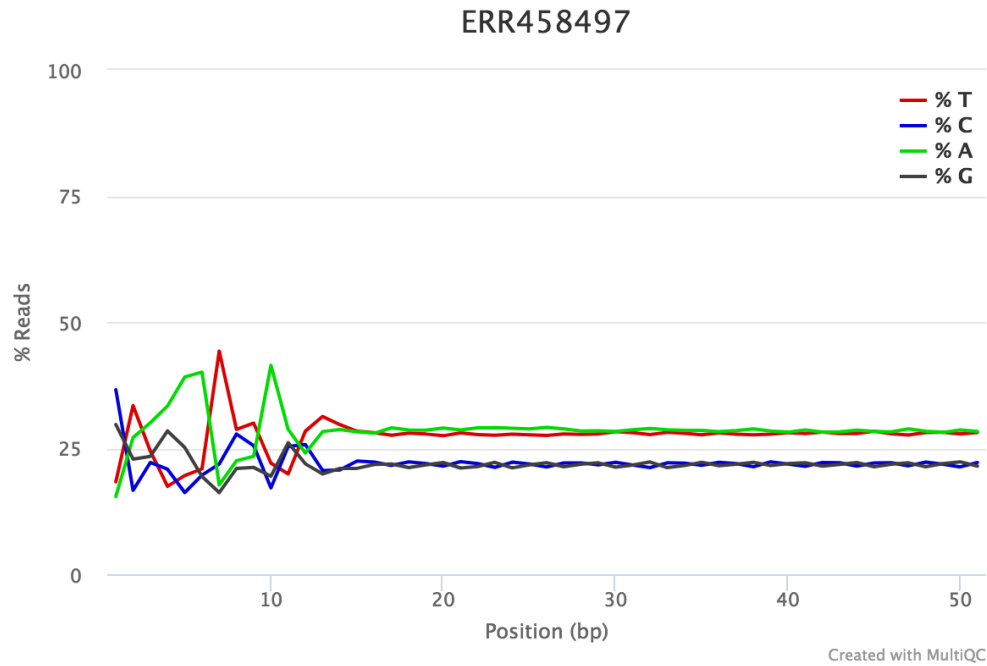
BAD: can indicate contamination with adapter dimers, or another species

# FastQC: Per Base Sequence Content

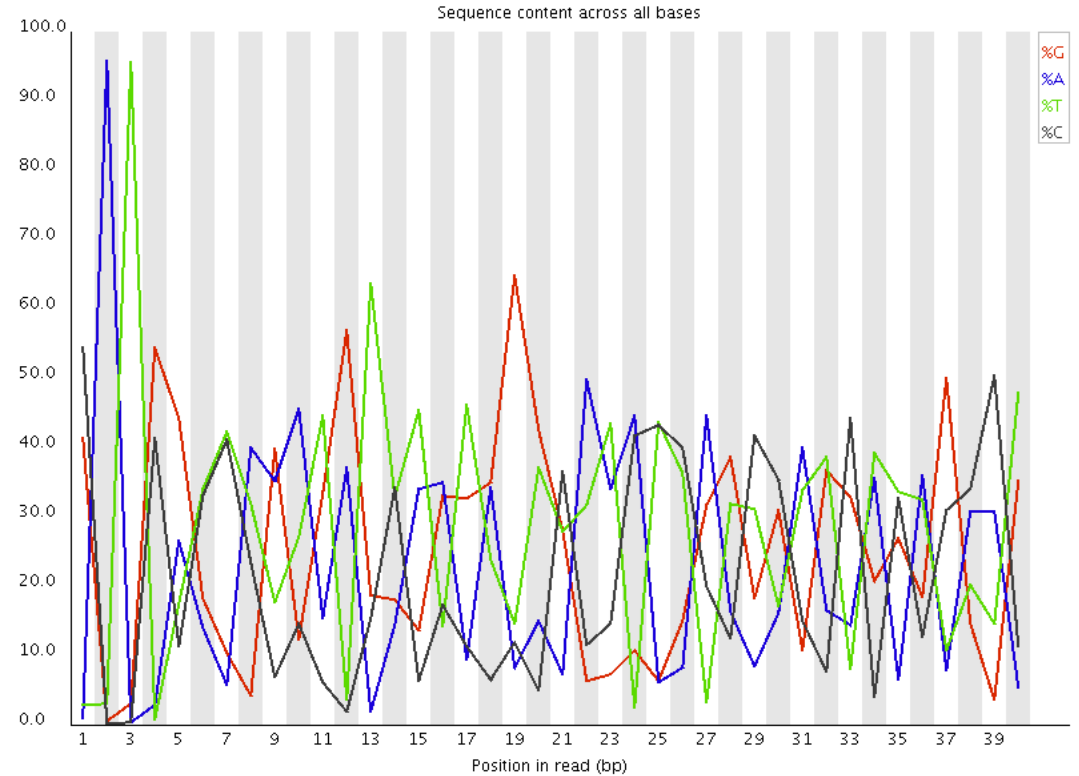


- Proportion of each position for which each DNA base has been called
- RNAseq data tends to show a positional sequence bias in the first ~12 bases
- The "random" priming step during library construction is not truly random and certain hexamers are more prevalent than others

# FastQC: Per Base Sequence Content



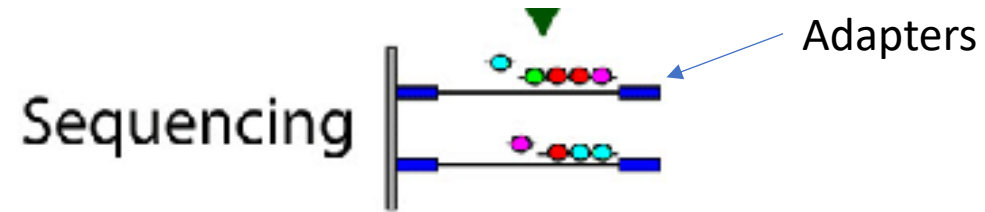
EXPECTED for RNAseq



**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

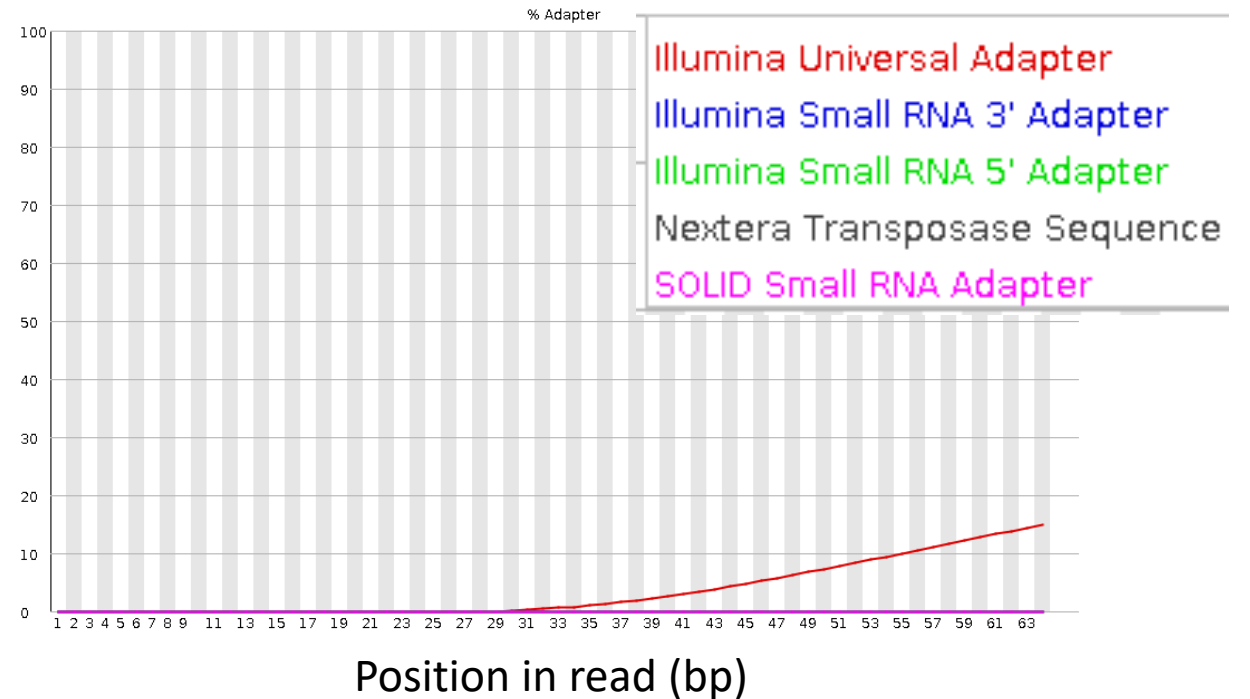
# FastQC: Adapter content



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

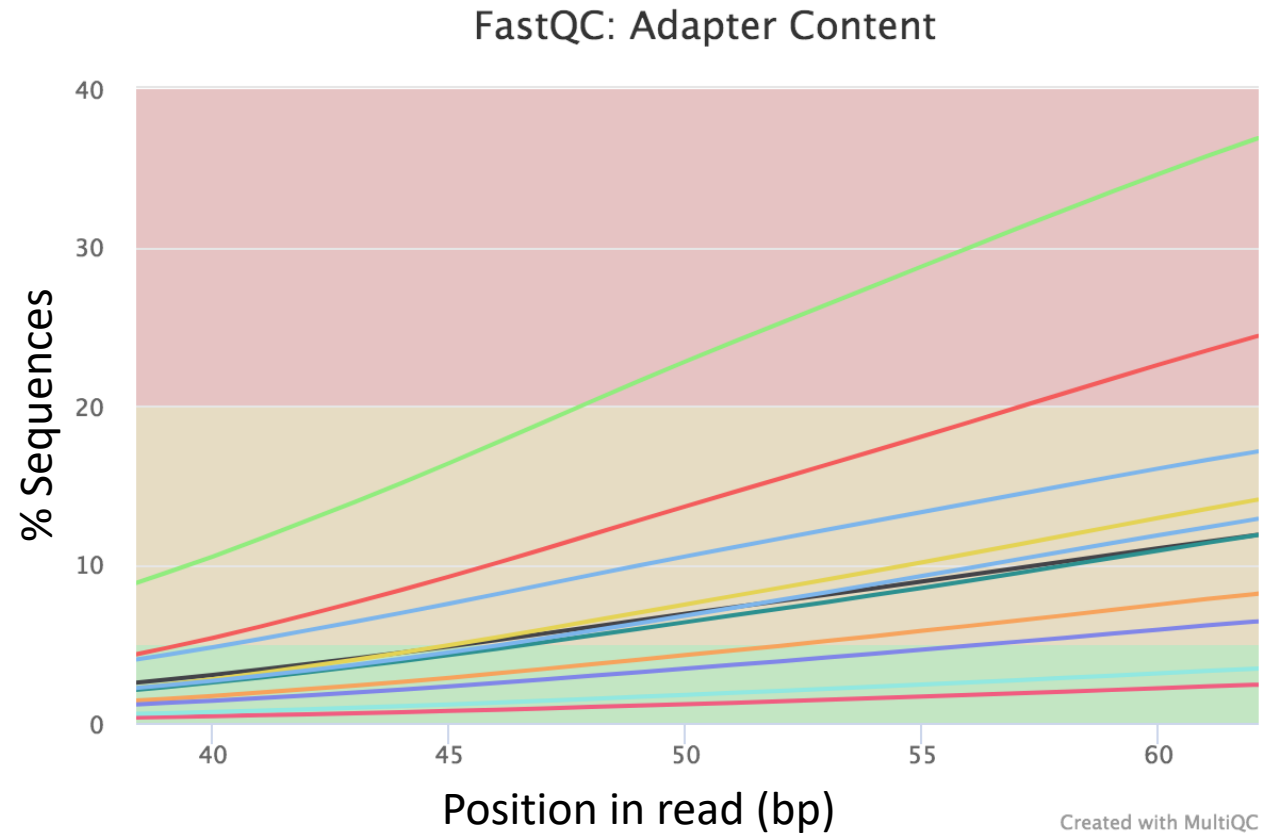
The plot shows that the adapter content rises over the course of the read

Solution – Adapter trimming!



# FastQC -> MultiQC

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



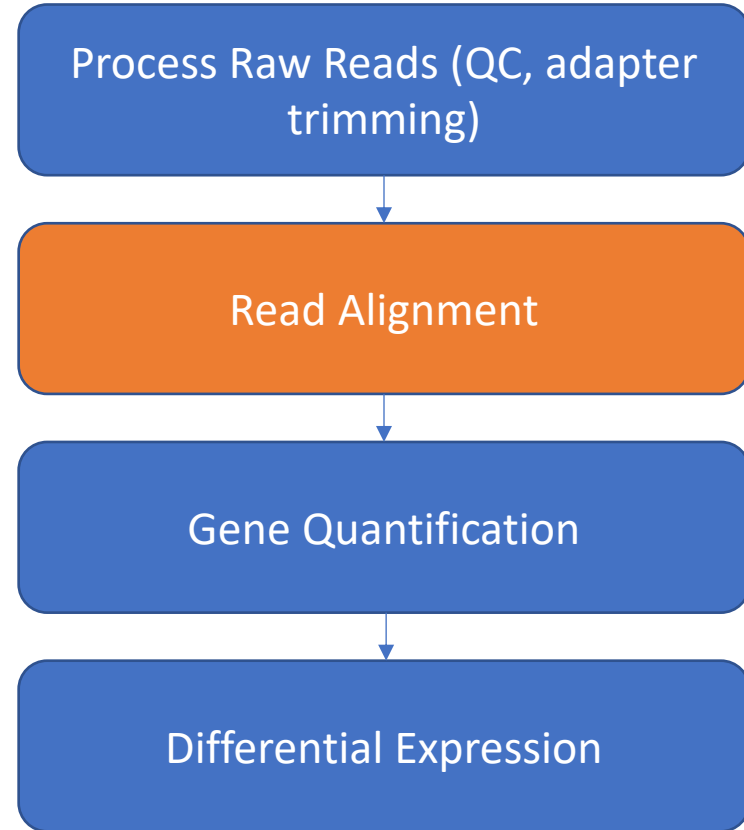


# Adapter trimming

Trim Galore! is a tool that:

- Scans and removes known Illumina or custom adapters
- Performs read trimming for low quality regions at the end of reads
- Removes reads that become too short in the trimming process

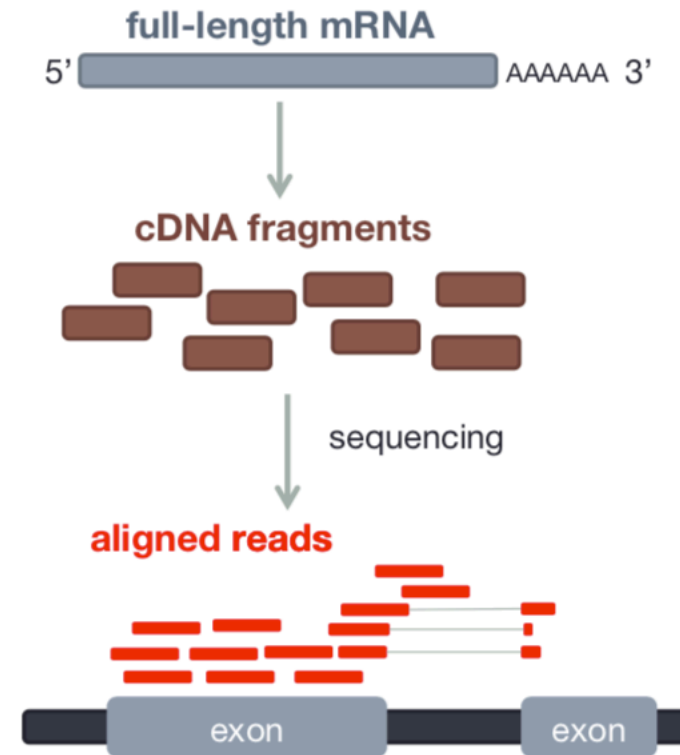
# Workflow



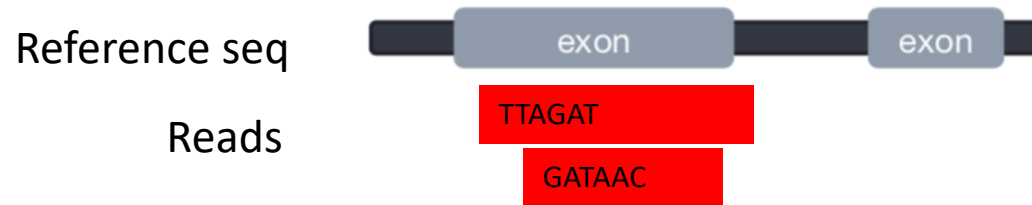
# Read Alignment

- 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



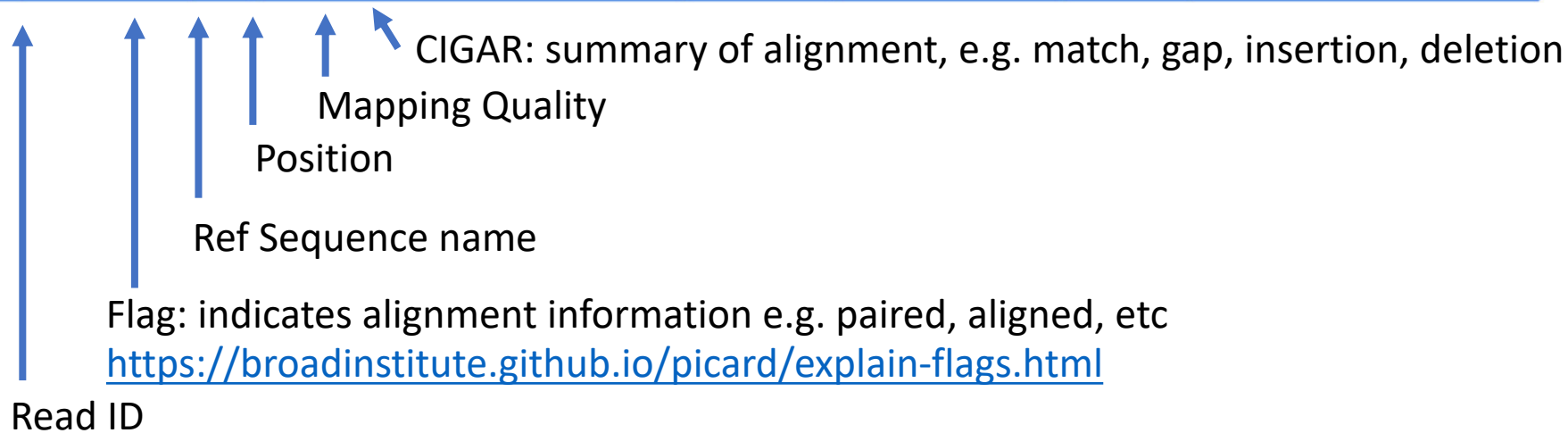
# Sequence Alignment Map (SAM)



```
@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



# 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

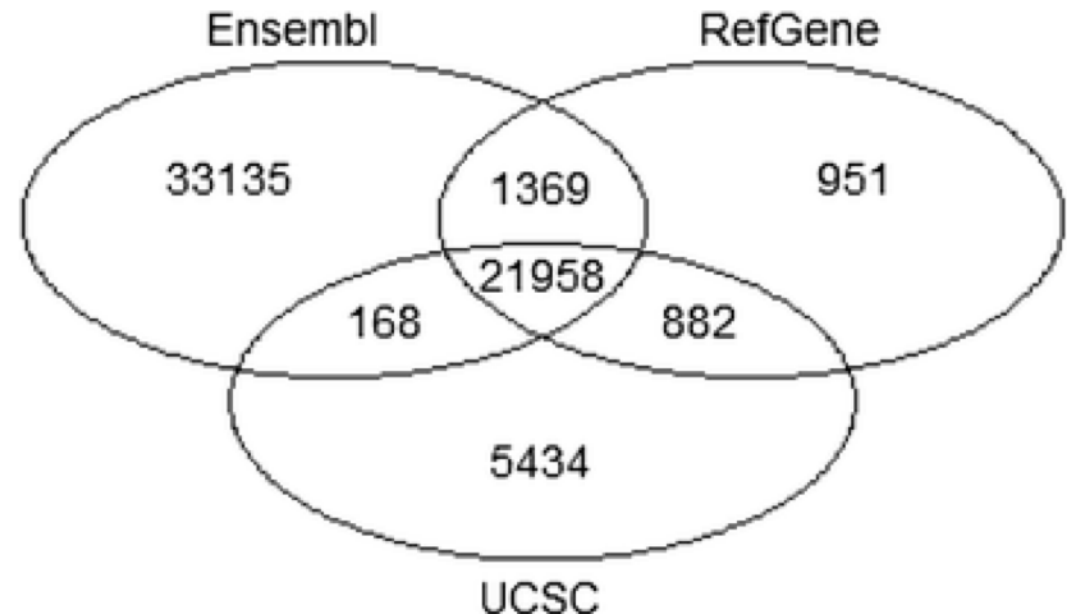
Sequence

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!



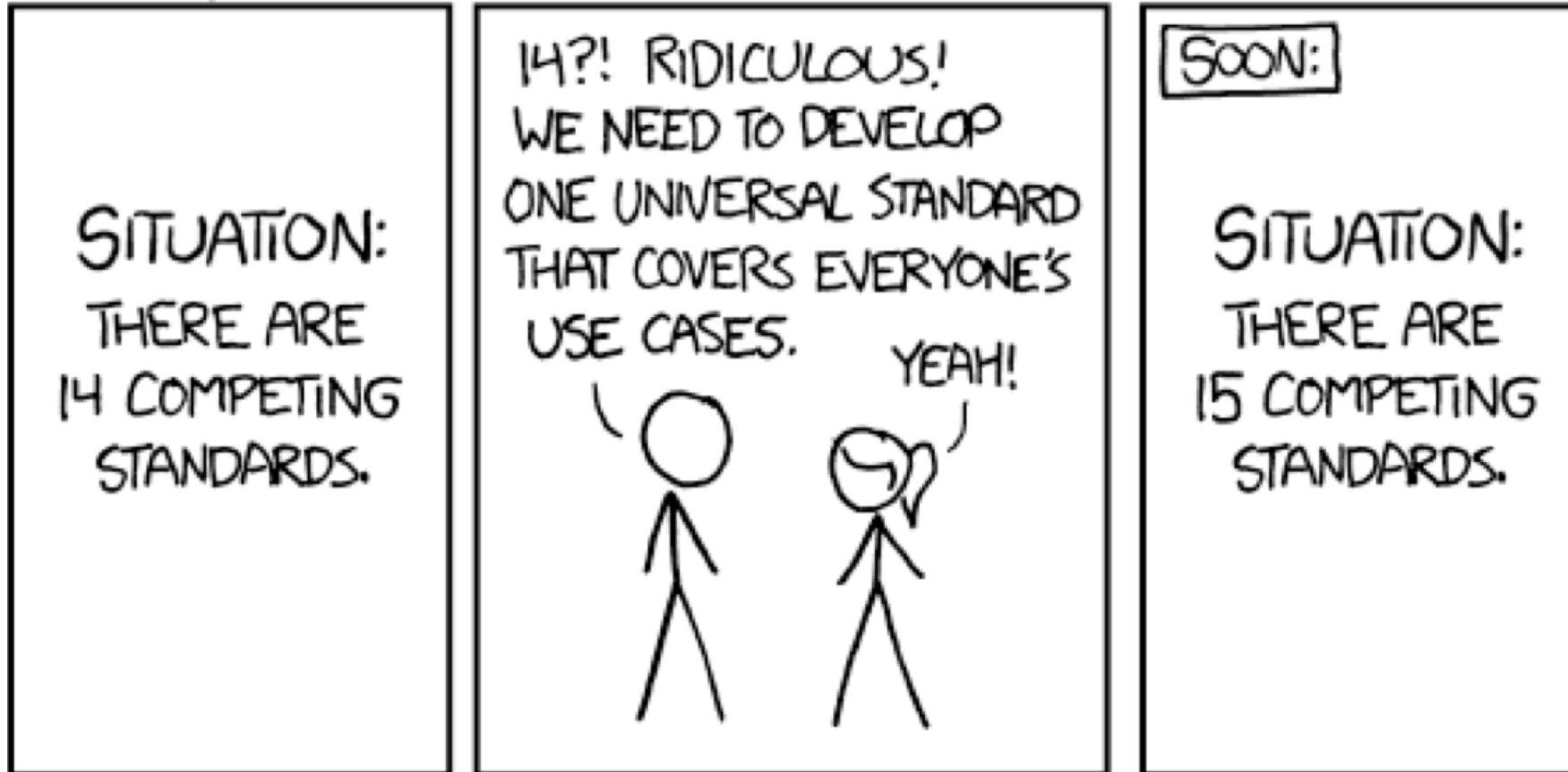
# Gene Annotation Format (GTF)

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	Frame			Attribute
					Strand	(Score)		
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";

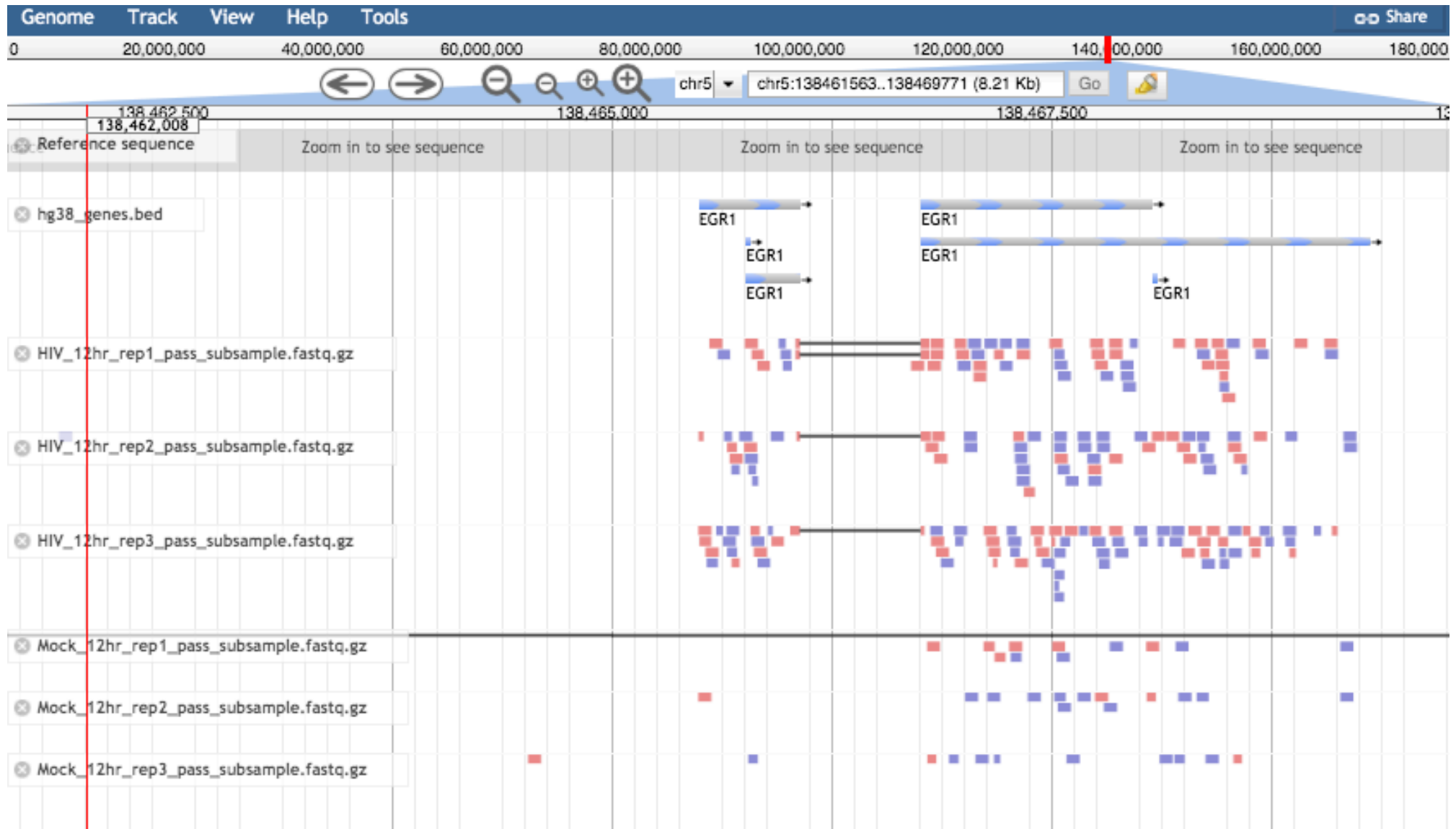
# A note on standards

HOW STANDARDS PROLIFERATE:  
(SEE: A/C CHARGERS, CHARACTER ENCODINGS, INSTANT MESSAGING, ETC)

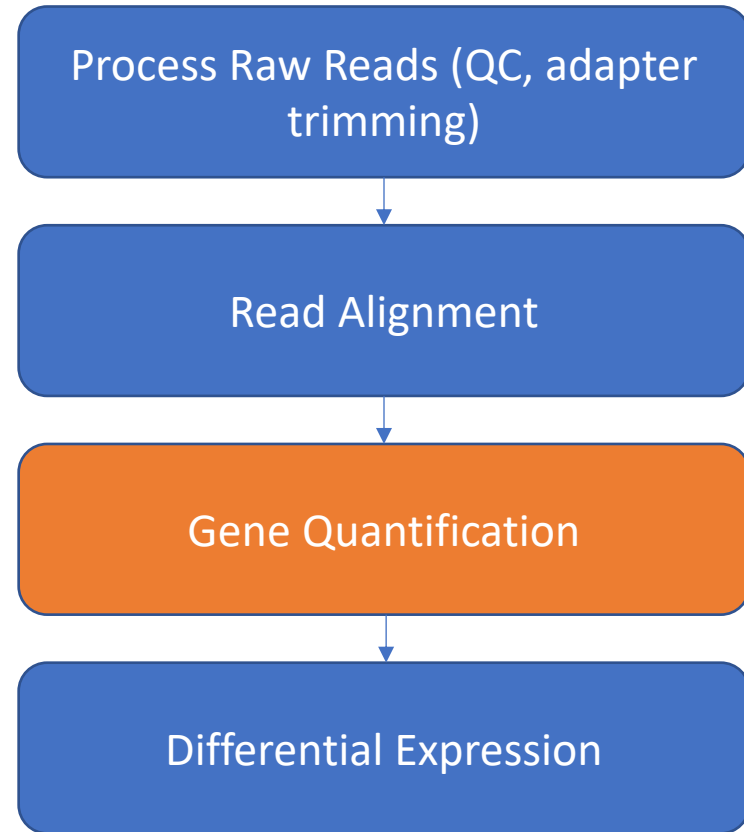




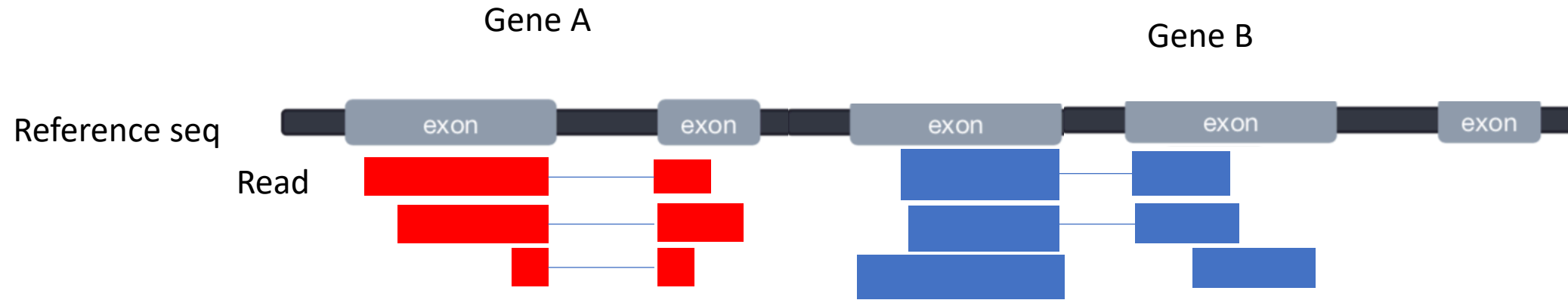
# Visualizing reads with JBrowse



# Workflow

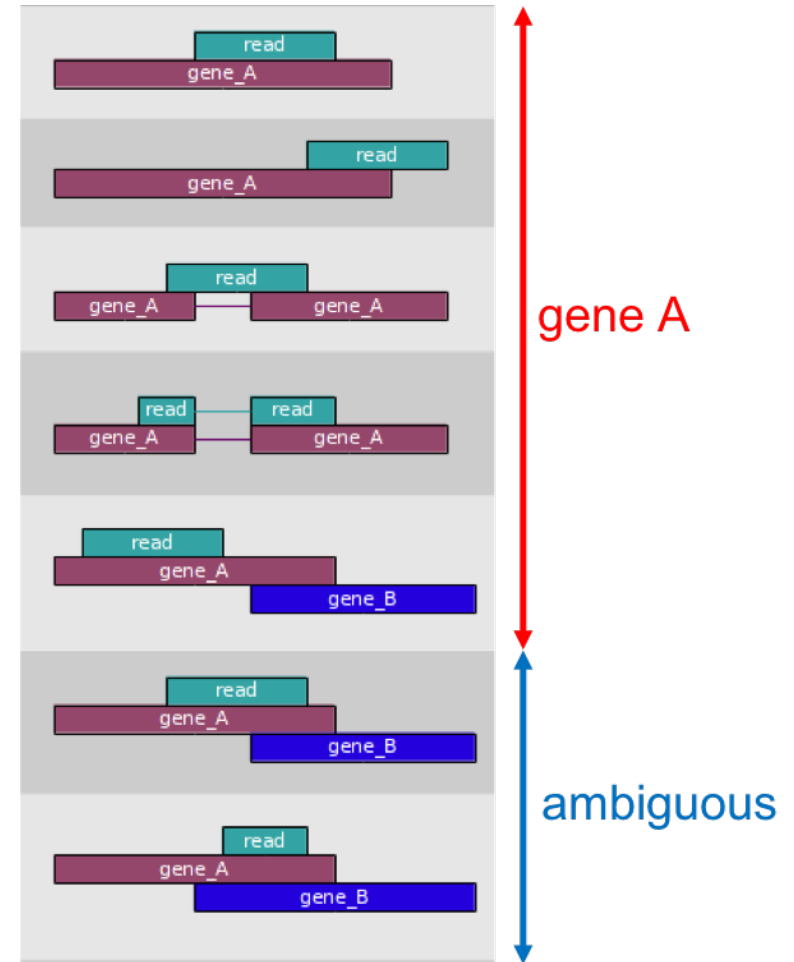


# Counting reads for each gene



# 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  $\geq 1$  bp are counted as belonging to that feature
- Ambiguous reads will be discarded

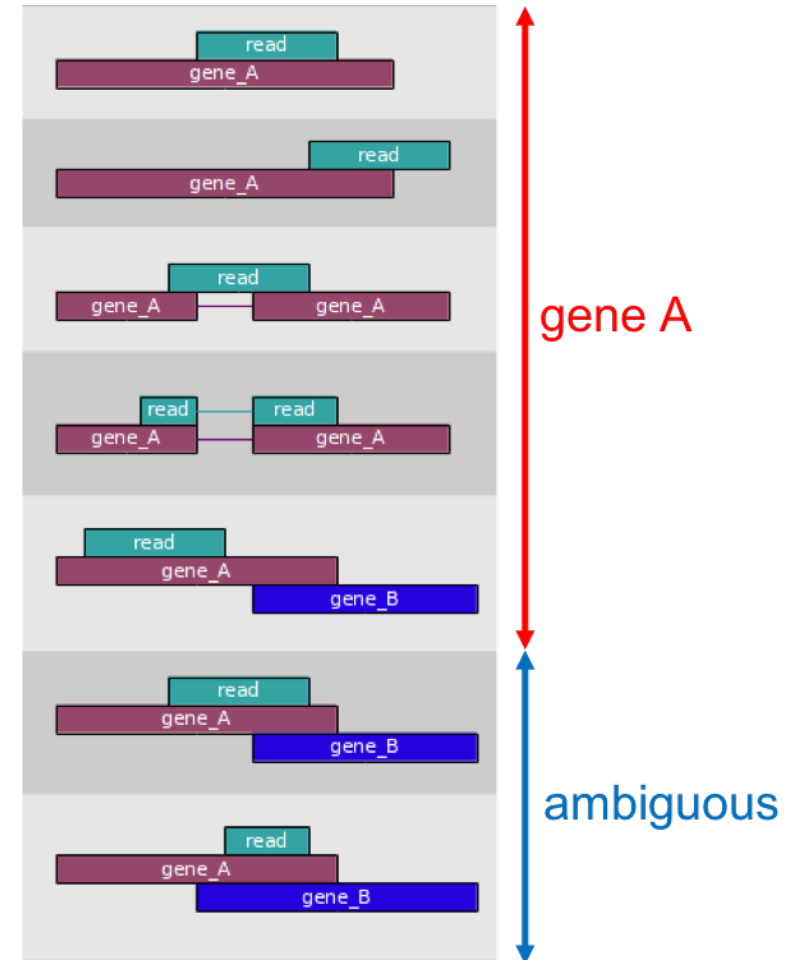


# Counting reads: featurecounts

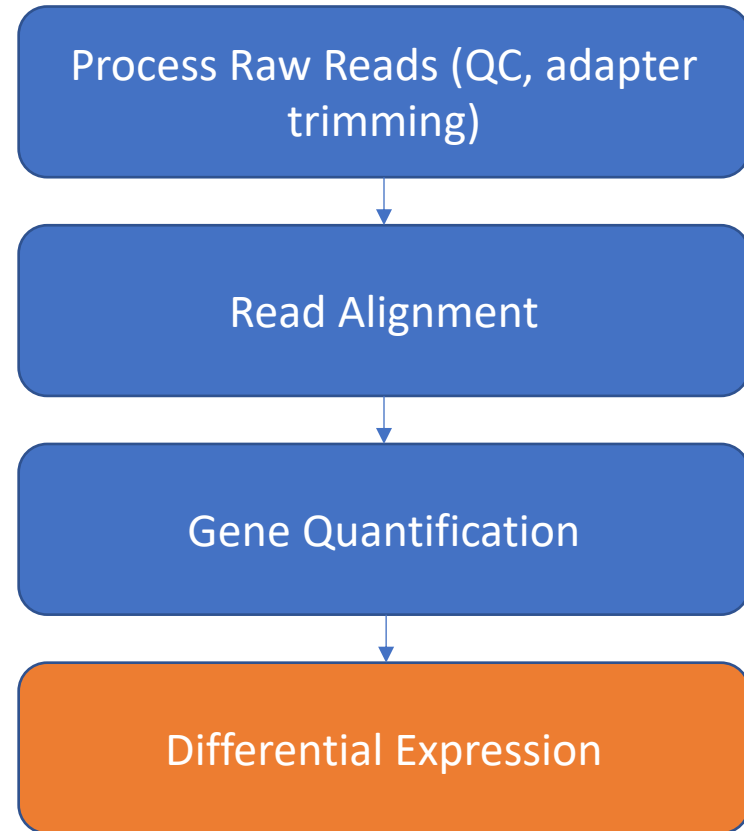
- The mapped coordinates of each read are compared with the features in the GTF file
- Reads that overlap with a gene by  $\geq 1$  bp are counted as belonging to that feature
- Ambiguous reads will be discarded

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

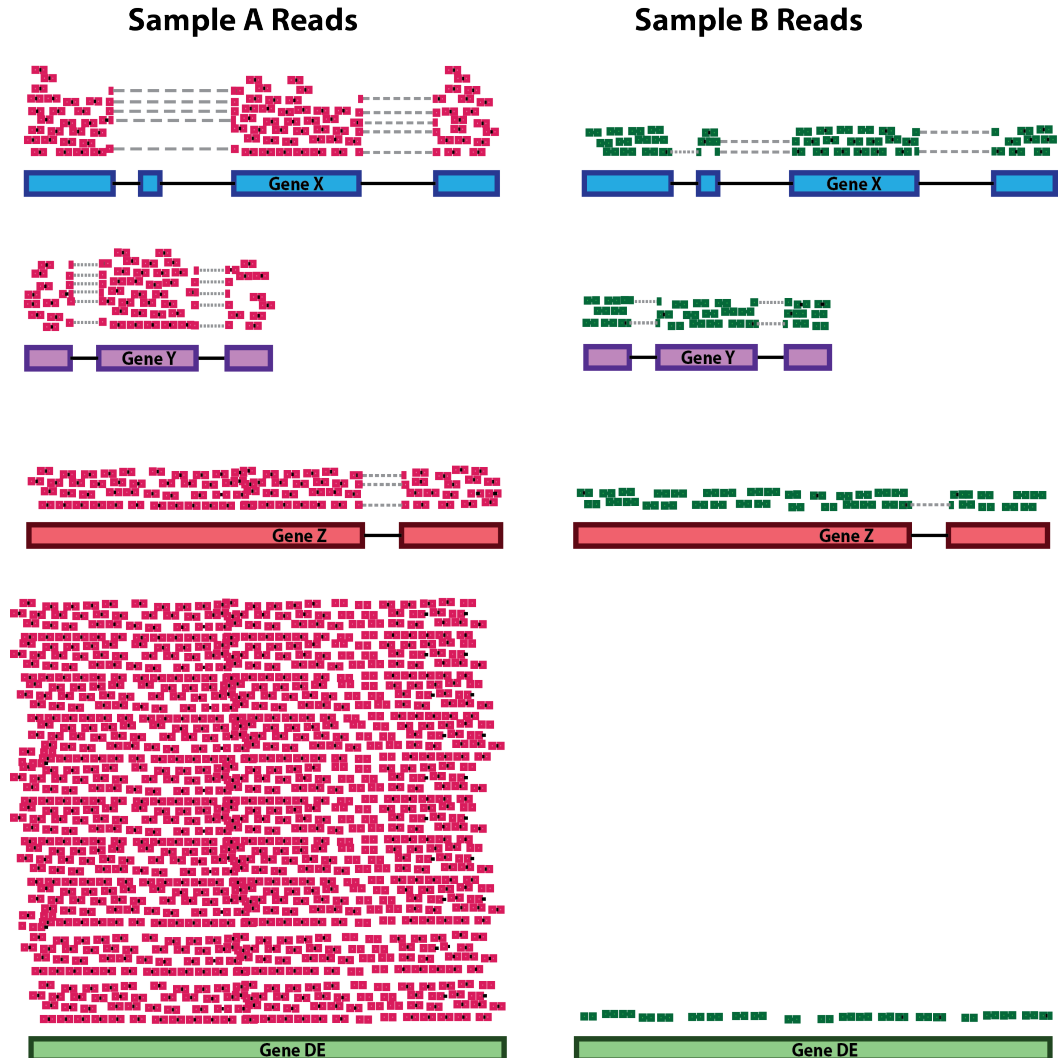


# Testing for Differential Expression

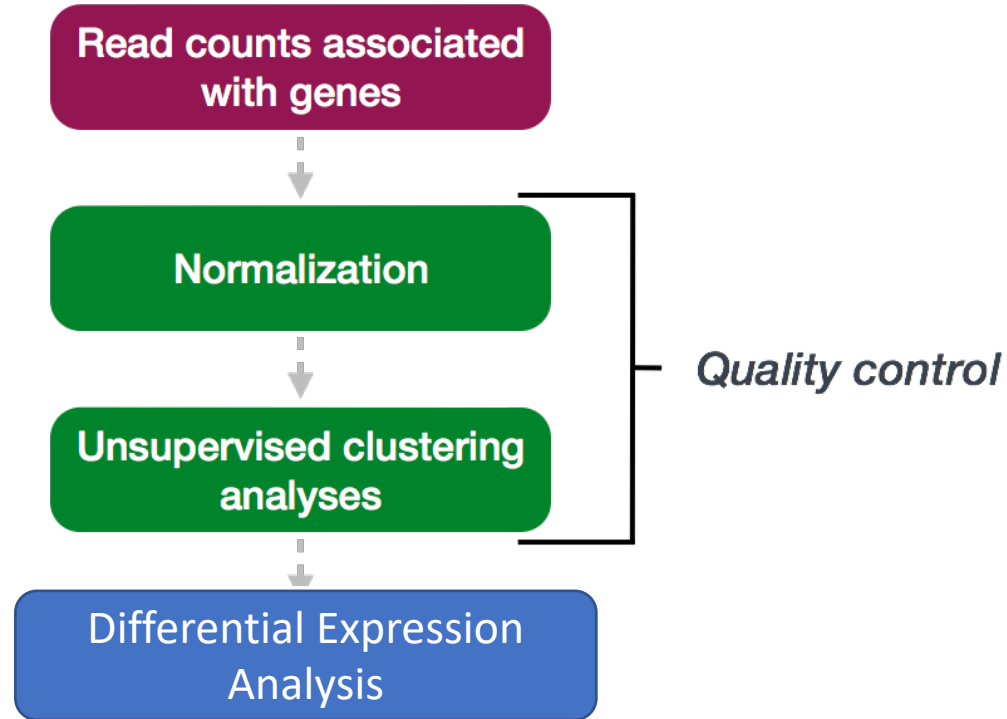
The goal of differential expression analysis (DE) is to find gene differences between conditions, developmental stages, treatments etc.

In particular DE has two goals:

- Estimate the *magnitude* of expression differences;
- Estimate the *significance* of expression differences.



# Differential Expression with DESeq2



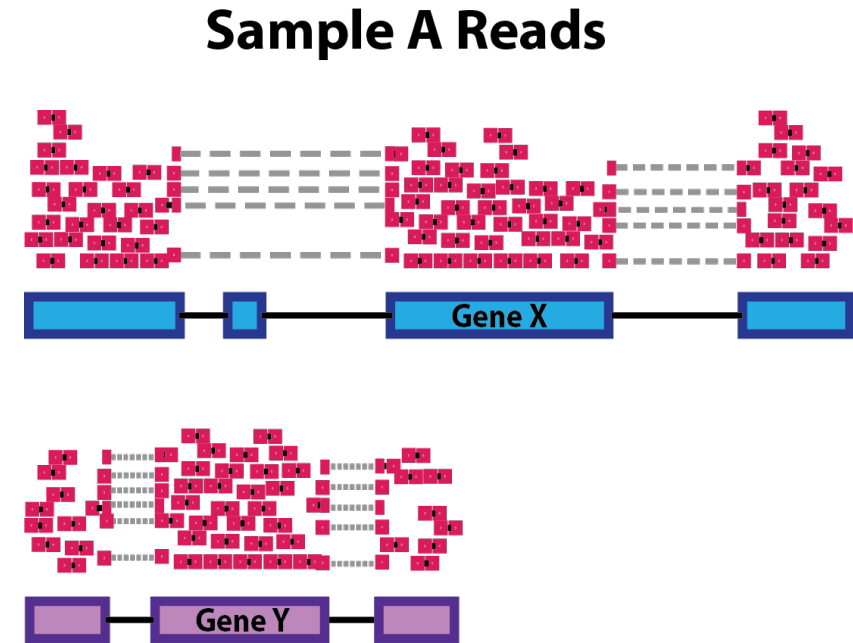
All steps are done with one click in Galaxy!



# Normalization

The number of sequenced reads mapped to a gene depends on

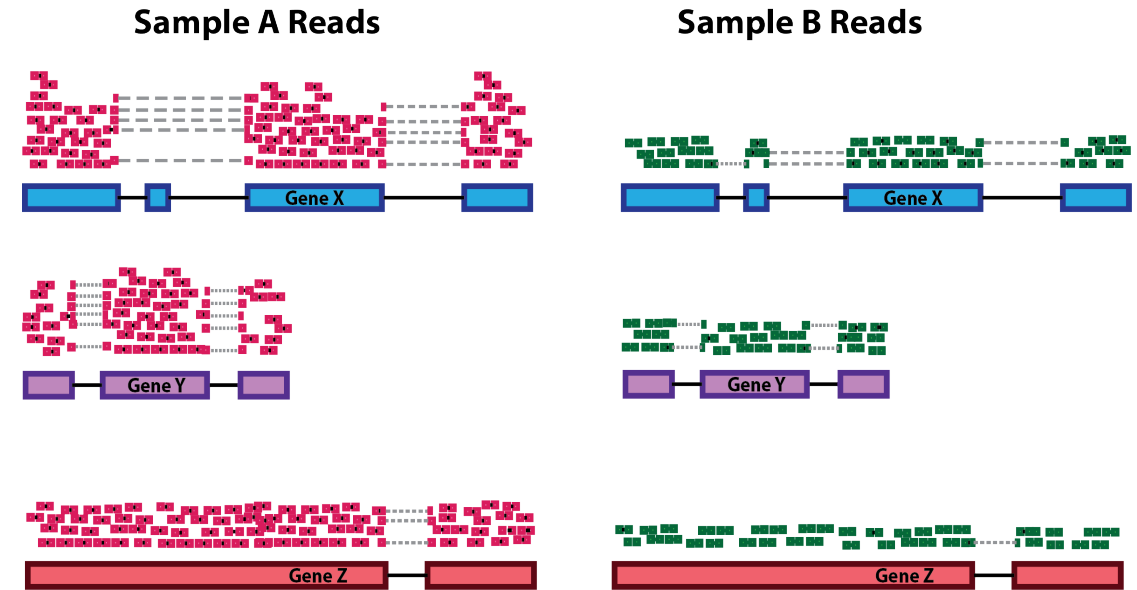
- **Gene Length**



# Normalization

The number of sequenced reads mapped to a gene depends on

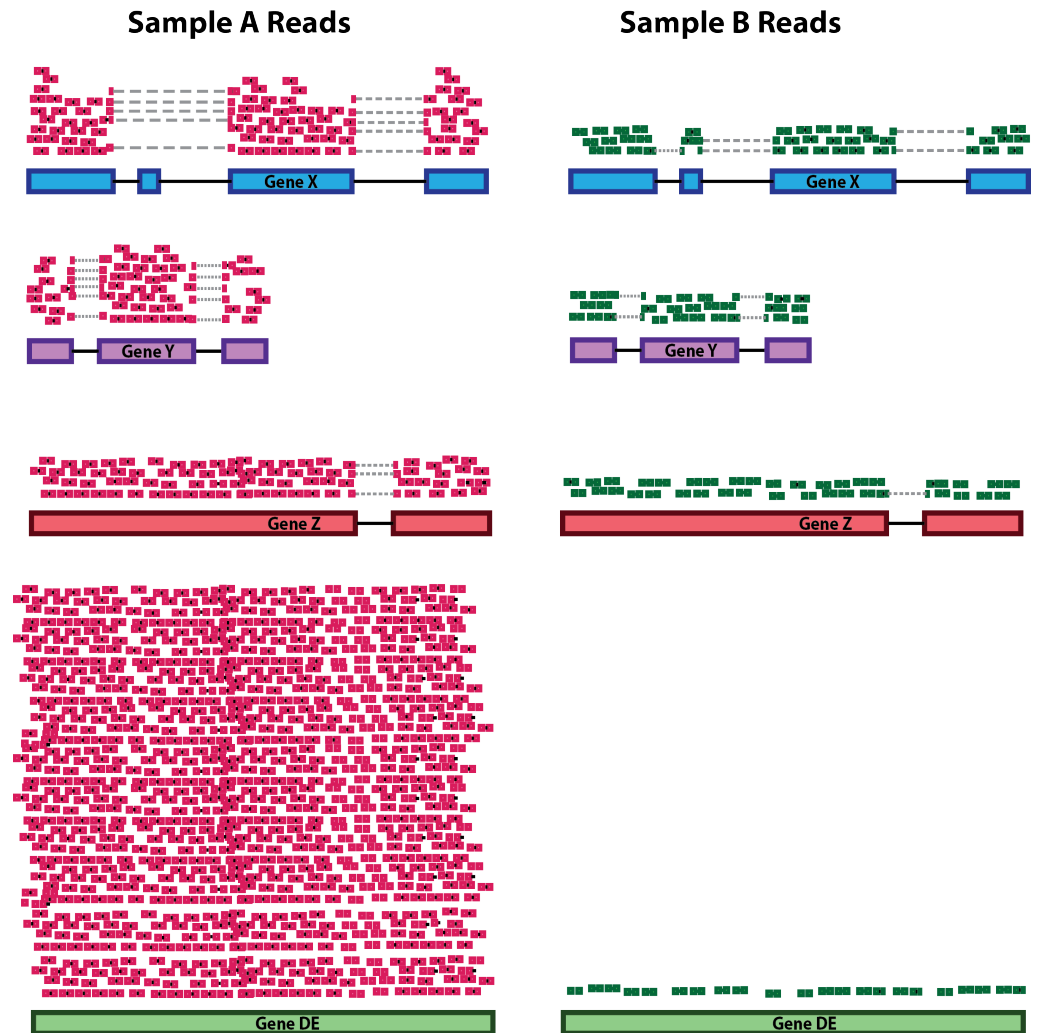
- Gene Length
- **Sequencing depth**



# Normalization

The number of sequenced reads mapped to a gene depends on

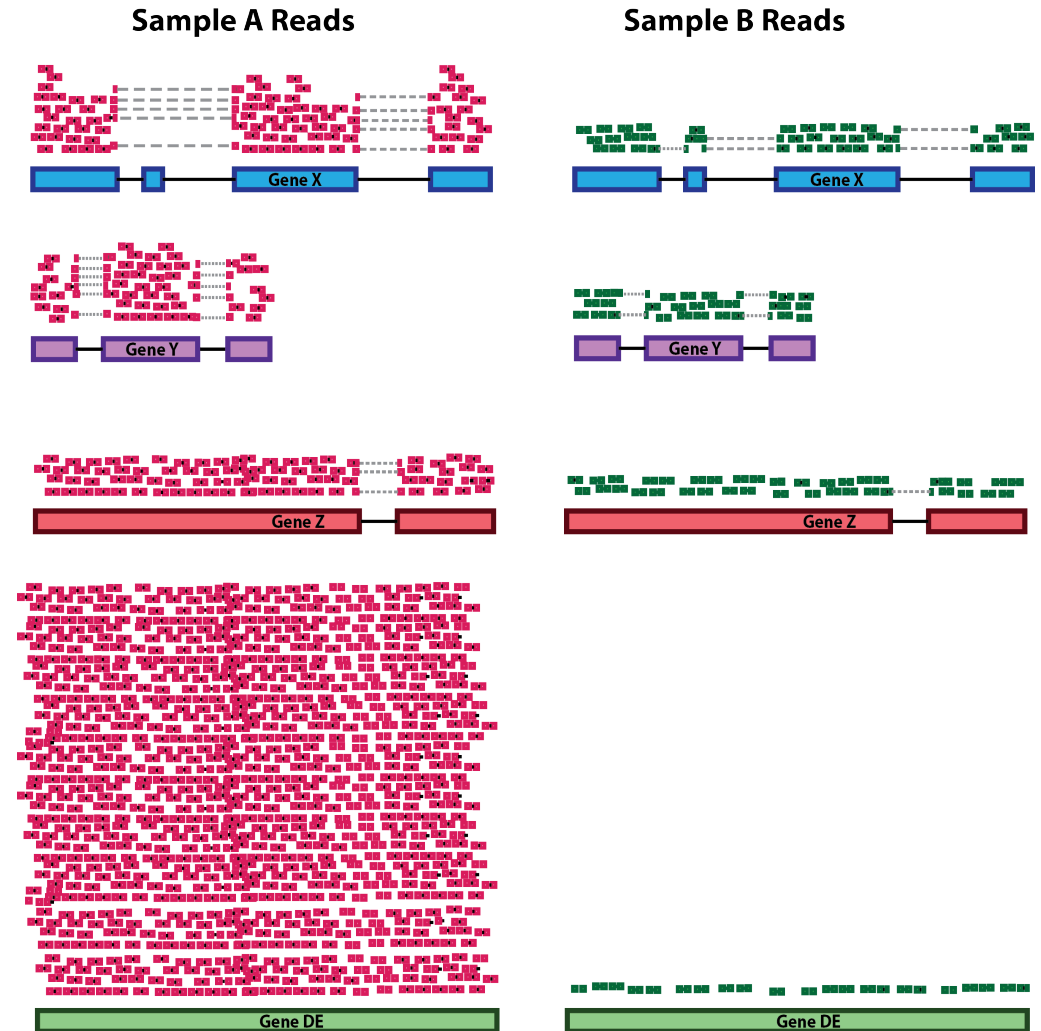
- Gene Length
- Sequencing depth
- **The expression level of other genes in the sample**



# Normalization

The number of sequenced reads mapped to a gene depends on

- ~~Gene Length~~
- ~~Sequencing depth~~
- ~~The expression level of other genes in the sample~~
- **It's own expression level**



Normalization eliminates the factors that are not of interest!

# Normalization: DESeq2 Median of Ratios

Accounts for both sequencing depth and composition

## Step 1: creates a pseudo-reference sample (row-wise geometric mean)

For each gene, a pseudo-reference sample is created that is equal to the geometric mean across all samples.

gene	sampleA	sampleB	pseudo-reference sample
1	1000	1000	$\sqrt{(1000 * 1000)} = 1000$
2	10	1	$\sqrt{(10 * 1)} = 3.16$
...	...	...	...

# Normalization: DESeq2 Median of Ratios

## Step 2: calculates ratio of each sample to the reference

Calculate the ratio of each sample to the pseudo-reference. Since most genes aren't differentially expressed, ratios should be similar.

gene	sampleA	sampleB	pseudo-reference sample	ratio of sampleA/ref	ratio of sampleB/ref
1	1000	1000	1000	$1000/1000 = \mathbf{1.00}$	$1000/1000 = \mathbf{1.00}$
2	10	1	3.16	$10/3.16 = \mathbf{3.16}$	$1/3.16 = \mathbf{0.32}$
...	...	...	...		

# Normalization: DESeq2 Median of Ratios

**Step 2: calculates ratio of each sample to the reference**

Calculate the ratio of each sample to the pseudo-reference.

gene	sampleA	sampleB	pseudo-reference sample	ratio of sampleA/ref	ratio of sampleB/ref
1	1000	1000	1000	$1000/1000 = 1.00$	$1000/1000 = 1.00$
2	10	1	3.16	$10/3.16 = 3.16$	$1/3.16 = 0.32$
...	...	...	...	...	...

Median = 2.08                      Median = 0.66

**Step 3: calculate the normalization factor for each sample (size factor)**

The median value of all ratios for a given sample is taken as the normalization factor (size factor) for that sample:

# Normalization: DESeq2 Median of Ratios

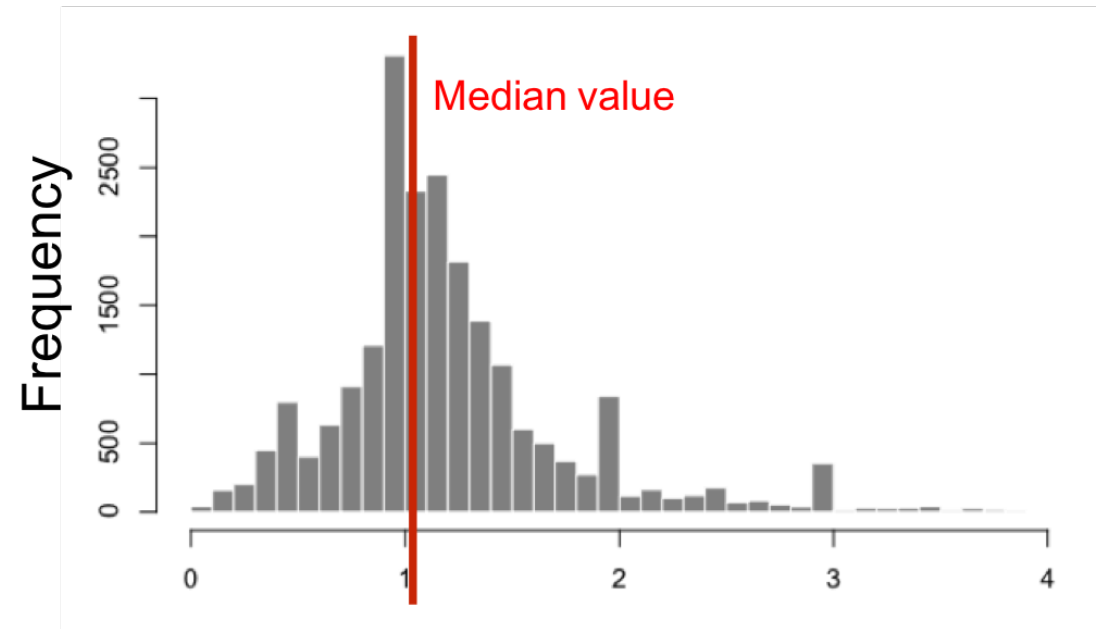
Visualization of normalization factor for a sample:

- Median should be  $\sim 1$  for each sample, otherwise data should be examined for the presence of large outliers
- This method is robust to imbalance in up-/down- regulation and large numbers of differentially expressed genes

Assumptions of this method:

Not all genes are differentially expressed

sample 1 / pseudo-reference sample





# Normalization: DESeq2 Median of Ratios

## Step 4: calculate the normalized count values using the normalization factor

This is performed by dividing each raw count value in a given sample by that sample's size factor to generate normalized count values.

SampleA normalization factor = 2.08

SampleB normalization factor = 0.66

### Raw Counts

gene	sampleA	sampleB
1	1000	1000
2	10	1

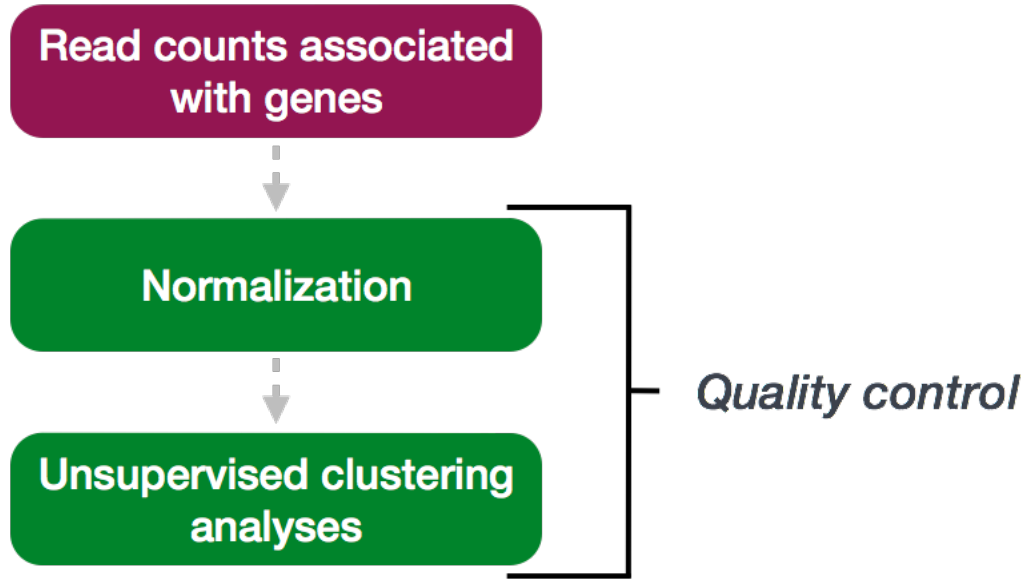
### Normalized Counts

gene	sampleA	sampleB
1	$1000/2.08 =$ <b>480.77</b>	$1000 / 0.66 =$ <b>1515.16</b>
2	$10/2.08 =$ <b>4.81</b>	$1 / 0.66 =$ <b>1.52</b>

# Normalization methods

Normalization method	Description	Accounted factors	For Differential Expression?
<b>CPM</b> (counts per million)	counts scaled by total number of reads in a sample	sequencing depth	NO
<b>TPM</b> (transcripts per kilobase million)	counts per length of transcript (kb) per million reads mapped	sequencing depth and gene length	NO
<b>RPKM/FPKM</b> (reads/fragments per kilobase of exon per million reads/fragments mapped)	similar to TPM	sequencing depth and gene length	NO
<b>DESeq2's median of ratios</b> [1]	counts divided by sample-specific size factors determined by median ratio of gene counts relative to geometric mean per gene	sequencing depth and RNA composition	YES

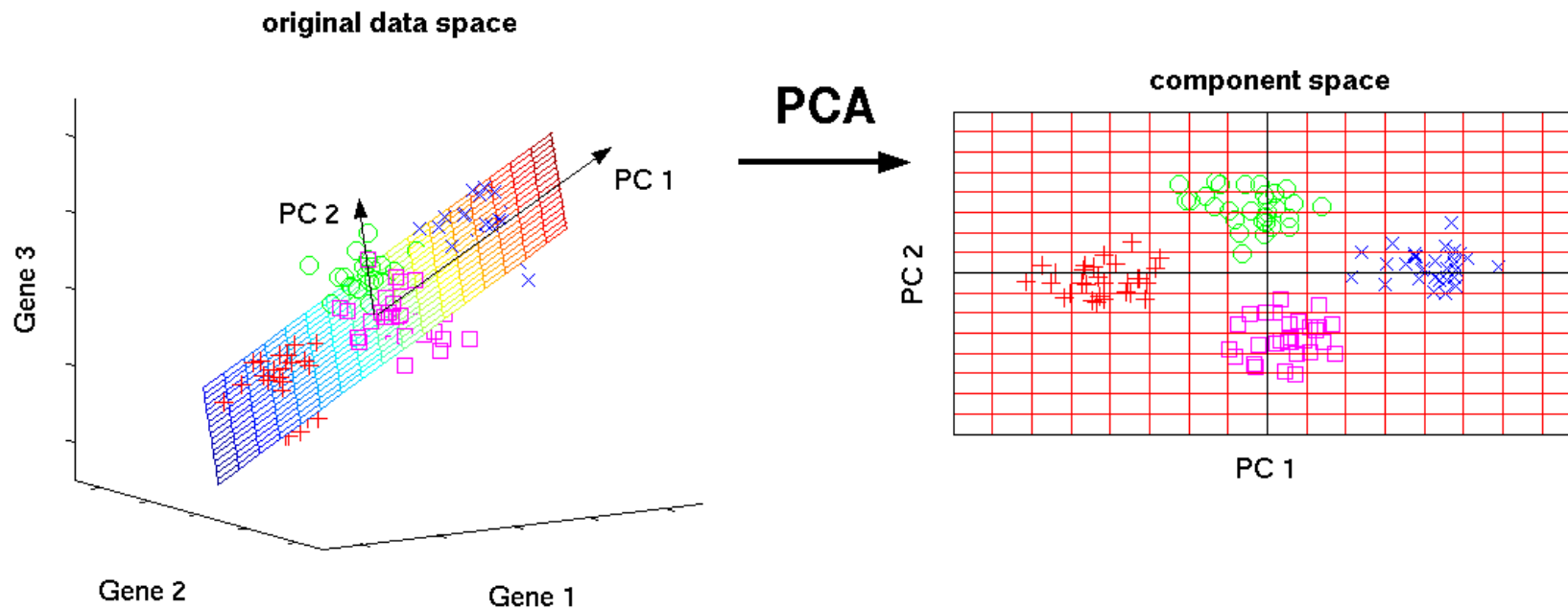
# Unsupervised Clustering



# Principle Component Analysis

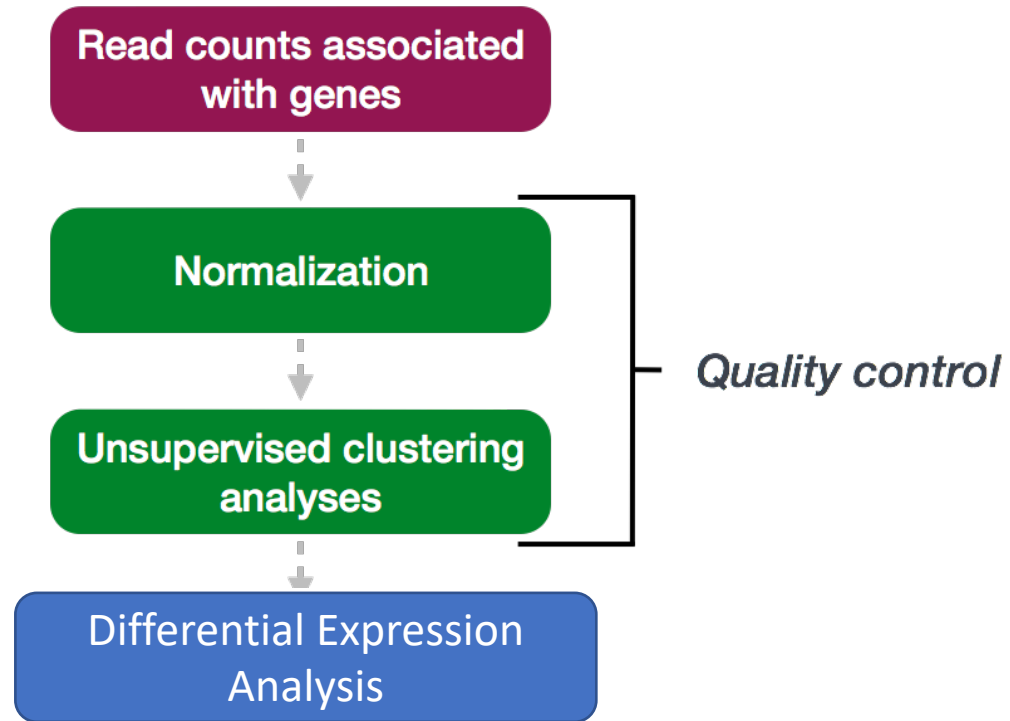
Here is an example with three genes measured in many samples:

Gene	Sample 1	Sample 2	Sample 3	Sample 4	
Gene 1	1000	1000	100	10	
Gene 2	10	1	5	6	...
Gene 3	10	1	10	20	

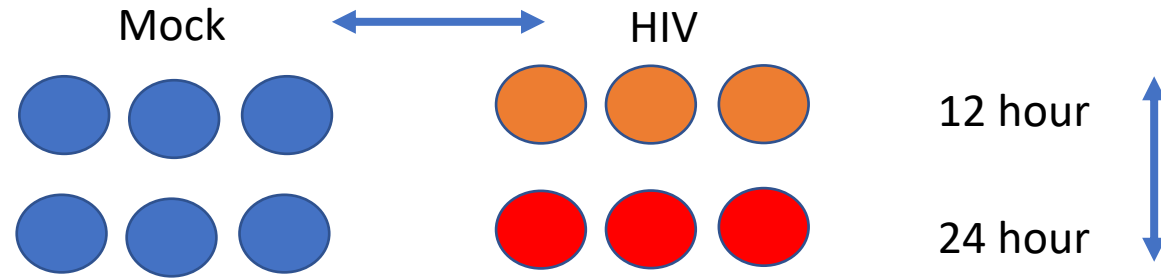


Do your samples cluster as expected?

# Differential Expression with DESeq2



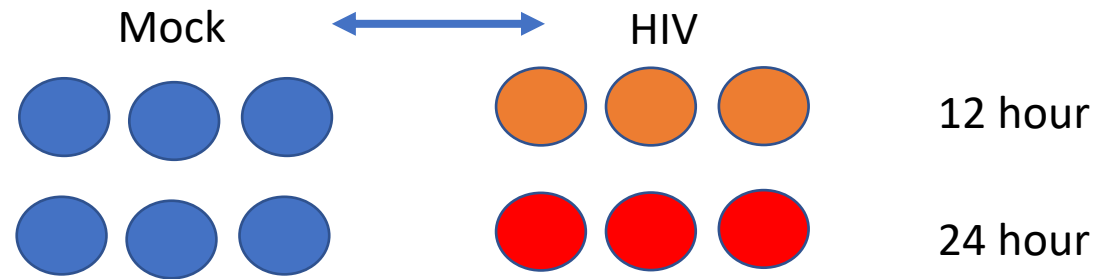
# Multi-factor experiment design



Factor 1:  
Infection status (Mock or HIV)

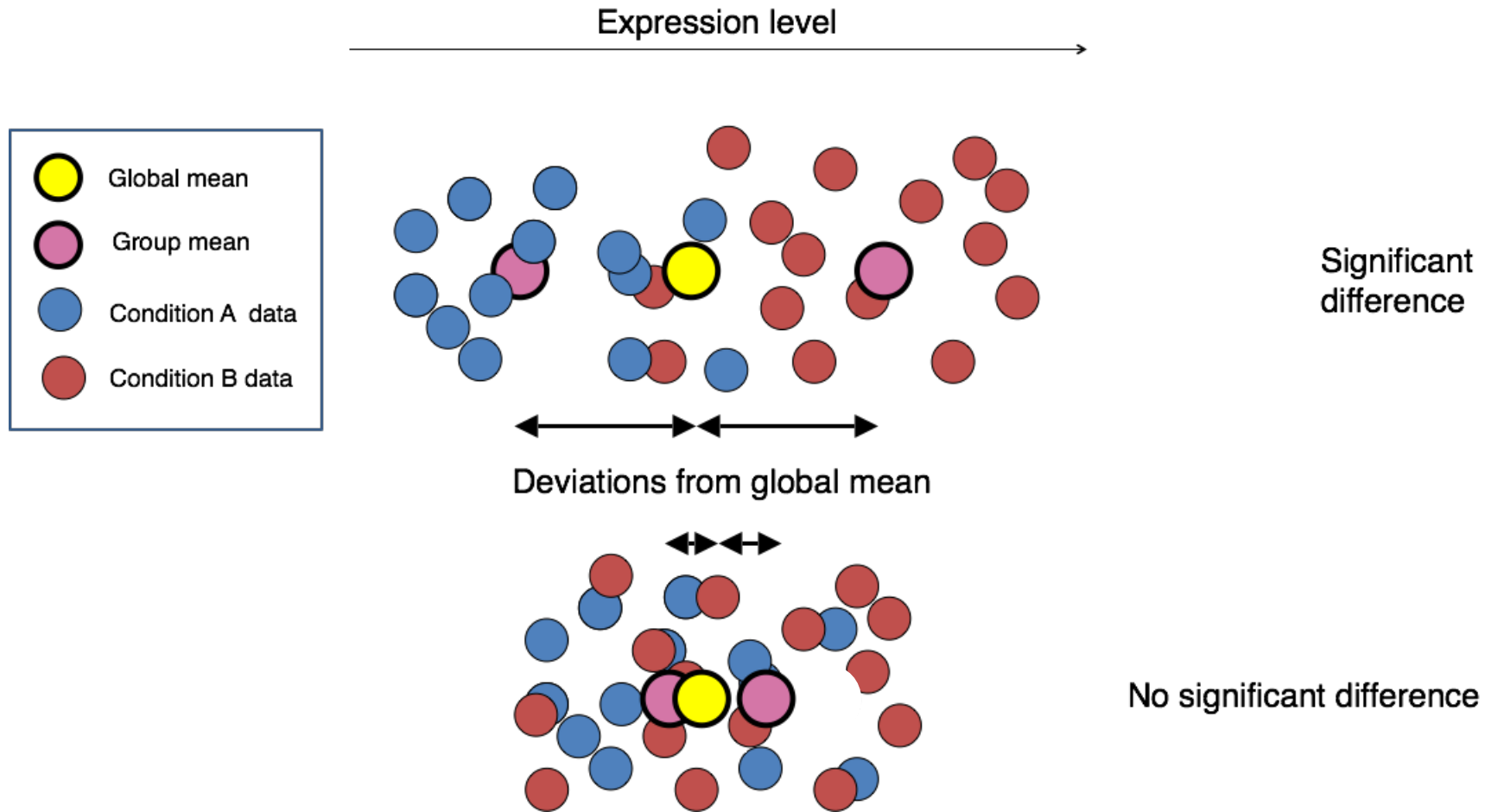
Factor 2:  
Time (12 or 24 hr)

# 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

# DESeq2 Test for Differential Expression



- DESeq2 models the gene counts for each gene as a [negative binomial distribution](#)
- One of the fitting parameters is the Log2foldChange for each gene

Image credit: Paul Pavlidis, UBC

[https://hbctraining.github.io/DGE\\_workshop/lessons/04\\_DGE\\_DESeq2\\_analysis.html](https://hbctraining.github.io/DGE_workshop/lessons/04_DGE_DESeq2_analysis.html)



# Wald Test

Statistical test (like T-test) used for hypothesis testing:

- Null hypothesis:  $\text{Log2foldChange}(\text{HIV counts} / \text{Mock counts}) == 0$
- Alternative hypothesis:  $\text{Log2FC}(\text{HIV counts} / \text{Mock counts}) \neq 0$

DESeq2 implements the Wald test by:

- Taking the  $\text{Log2foldChange}$  and dividing it by its standard error, resulting in a z-statistic
- The z-statistic is compared to a standard normal distribution, and a p-value is computed reporting the probability that a z-statistic at least as extreme as the observed value would be selected at random
- If the p-value is small we reject the null hypothesis and state that there is evidence against the null (i.e. the gene is differentially expressed).

# DESeq2 Results table

GeneID	Base mean	log2(FC)	StdErr	Wald-Stats	P-value	P-adj
EGR1	<b>1273.65</b>	-2.22	0.12	-18.65	1.25E-77	1.44E-73
MYC	<b>5226.12</b>	1.41	0.11	12.53	4.95E-36	2.87E-32
OPRK1	<b>78.35</b>	-1.83	0.17	-10.57	4.11E-26	1.59E-22
CCNI2	<b>7427.12</b>	0.93	0.10	9.43	4.27E-21	1.24E-17
STRA6	<b>785.78</b>	0.97	0.11	8.61	7.29E-18	1.69E-14

- Mean of normalized counts – averaged over all samples from two conditions

# DESeq2 Results table

GeneID	Base mean	log <sub>2</sub> (FC)	StdErr	Wald-Stats	P-value	P-adj
EGR1	1273.65	<b>-2.22</b>	0.12	-18.65	1.25E-77	1.44E-73
MYC	5226.12	<b>1.41</b>	0.11	12.53	4.95E-36	2.87E-32
OPRK1	78.35	<b>-1.83</b>	0.17	-10.57	4.11E-26	1.59E-22
CCNI2	7427.12	<b>0.93</b>	0.10	9.43	4.27E-21	1.24E-17
STRA6	785.78	<b>0.97</b>	0.11	8.61	7.29E-18	1.69E-14

- Mean of normalized counts – averaged over all samples from two conditions
- **Log of the fold change between two conditions** =  $\text{Log}_2(\text{HIV counts} / \text{Mock counts})$

# DESeq2 Results table

GeneID	Base mean	log2(FC)	StdErr	Wald-Stats	P-value	P-adj
EGR1	1273.65	-2.22	<b>0.12</b>	-18.65	1.25E-77	1.44E-73
MYC	5226.12	1.41	<b>0.11</b>	12.53	4.95E-36	2.87E-32
OPRK1	78.35	-1.83	<b>0.17</b>	-10.57	4.11E-26	1.59E-22
CCNI2	7427.12	0.93	<b>0.10</b>	9.43	4.27E-21	1.24E-17
STRA6	785.78	0.97	<b>0.11</b>	8.61	7.29E-18	1.69E-14

- Mean of normalized counts – averaged over all samples from two conditions
- Log2 of the fold change between two conditions
- **StdErr, standard error:**  $\text{Log}_2(\text{HIV counts}/\text{Mock counts}) = [-2.22 - 0.12, -2.22 + 0.12]$

# DESeq2 Results table

GeneID	Base mean	log2(FC)	StdErr	Wald-Stats	P-value	P-adj
EGR1	1273.65	-2.22	0.12	<b>-18.65</b>	1.25E-77	1.44E-73
MYC	5226.12	1.41	0.11	<b>12.53</b>	4.95E-36	2.87E-32
OPRK1	78.35	-1.83	0.17	<b>-10.57</b>	4.11E-26	1.59E-22
CCNI2	7427.12	0.93	0.10	<b>9.43</b>	4.27E-21	1.24E-17
STRA6	785.78	0.97	0.11	<b>8.61</b>	7.29E-18	1.69E-14

- Mean of normalized counts – averaged over all samples from two conditions
- Log of the fold change between two conditions
- **StdErr, standard error**
- **Wald Stats – statistical test for hypothesis testing**

# DESeq2 Results table

GeneID	Base mean	log2(FC)	StdErr	Wald-Stats	P-value	P-adj
EGR1	1273.65	-2.22	0.12	-18.65	<b>1.25E-77</b>	1.44E-73
MYC	5226.12	1.41	0.11	12.53	<b>4.95E-36</b>	2.87E-32
OPRK1	78.35	-1.83	0.17	-10.57	<b>4.11E-26</b>	1.59E-22
CCNI2	7427.12	0.93	0.10	9.43	<b>4.27E-21</b>	1.24E-17
STRA6	785.78	0.97	0.11	8.61	<b>7.29E-18</b>	1.69E-14

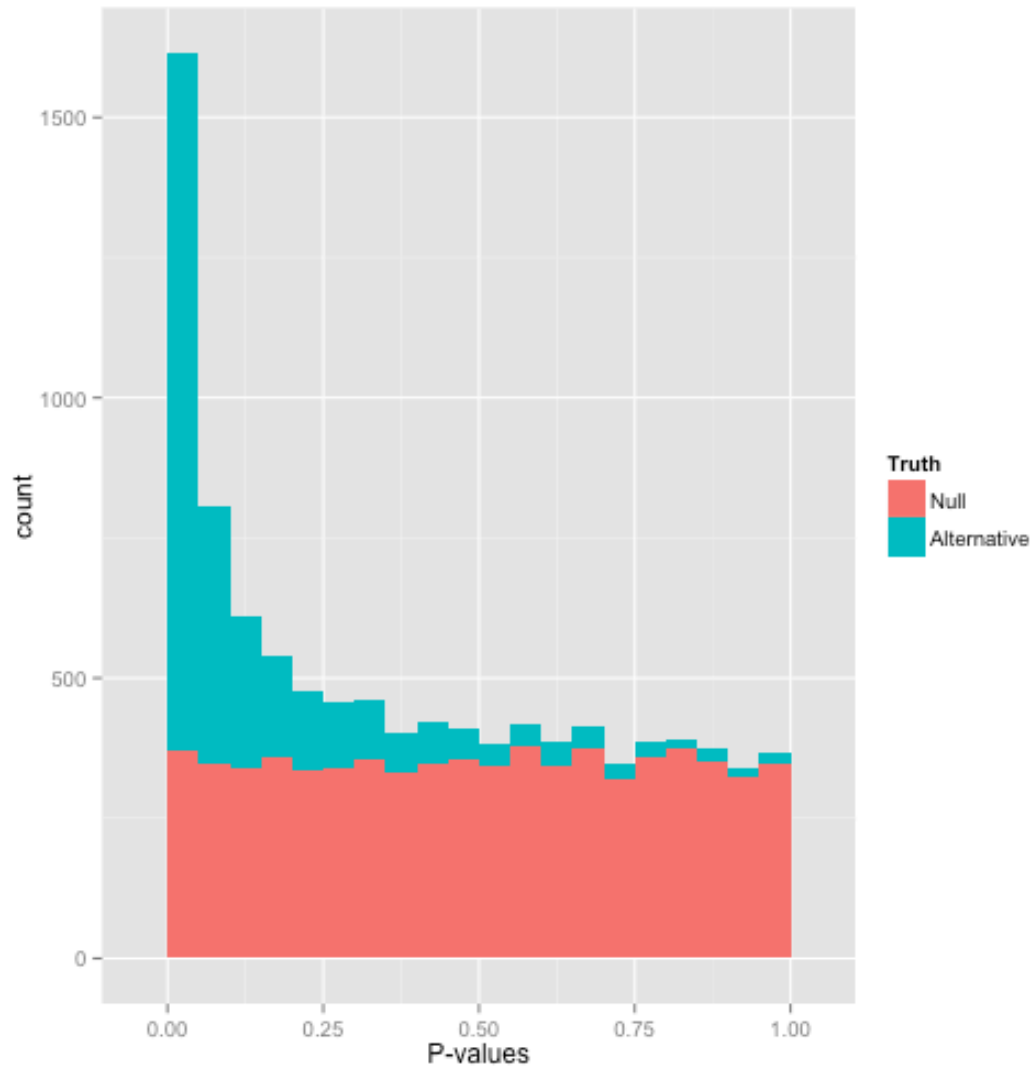
- Mean of normalized counts – averaged over all samples from two conditions
- Log2 of the fold change between two conditions
- Wald Stats – statistical test for hypothesis testing
- **P-value – the probability that the Wald statistic is as extreme as observed if the null hypothesis were true**
- Adjusted P value – accounting for multiple testing correction

# DESeq2 Results table

GeneID	Base mean	log2(FC)	StdErr	Wald-Stats	P-value	P-adj
EGR1	1273.65	-2.22	0.12	-18.65	1.25E-77	<b>1.44E-73</b>
MYC	5226.12	1.41	0.11	12.53	4.95E-36	<b>2.87E-32</b>
OPRK1	78.35	-1.83	0.17	-10.57	4.11E-26	<b>1.59E-22</b>
CCNI2	7427.12	0.93	0.10	9.43	4.27E-21	<b>1.24E-17</b>
STRA6	785.78	0.97	0.11	8.61	7.29E-18	<b>1.69E-14</b>

- Mean of normalized counts – averaged over all samples from two conditions
- Log of the fold change between two conditions
- Wald Stats – statistical test for hypothesis testing
- P-value – the probability that the Wald statistic is as extreme as observed if the null hypothesis were true
- **Adjusted P value – accounting for multiple testing correction**

# DESeq2 P-value histogram



- Histogram of raw p-values for all genes examined
- P-value: Probability of getting a  $\log_2\text{FoldChange}$  as extreme as observed if the true  $\log_2\text{FoldChange} = 0$  for that gene (null hypothesis)

How to interpret:

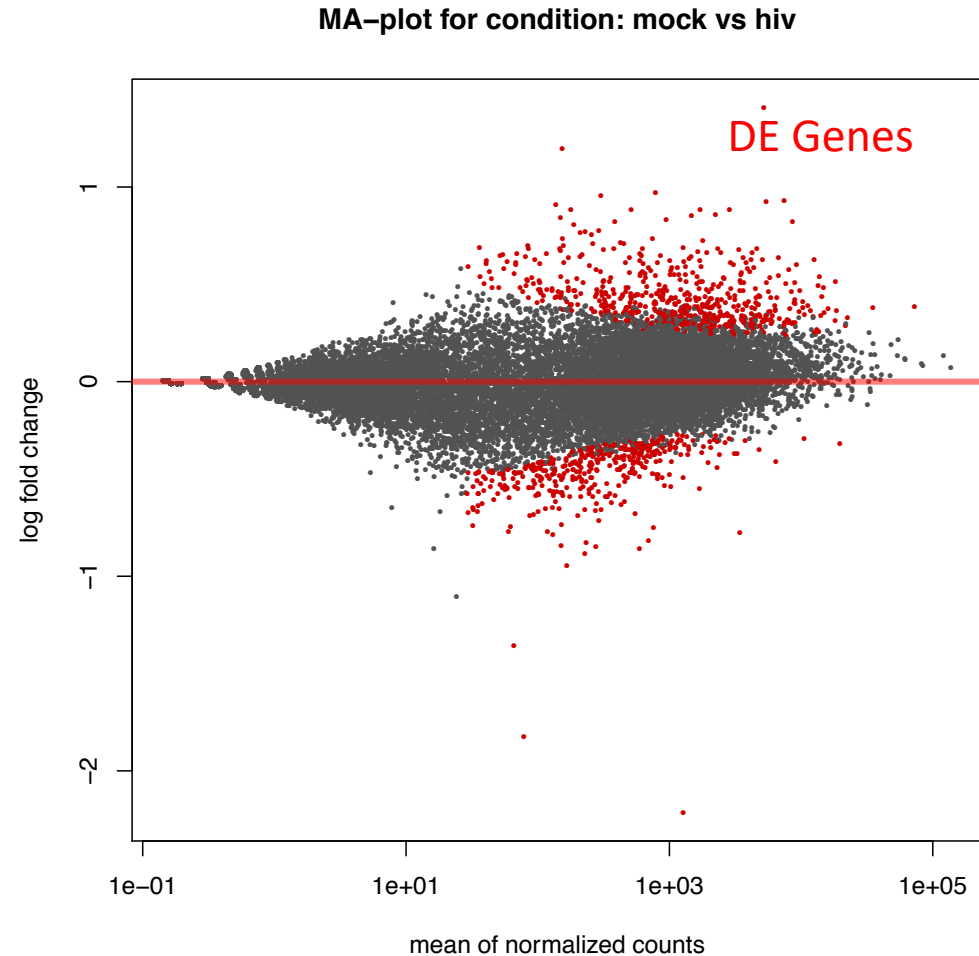
- Random P-values are expected to be uniform, if you have true positives you should see a peak close to zero



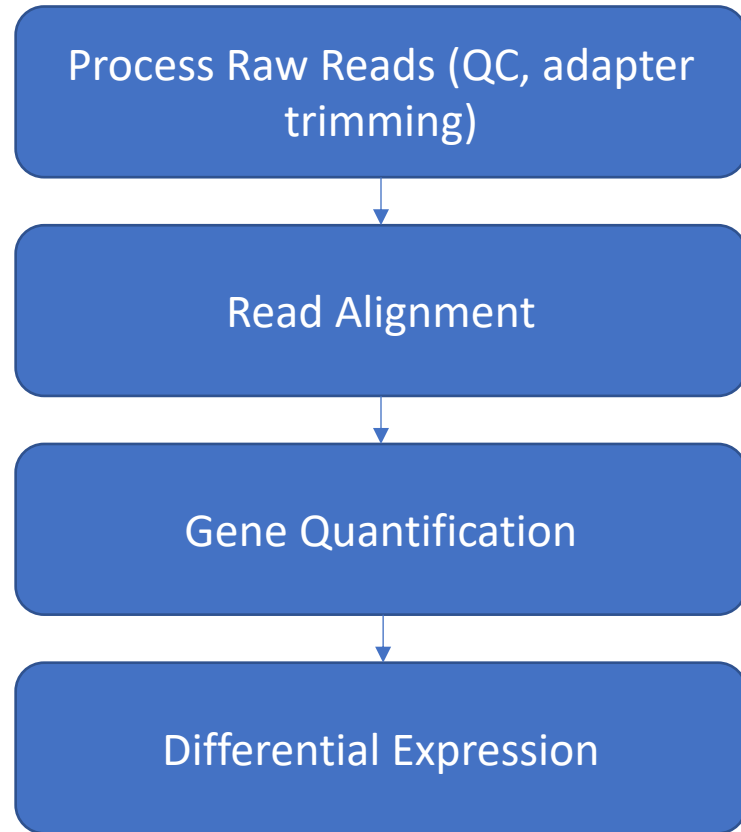
# DESeq2 MA plot

Shows the relationship between

- M: The difference in expression  
 $\text{Log}(\text{HIV}) - \text{Log}(\text{Mock}) = \text{Log}(\text{HIV}/\text{Mock})$
- A: Average expression strength  $\text{Average}(\text{Mock}, \text{HIV})$
- Genes with adjusted  $p$ -value  $< 0.1$  are in red
- Gives an overview of your results



# Conclusions



# 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](https://hbctraining.github.io/DGE_workshop)

Galaxy Training:

[https://galaxyproject.org/tutorials/rb\\_rnaseq/](https://galaxyproject.org/tutorials/rb_rnaseq/)

# Outline

Bulk and single cell  
RNA sequencing

Intro to Galaxy Platform for  
Bioinformatics (Tufts network or  
VPN required)

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

Work through RNAseq  
example together on Galaxy

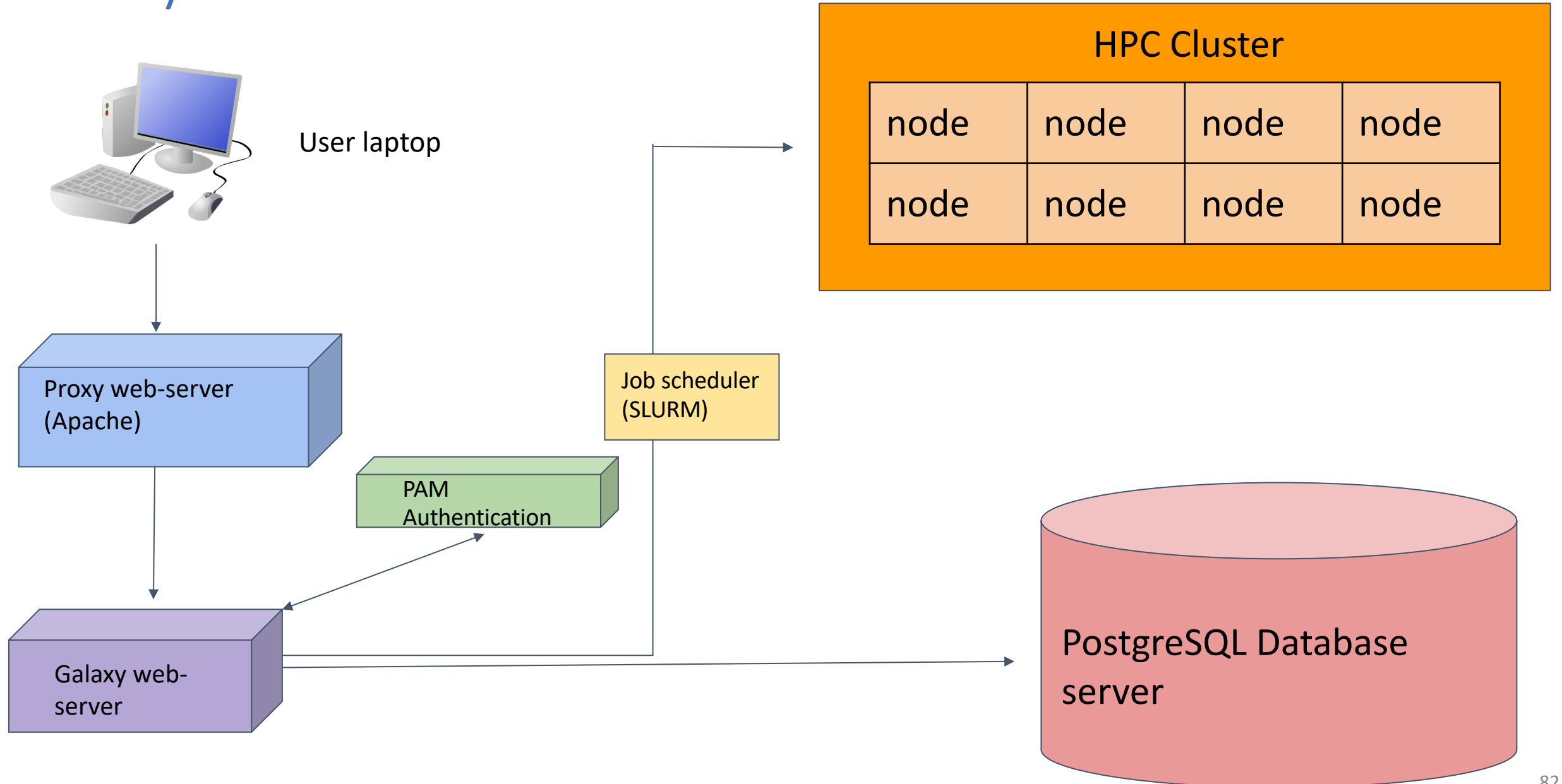
[https://rbatorsky.github.io/in  
tro-to-rnaseq-with-galaxy/](https://rbatorsky.github.io/intro-to-rnaseq-with-galaxy/)

Turn in workshop  
questions on Canvas



- ❖ **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

# Galaxy on the Tufts HPC



# User Interface

The screenshot displays the Galaxy Tufts user interface. At the top, a dark navigation bar contains the 'Galaxy Tufts' logo, navigation links for 'Analyze Data', 'Workflow', 'Visualize', 'Shared Data', 'Admin', 'Help', and 'User', and a 'Using 20%' status indicator. On the left, a sidebar lists tool categories: 'Tools' (with search and star icons), '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', and 'Scratchbook'. Below this, there are two paragraphs of text: 'For information about using Galaxy at Tufts, reference Galaxy documentation, or visit the official GalaxyProject support page.' and 'For more information about Research Technology bioinformatics services, visit the Biotools or email [tts-research@tufts.edu](mailto:tts-research@tufts.edu).' An aerial photograph of the Tufts University campus is shown at the bottom. On the right, a 'History' panel shows 'Unnamed history (empty)' and a message: 'This history is empty. You can load your own data or get data from an external source'.

# User Interface

## TOP MENU BAR

The screenshot displays the Galaxy Tufts web interface. At the top is a blue navigation bar with the 'Galaxy Tufts' logo on the left and a menu of options: 'Analyze Data', 'Workflow', 'Visualize', 'Shared Data', 'Admin', 'Help', and 'User'. A green status indicator on the right shows 'Using 20%'. Below the navigation bar is a left sidebar containing a 'Tools' section with a search bar and a list of tool 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 large heading 'Welcome to Galaxy on the Tufts University High Performance Compute Cluster!' and a 'Tufts Galaxy Support»' button. Below this, there are links for 'Take an interactive tour: Galaxy UI', 'History', and 'Scratchbook'. Two paragraphs of text provide information about using Galaxy at Tufts and contacting research services. An aerial photograph of the Tufts University campus is shown at the bottom. The right sidebar contains a 'History' panel with a search bar, an 'Unnamed history' section that is currently empty, and an informational message: 'This history is empty. You can load your own data or get data from an external source'.



# User Interface

## TOP MENU BAR

The screenshot displays the Galaxy user interface. At the top, a blue navigation bar contains the following items: a hamburger menu icon, the word "TOOLS" in green, and a series of menu items: "Analyze Data", "Workflow", "Visualize", "Shared Data", "Admin", "Help", "User", and a grid icon. On the far right of the navigation bar, the text "g 20%" is visible. Below the navigation bar, the main content area is divided into three sections. On the left is a "Tools" sidebar with a search bar and a list of tool 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 "RNA-seq" and "SAMTOOLS" categories are highlighted with a green border. The central content area features a "Welcome to Galaxy on the Tufts University High Performance Compute Cluster!" message, a "Tufts Galaxy Support»" button, and a "Take an interactive tour:" section with buttons for "Galaxy UI", "History", and "Scratchbook". Below this, there are two paragraphs of text and a large aerial photograph of the Tufts University campus. On the right side, a "History" panel is highlighted with a red border. It contains a "search datasets" search bar, the text "Unnamed history (empty)", and a blue information box that reads: "This history is empty. You can load your own data or get data from an external source".

# User Interface

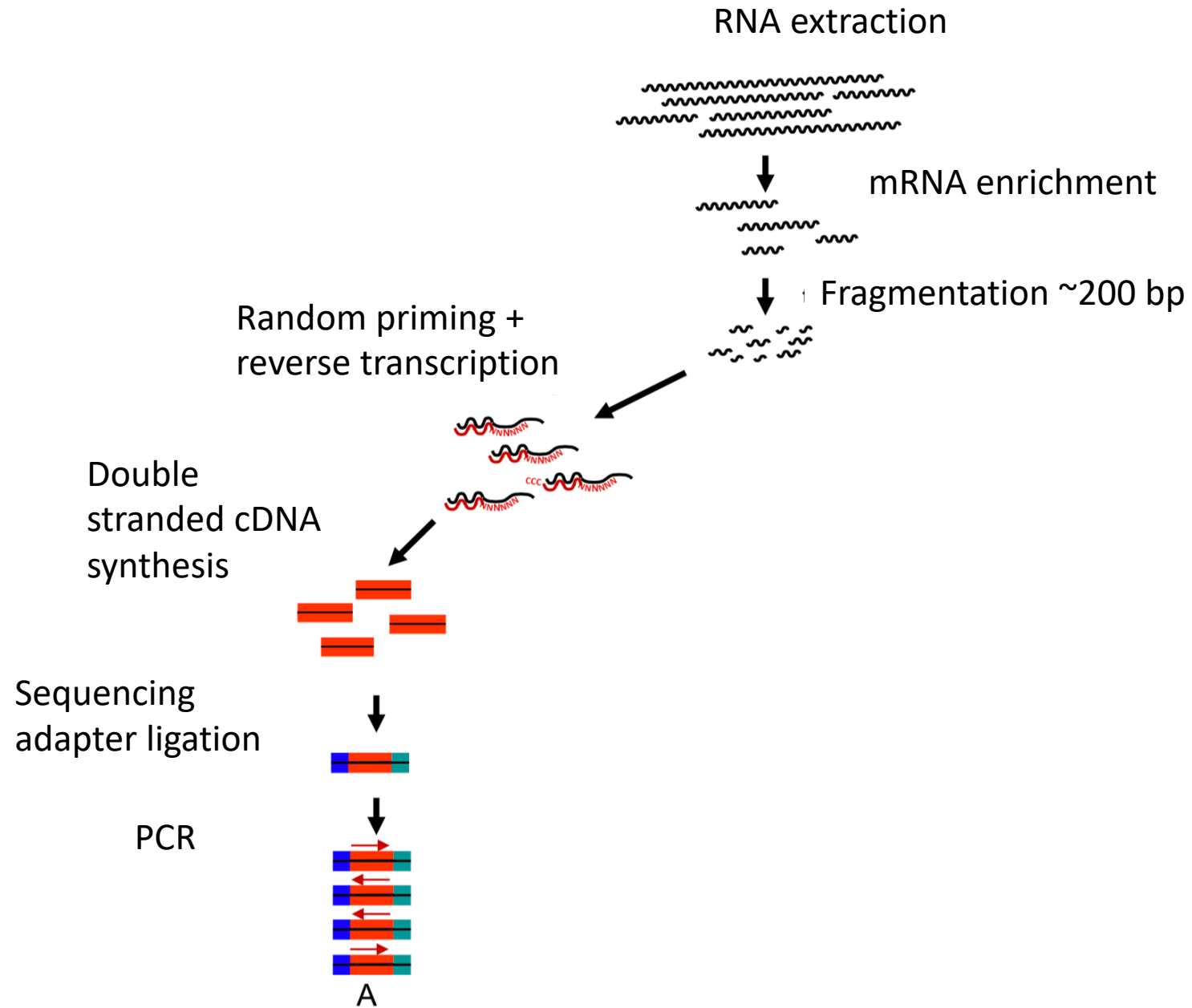
## TOP MENU BAR

The screenshot displays the Galaxy web interface. At the top, a dark blue navigation bar contains the following menu items: Analyze Data, Workflow, Visualize, Shared Data, Admin, Help, and User. A hamburger menu icon is on the left, and a grid icon is on the right. A '20%' zoom indicator is visible in the top right corner.

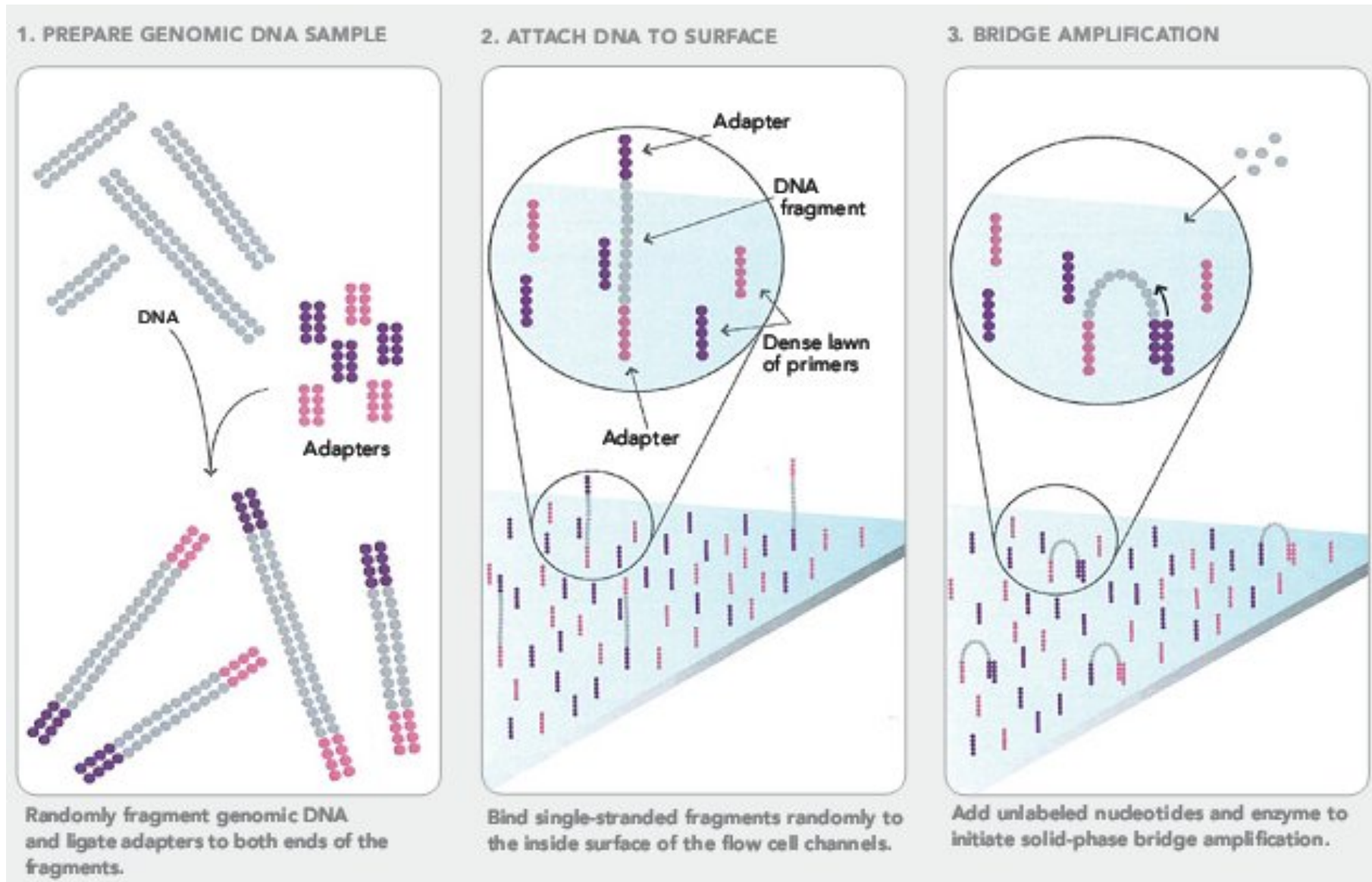
On the left side, a 'TOOLS' sidebar is highlighted with a green border. It features a search bar labeled 'search tools' and a list of tool categories: Get Data, Send Data, Collection Operations, Expression Tools, Lift-Over, Text Manipulation, Convert Formats, Filter and Sort, Join, Subtract and Group, Fetch Alignments/Sequences, Operate on Genomic Intervals, Statistics, Graph/Display Data, Phenotype Association, FASTQ Quality Control, RNA-seq, and SAMTOOLS. The 'RNA-seq' and 'SAMTOOLS' items are highlighted with a green box.

The central 'MAIN' area is highlighted with a purple border. It contains a welcome message: 'Welcome to Galaxy on the Tufts University High Performance Compute Cluster!'. Below this is a 'Tufts Galaxy Support»' button. A section titled 'Take an interactive tour:' includes buttons for 'Galaxy UI', 'History', and 'Scratchbook'. Further down, there are two paragraphs of text: 'For information about using Galaxy at Tufts, reference Galaxy documentation, or visit the official GalaxyProject support page.' and 'For more information about Research Technology bioinformatics services, visit the Biotools or email [tts-research@tufts.edu](mailto:tts-research@tufts.edu)'. At the bottom of this section is an aerial photograph of the Tufts University campus.

On the right side, a 'HISTORY' sidebar is highlighted with a red border. It has a search bar labeled 'search datasets' and shows 'Unnamed history (empty)'. A blue information box contains the text: 'This history is empty. You can load your own data or get data from an external source'. The sidebar also includes icons for refresh, add, and settings.

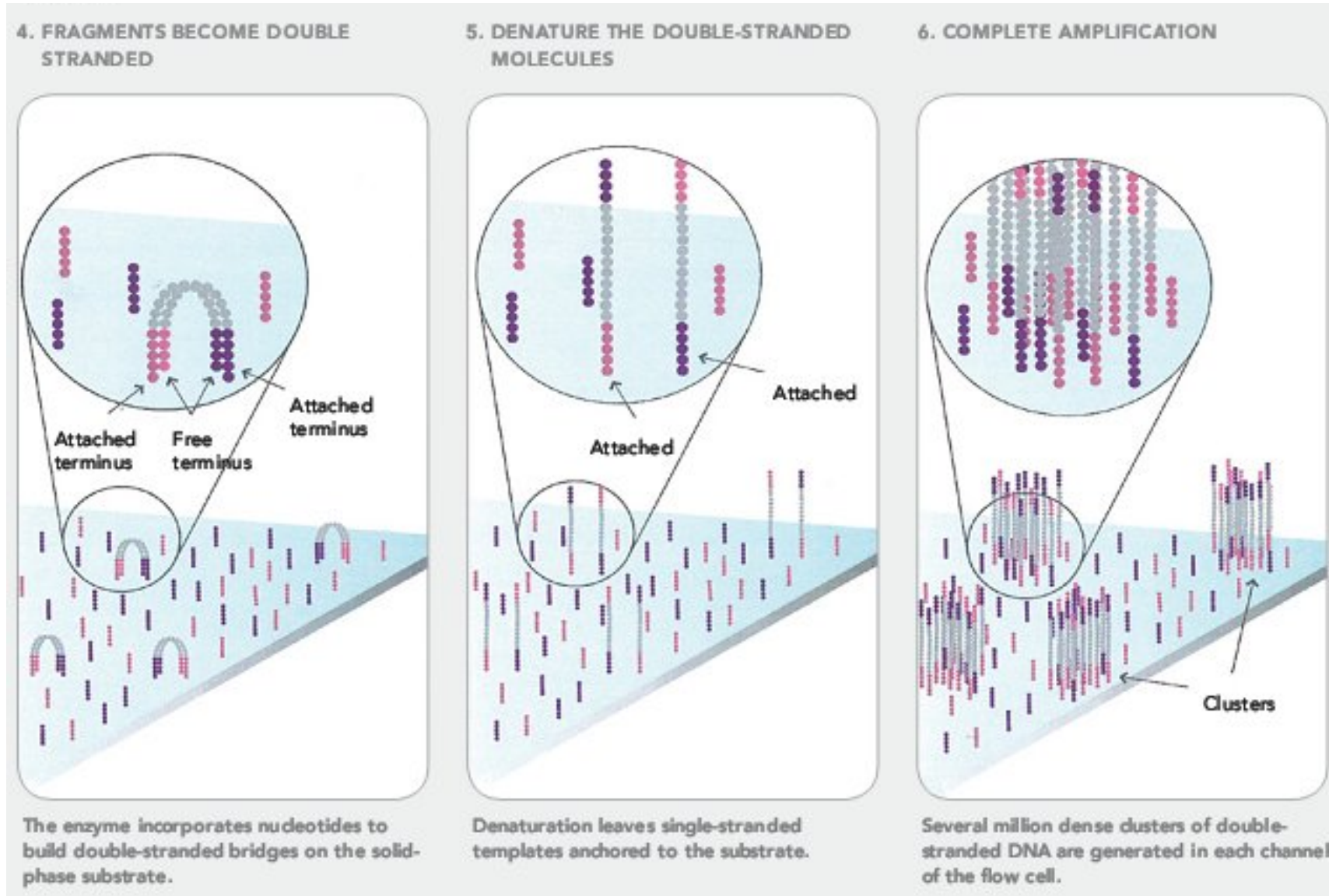


# Next Generation Sequencing (NGS)

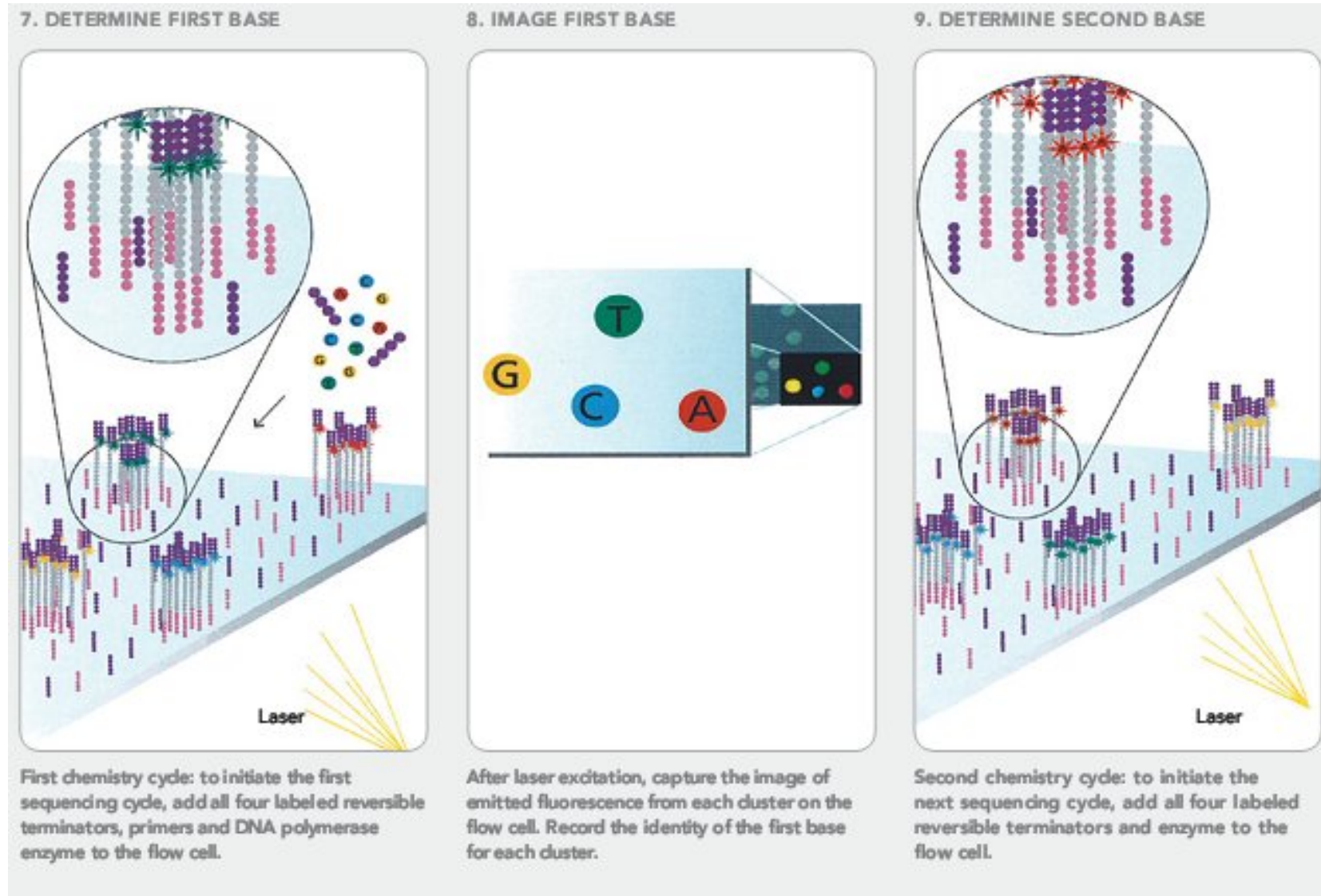




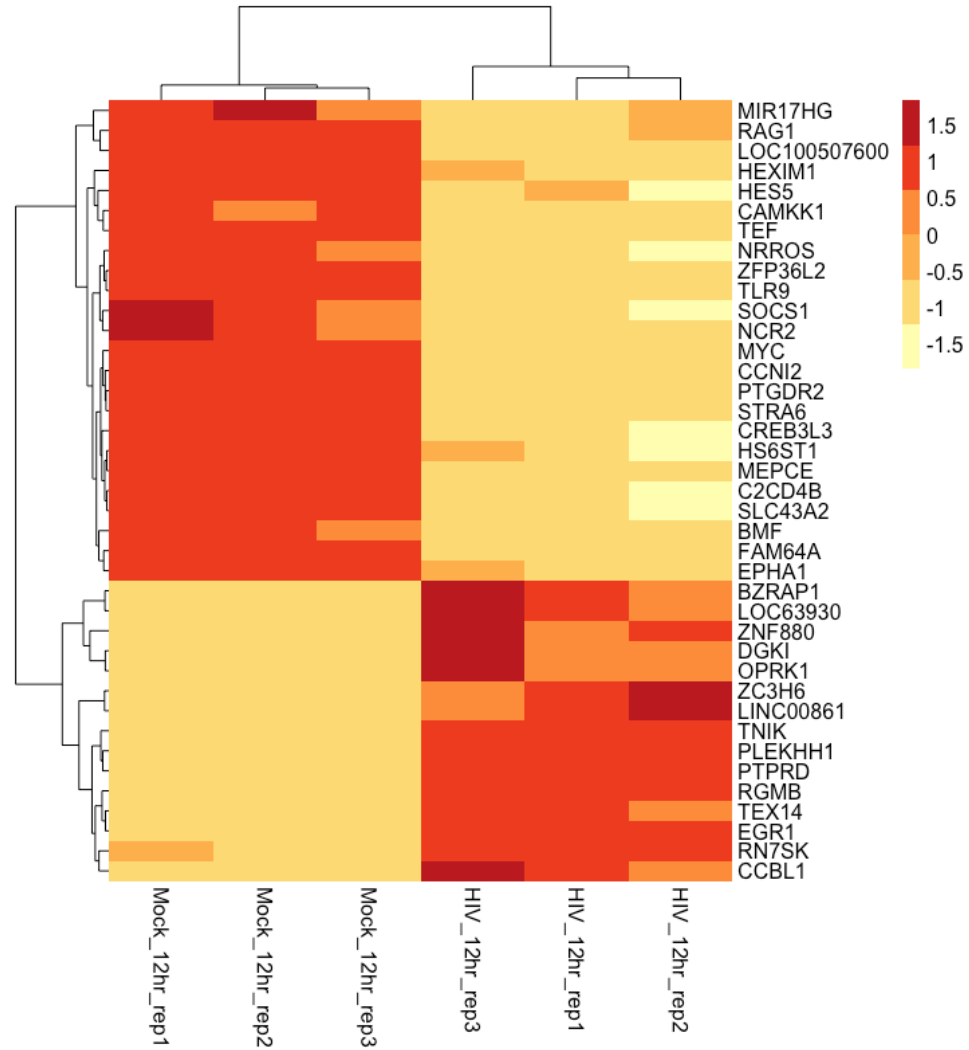
# Next Generation Sequencing (NGS)



# Next Generation Sequencing (NGS)



# Final Heatmap – not part of DESeq2 output



## Common RNAseq analysis goals

- Novel transcript discovery
- Transcriptome assembly
- Single cell analysis
- Quantify alternative splicing
- **Differential Expression**

Replace with actual heatmap

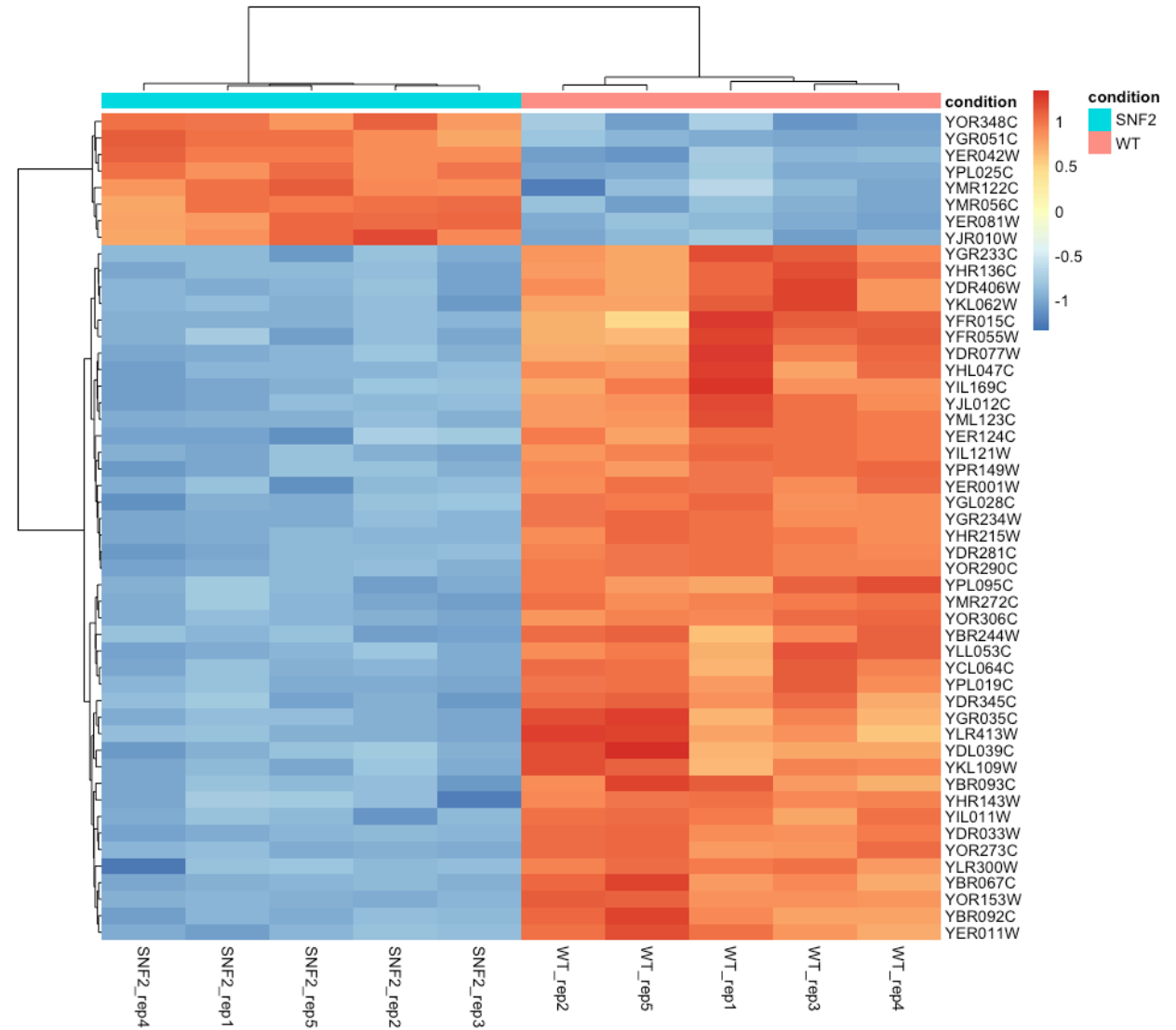
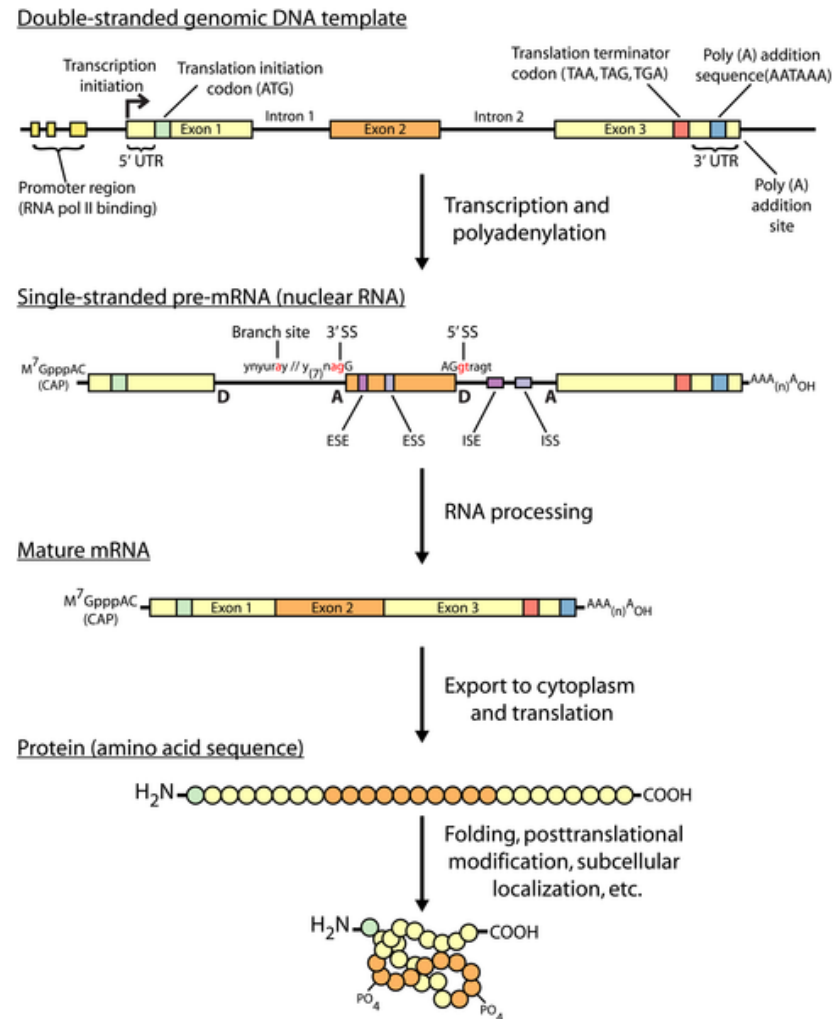


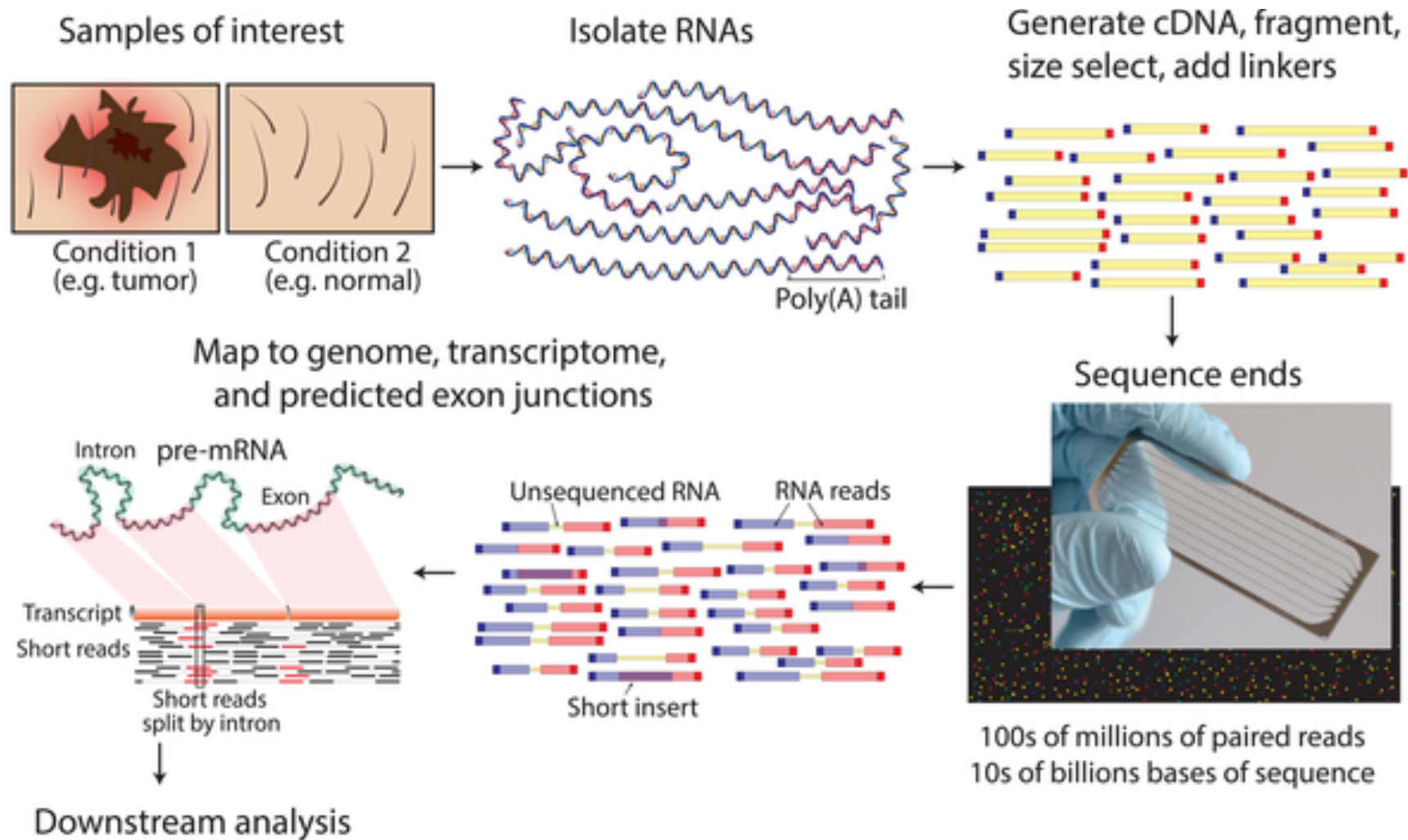


Fig 1. An overview of the central dogma of molecular biology.



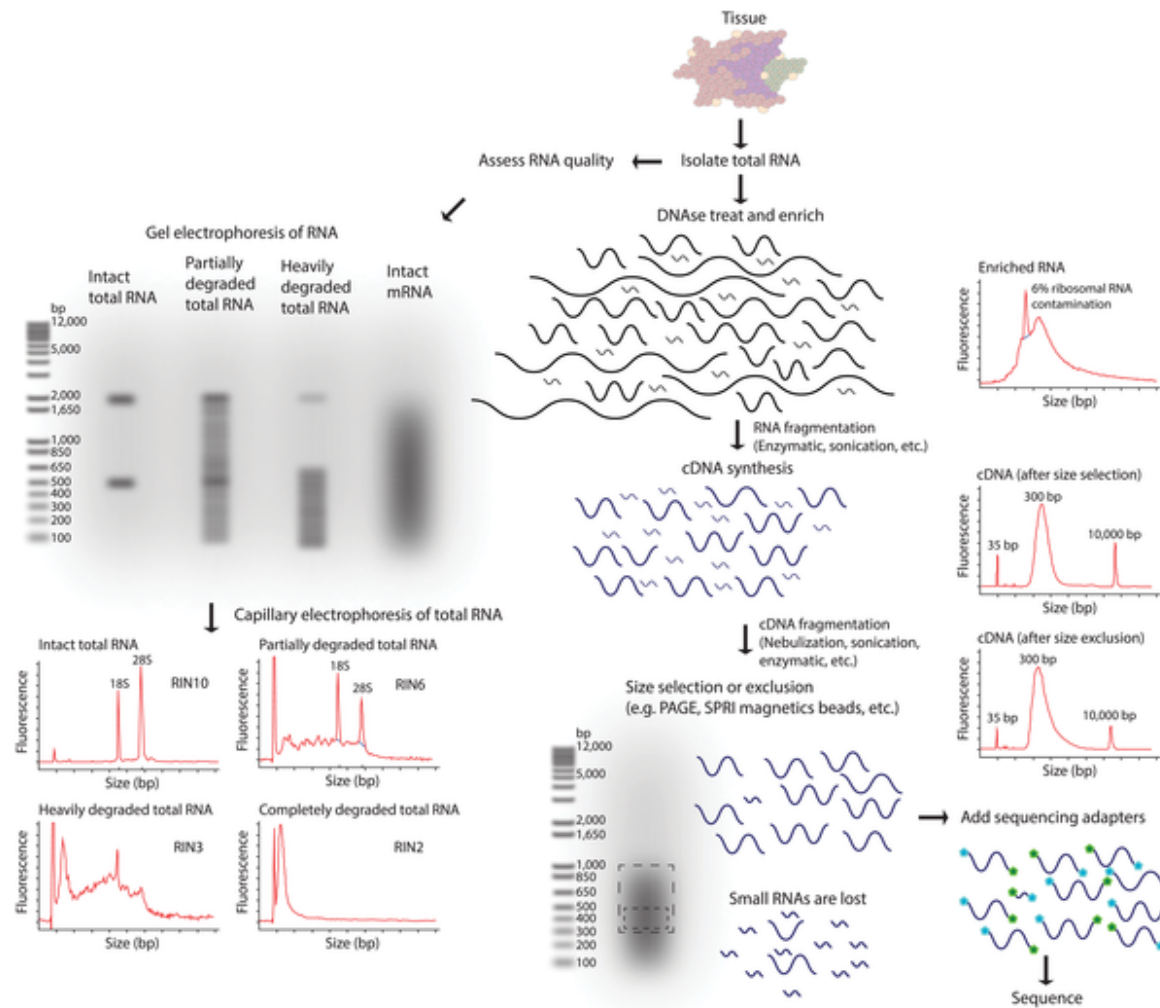
Griffith M, Walker JR, Spies NC, Ainscough BJ, Griffith OL (2015) Informatics for RNA Sequencing: A Web Resource for Analysis on the Cloud. PLOS Computational Biology 11(8): e1004393. <https://doi.org/10.1371/journal.pcbi.1004393>  
<https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1004393>

Fig 2. RNA-seq data generation.



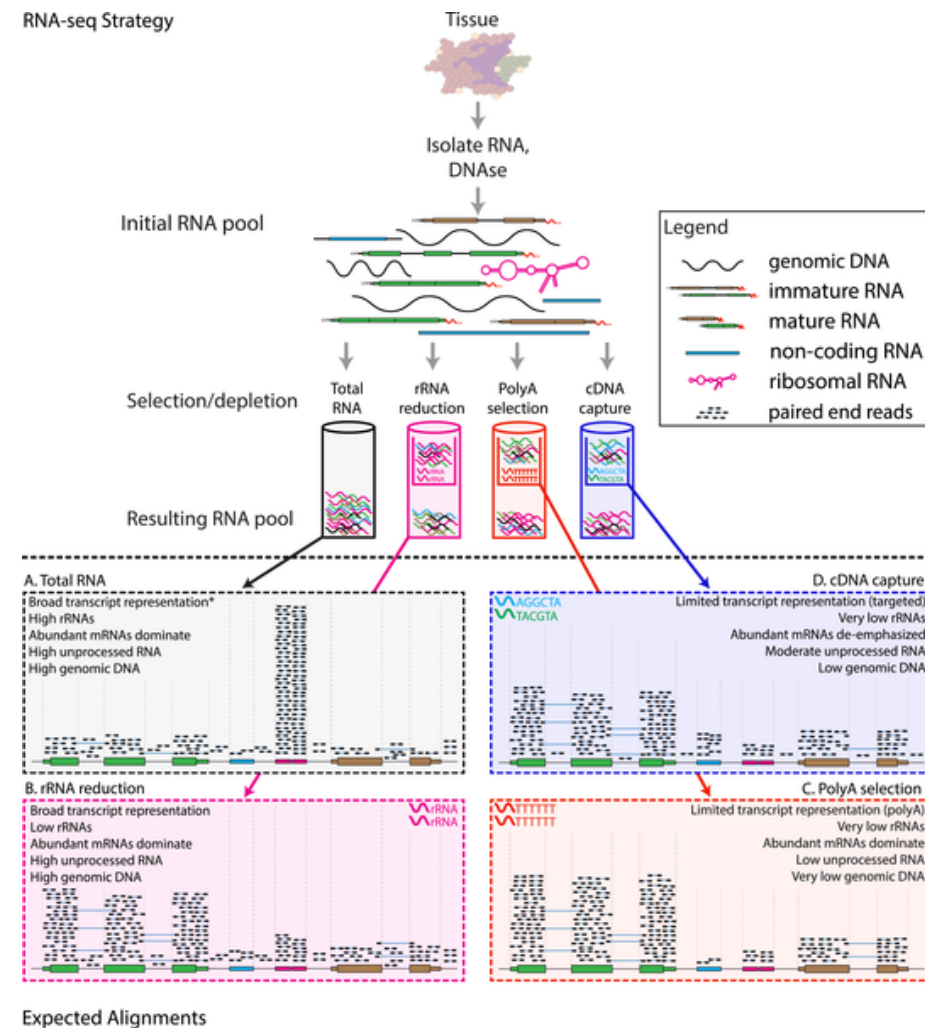
Griffith M, Walker JR, Spies NC, Ainscough BJ, Griffith OL (2015) Informatics for RNA Sequencing: A Web Resource for Analysis on the Cloud. PLOS Computational Biology 11(8): e1004393. <https://doi.org/10.1371/journal.pcbi.1004393>  
<https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1004393>

Fig 3. RNA-seq library fragmentation and size selection strategies that influence interpretation and analysis.



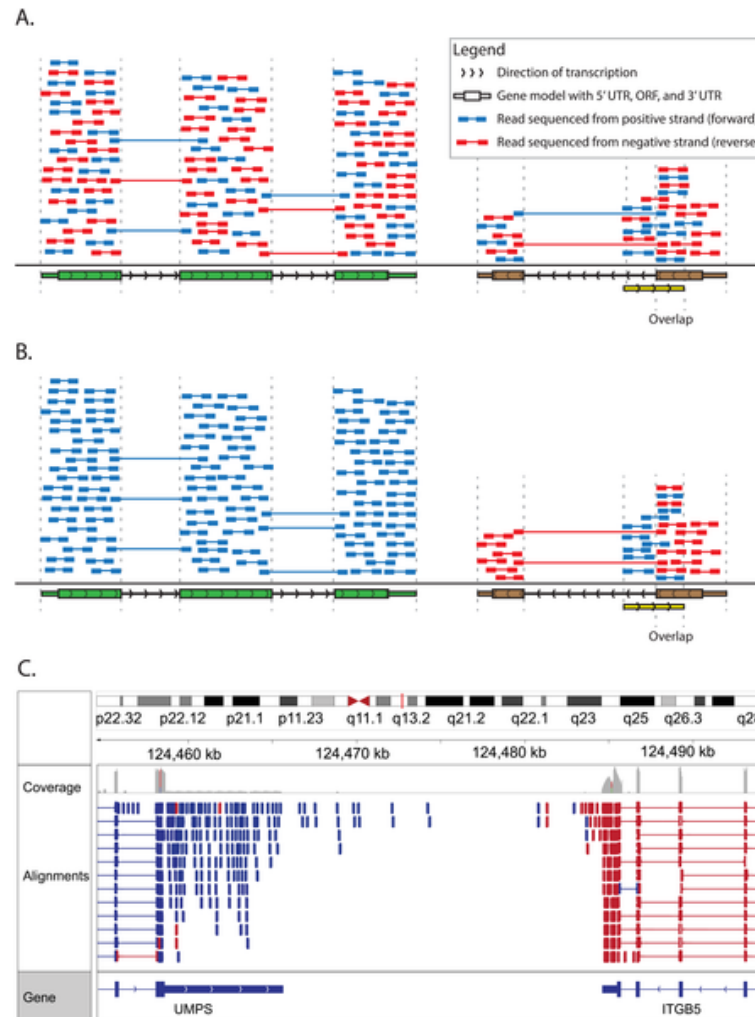
Griffith M, Walker JR, Spies NC, Ainscough BJ, Griffith OL (2015) Informatics for RNA Sequencing: A Web Resource for Analysis on the Cloud. PLOS Computational Biology 11(8): e1004393. <https://doi.org/10.1371/journal.pcbi.1004393>  
<https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1004393>

**Fig 4. RNA-seq library enrichment strategies that influence interpretation and analysis.**



Griffith M, Walker JR, Spies NC, Ainscough BJ, Griffith OL (2015) Informatics for RNA Sequencing: A Web Resource for Analysis on the Cloud. PLOS Computational Biology 11(8): e1004393. <https://doi.org/10.1371/journal.pcbi.1004393>  
<https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1004393>

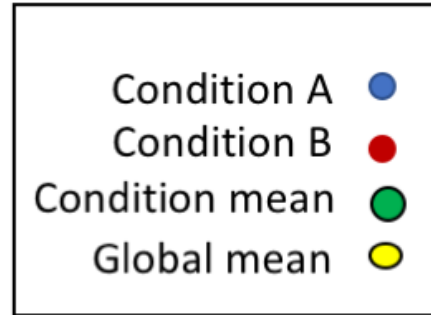
**Fig 6. Comparison of stranded and unstranded RNA-seq library methods and their influence on interpretation and analysis.**



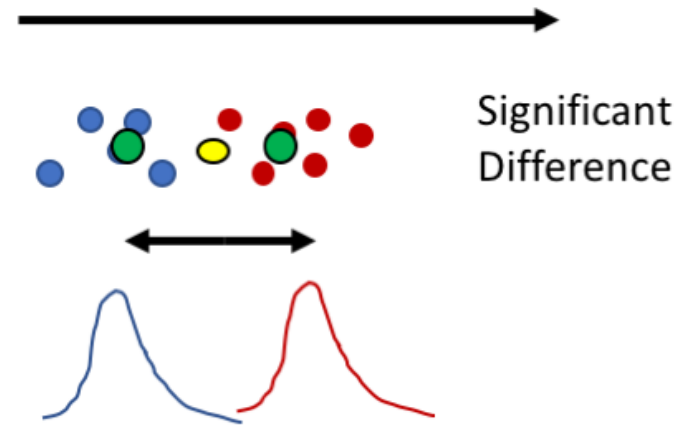
Griffith M, Walker JR, Spies NC, Ainscough BJ, Griffith OL (2015) Informatics for RNA Sequencing: A Web Resource for Analysis on the Cloud. PLOS Computational Biology 11(8): e1004393. <https://doi.org/10.1371/journal.pcbi.1004393>  
<https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1004393>

# Test for Differential Expression

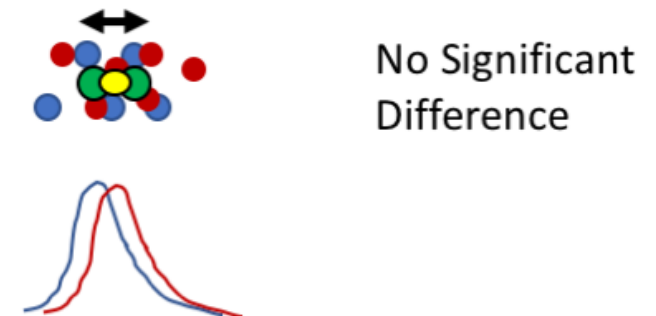
DESeq2 will seek to fit a probability distribution to each gene we measured and perform a statistical test to determine whether there is a difference between conditions



Expression Level of a Gene



Deviation from global mean



# Reference-based vs Reference-free RNAseq

RNAseq can be roughly divided into two "types":

- **Reference genome-based** - an assembled genome exists for a species for which an RNAseq experiment is performed. It allows reads to be aligned against the reference genome and significantly improves our ability to reconstruct transcripts. This category would obviously include humans and most model organisms
- **Reference genome-free** - no genome assembly for the species of interest is available. In this case one would need to assemble the reads into transcripts using *de novo* approaches. This type of RNAseq is as much of an art as well as science because assembly is heavily parameter-dependent and difficult to do well.

In this lesson we will focus on the **Reference genome-based** type of RNA seq.



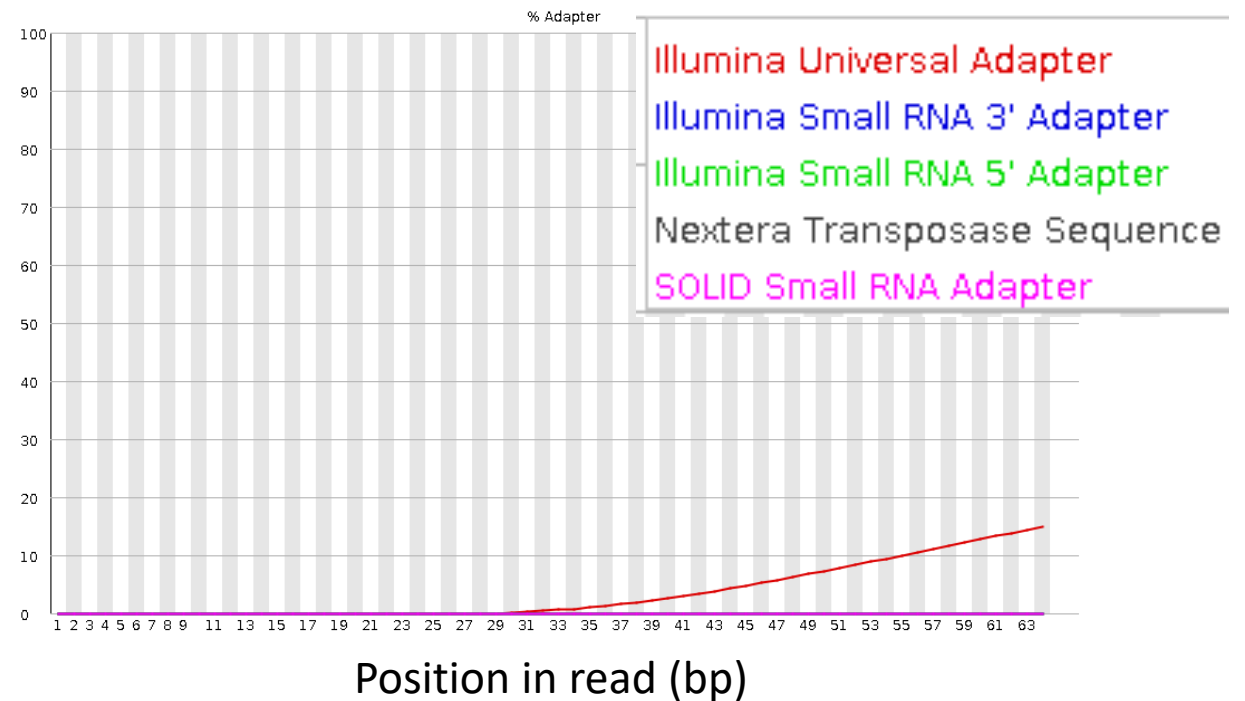
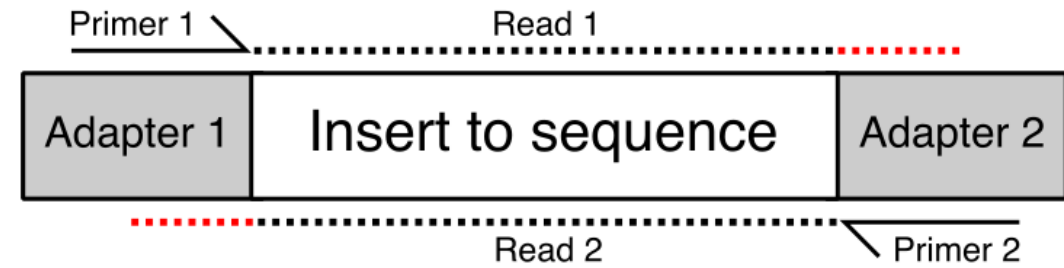
# FastQC: Adapter content

The cause: The “insert” sequence is shorter than the read, and the read contains part of the adapter sequence

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!

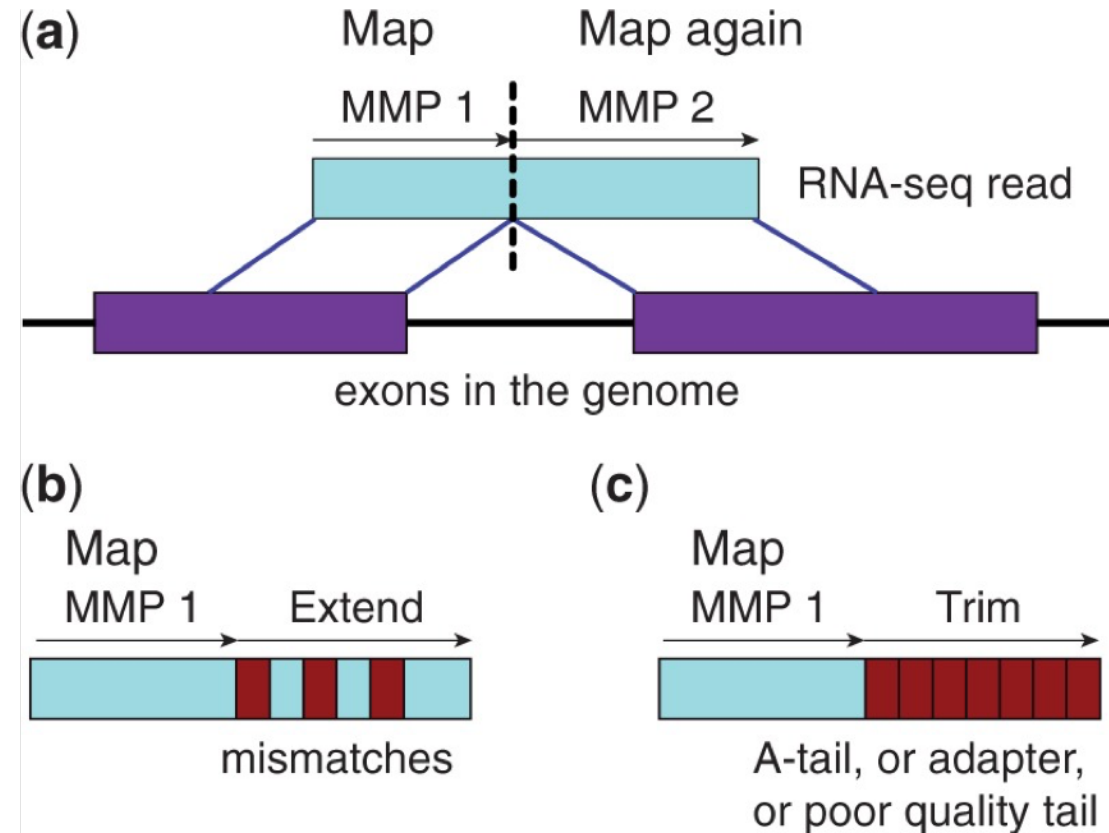




# STAR Aligner (Spliced Transcripts Alignment to a Reference)

Highly accurate, memory intensive **aligner**  
Two phase mapping process

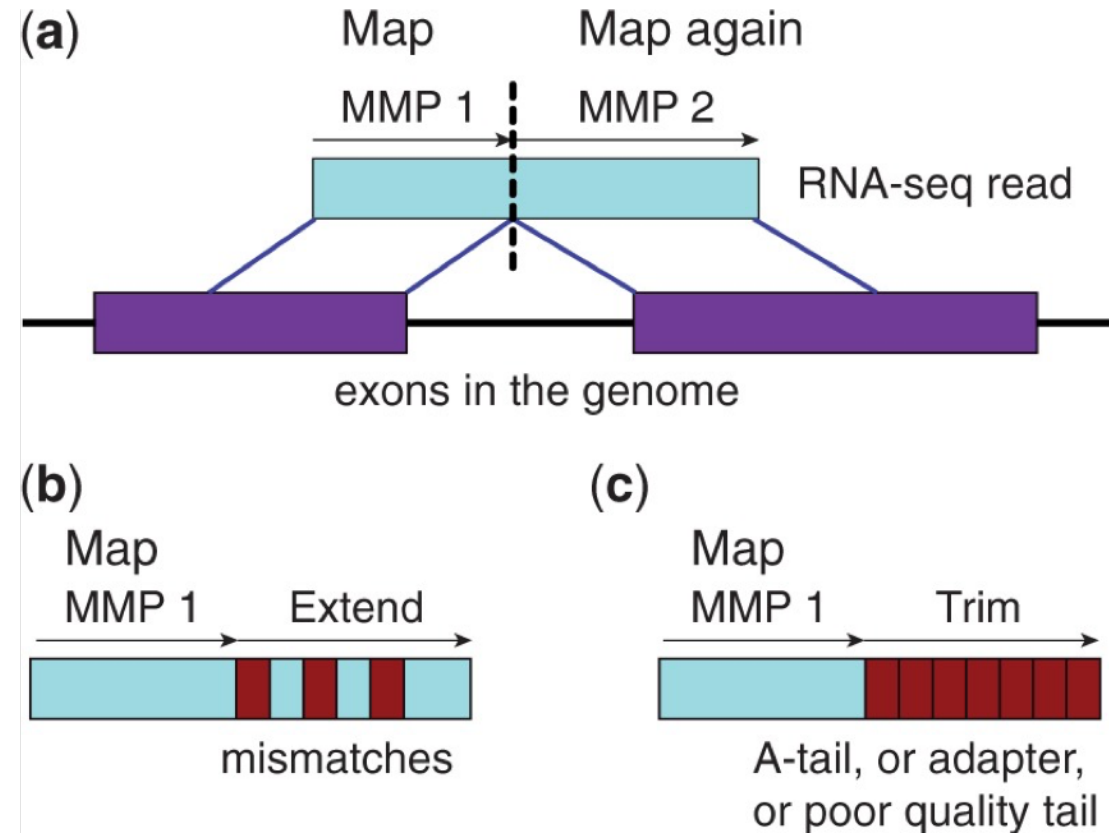
1. Find Maximum Mappable Prefixes (MMP) in a read. MMP can be extended by
  - mismatches
  - Indels
  - soft-clipping



# STAR Aligner (Spliced Transcripts Alignment to a Reference)

Highly accurate, memory intensive **aligner**  
Two phase mapping process

1. Find Maximum Mappable Prefixes (MMP) in a read. MMP can be extended by
  - mismatches
  - Indels
  - soft-clipping
2. Clustering MMP, stitching and scoring to determine final read location

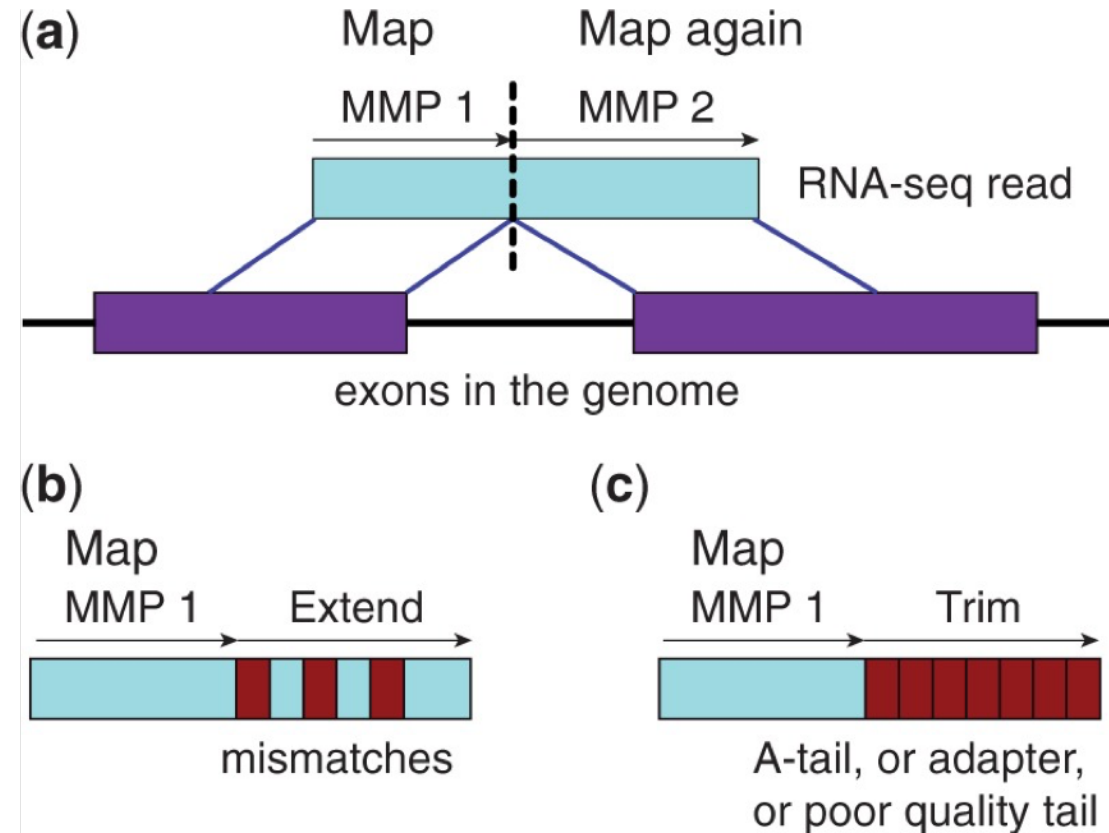


# STAR Aligner (Spliced Transcripts Alignment to a Reference)

Highly accurate, memory intensive **aligner**  
Two phase mapping process

1. Find Maximum Mappable Prefixes (MMP) in a read. MMP can be extended by
  - mismatches
  - Indels
  - soft-clipping
2. Clustering MMP, stitching and scoring to determine final read location

Output is a Sequence Alignment Map (SAM) file

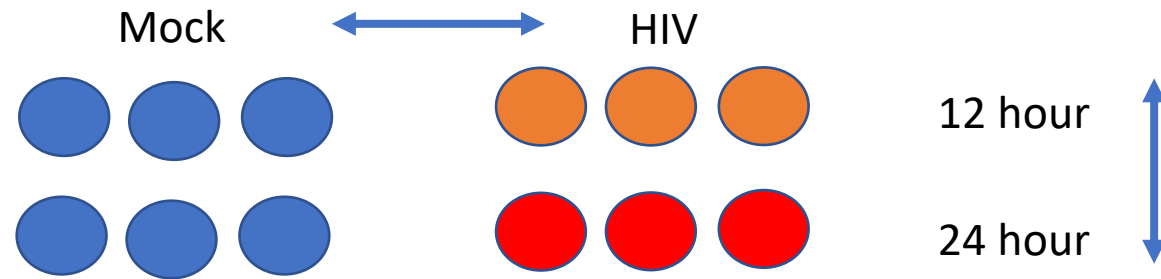


# Tracking read numbers

Revisit quality control after each processing step!

Number of Reads	Source	Result
Raw reads	FastQC run 1	8 M
After Trimming	FastQC run 2	7.1 M
Aligned to genome	STAR log	6 M
Associated with genes	FeatureCounts log	5.4 M

# Multi-factor experiment design

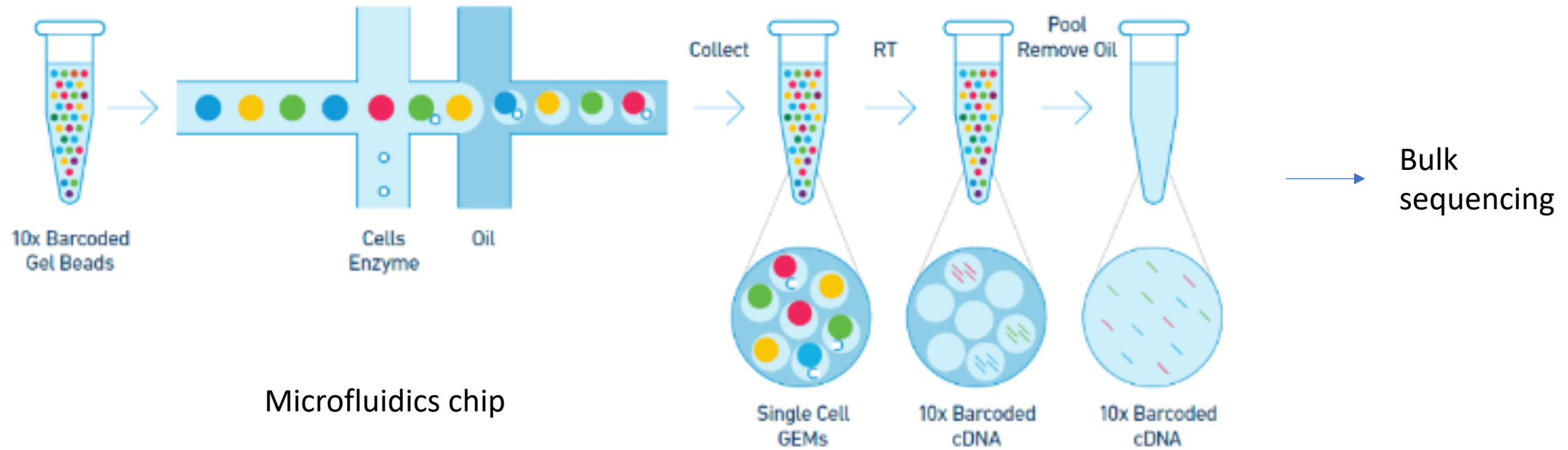


Factor 1:  
Infection status (Mock or HIV)

Factor 2:  
Time (12 or 24 hr)

Sample	Condition	Time
1	Mock	12
2	Mock	12
3	Mock	12
4	Mock	24
5	Mock	24
6	Mock	24
7	HIV	12
8	HIV	12
9	HIV	12
10	HIV	24
11	HIV	24
12	HIV	24

# 10x single cell technology



Beads are barcoded, and RT occurs inside the GEM -> all reads from a given GEM will have the same barcode

