

An Introduction to RNA-sequencing

GSND 5340Q, BMDA

W. Evan Johnson, Ph.D.
Professor, Division of Infectious Disease
Director, Center for Data Science
Rutgers University – New Jersey Medical School

2024-05-19

Installing R Packages:

Install the following tools: Rsubread, Rsamtools, edgeR, DESeq2, sva
SummarizedExperiment, ComplexHeatmap, umap, and the
TBSignatureProfiler. We will also need help from the tidyverse.

```
if (!requireNamespace("BiocManager", quietly = TRUE))
  install.packages("BiocManager")
BiocManager::install(c("Rsubread", "Rsamtools", "tidyverse",
  "SummarizedExperiment", "edgeR", "DESeq2", "sva", "ComplexHeatmap",
  "TBSignatureProfiler", "umap"))
```

Installing and using the SCTK

```
install.packages("devtools")
devtools::install_github("wewanjohnson/singleCellTK")
library(singleCellTK)
singleCellTK()

### Example: open downstream_analysis/
### features_combined.txt and meta_data.txt
```

Load Packages for RNA-Seq

We will be using the following packages for our RNA-seq lecture:

```
library(tidyverse) ## tools for data wrangling
library(Rsubread) ## alignment and feature counts
library(Rsamtools) ## managing .sam and .bam files
library(SummarizedExperiment) ## managing counts data
library(edgeR) ## differential expression
library(DESeq2) ## differential expression
library(ComplexHeatmap) ## Heatmap visualization
library(TBSignatureProfiler) ## TB signature analysis
library(umap) ## dimension reduction and plotting data
```

Objective

- Disclaimer: non-comprehensive introduction to RNA-sequencing
- Introduce preprocessing steps
- Visualization
- Analytical methods
- Common software tools

Steps to an RNA-seq Analysis (Literacy)

① Preprocessing and QC:

- Fasta and Fastq files
- FastQC: good vs. bad examples
- Visualization

② Alignment

- Obtaining genome sequence and annotation
- Software: Bowtie, TopHat, STAR, Subread/Rsubread

③ Expression Quantification

- Count reads hitting genes, etc
- Approaches/software: HT-Seq, STAR, Cufflinks, RPKM FPKM or CPM, RSEM, edgeR, findOverlaps (GenomicRanges). featureCounts (Rsubread)

Steps to an RNA-seq Analysis (Literacy)

④ More visualization

- Heatmaps, boxplots, PCA, t-SNE, UMAP

⑤ Differential Expression

- Batch correction
- Overdispersion
- General Workflow
- Available tools: edgeR, DESeq, Limma/voom
- Even more visualization!!

Illumina Sequencing Workflow

1

Library Preparation



Fragment DNA
Repair ends
Add A overhang
Ligate adapters
Purify

2

Cluster Generation



Hybridize to flow cell
Extend hybridized template
Perform bridge amplification
Prepare flow cell for sequencing



3

Sequencing



Perform sequencing
Generate base calls



4

Data Analysis



Images
Intensities
Reads
Alignments

Sequencing Data Formats

Genome sequencing data is often stored in one of two formats, FASTA and FASTQ text files. For example a FASTA file looks like the following:

```
>chrX
ttgaactcctgacctcaggtgatccgcggcctgacccaaaggcgct
cctgcctcagcctcccggtagctggactacaggtcgtccaccatgc
....
caggctaattttgtattttagtagagacgggtttcaccatgttagc
caggatggtctcaatctcctgacccatgatccgcctgcctcgccccc
>chrY
tgtacacttaaatgggtgaatttatggaatgtgaattataCGTGTG
CTTGTAAAAAAATGATGGAGATGGAGACGTGACTCTAGCGTGAAGGGG
...
GTGGGGAGAGTAGATCTAGAGTGGAGACACCACTTTAGGAGGTATGATC
cctgccaccatgcgttggtaattttgtattttagtagagacagggtt
```

FASTQ Files

We can also store confidence or quality scores using a FASTQ format:

FASTQ Encoding

In order to translate FASTQ 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, 40)
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)

FASTQ Probability

And now converting to confidence probabilities:

> Sanger	> Solexa	> Illumina1.3	> Illumina1.5
ASCII Quality	ASCII Quality	ASCII Quality	ASCII Quality
! 33 0.0000	; 59 0.2403	@ 64 0.5000	B 66 0.3690
" 34 0.2057	< 60 0.2847	A 65 0.5573	C 67 0.4988
# 35 0.3690	= 61 0.3339	B 66 0.6131	D 68 0.6019
\$ 36 0.4988	> 62 0.3869	C 67 0.6661	E 69 0.6838
% 37 0.6019	? 63 0.4427	D 68 0.7153	F 70 0.7488
& 38 0.6838	@ 64 0.5000	E 69 0.7597	G 71 0.8005
' 39 0.7488	A 65 0.5573	F 70 0.7992	H 72 0.8415
(40 0.8005	B 66 0.6131	G 71 0.8337	I 73 0.8741
) 41 0.8415	C 67 0.6661	H 72 0.8632	J 74 0.9000
* 42 0.8741	D 68 0.7153	I 73 0.8882	K 75 0.9206
+ 43 0.9000	E 69 0.7597	J 74 0.9091	L 76 0.9369
, 44 0.9206	F 70 0.7992	K 75 0.9264	M 77 0.9499
- 45 0.9369	G 71 0.8337	L 76 0.9406	N 78 0.9602
. 46 0.9499	H 72 0.8632	M 77 0.9523	O 79 0.9684
/ 47 0.9602	I 73 0.8882	N 78 0.9617	P 80 0.9749
0 48 0.9684	J 74 0.9091	O 79 0.9693	Q 81 0.9800
	K 75 0.9264	P 80 0.9755	R 82 0.9842
	L 76 0.9406	Q 81 0.9804	S 83 0.9874

Preprocessing and QC using FASTQC

FastQC provides a simple way to do QC checks on raw sequence data:

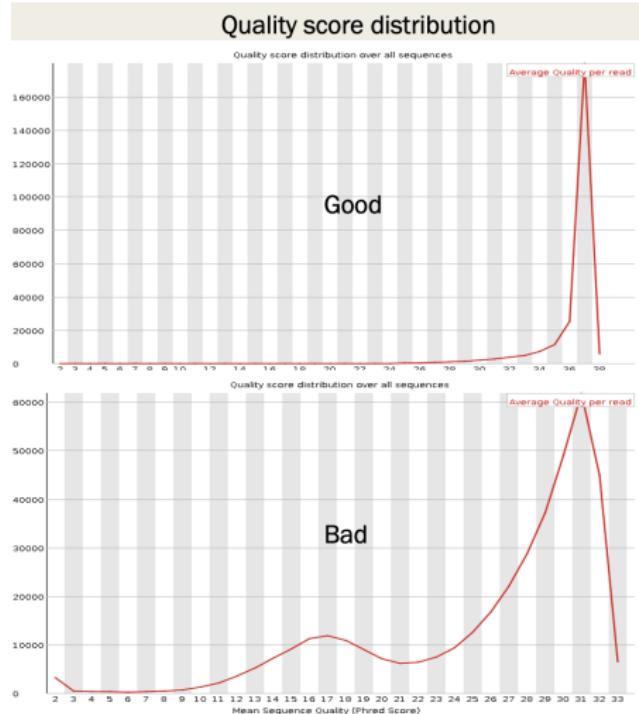
- Import of data from BAM, SAM or FastQ files
- Quick overview and summary graphs and tables to quickly assess your data
- Export of results to an HTML based permanent report
- Offline operation to allow automated generation of reports without running the interactive application

Preprocessing and QC using FASTQC

To run FastQC you can launch the GUI app, or run from the command line:

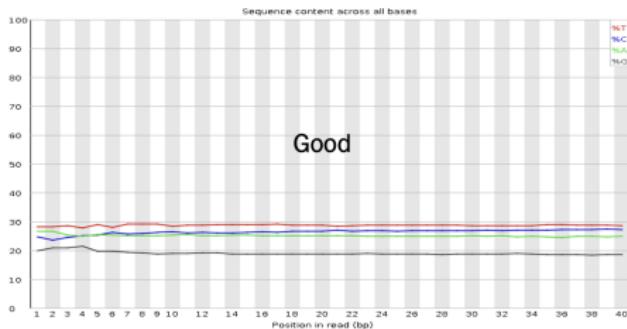
```
rna_seq/FastQC/./fastqc \
  rna_seq/reads/R01_10_short500K.fq.gz
```

FastQC Score Distribution

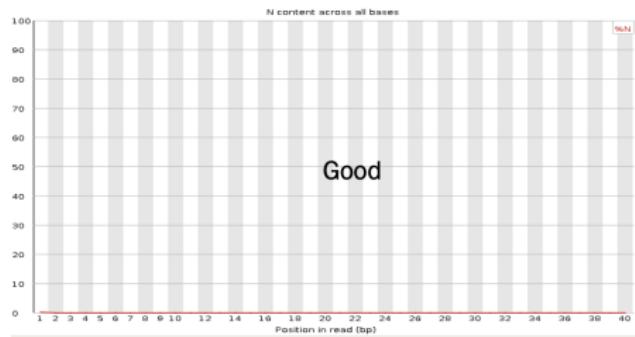


FastQC Base and N Distribution

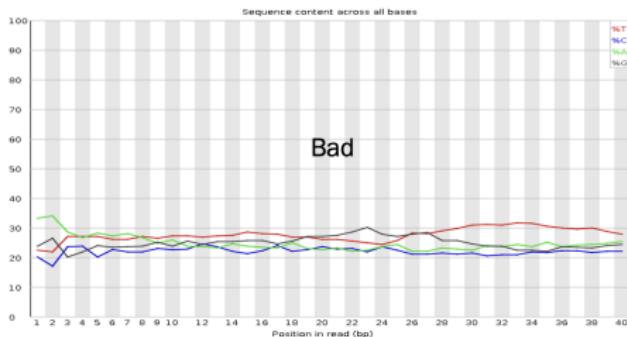
Sequence Base Content



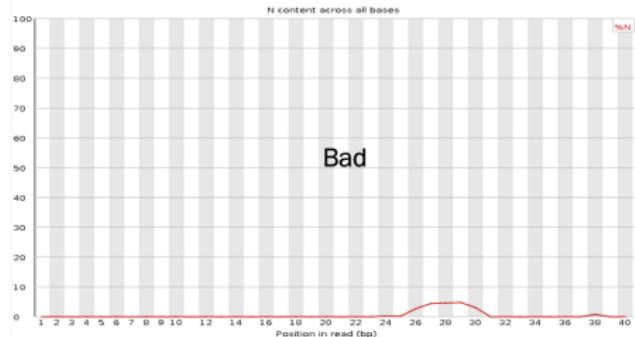
N base content



Sequence content across all bases



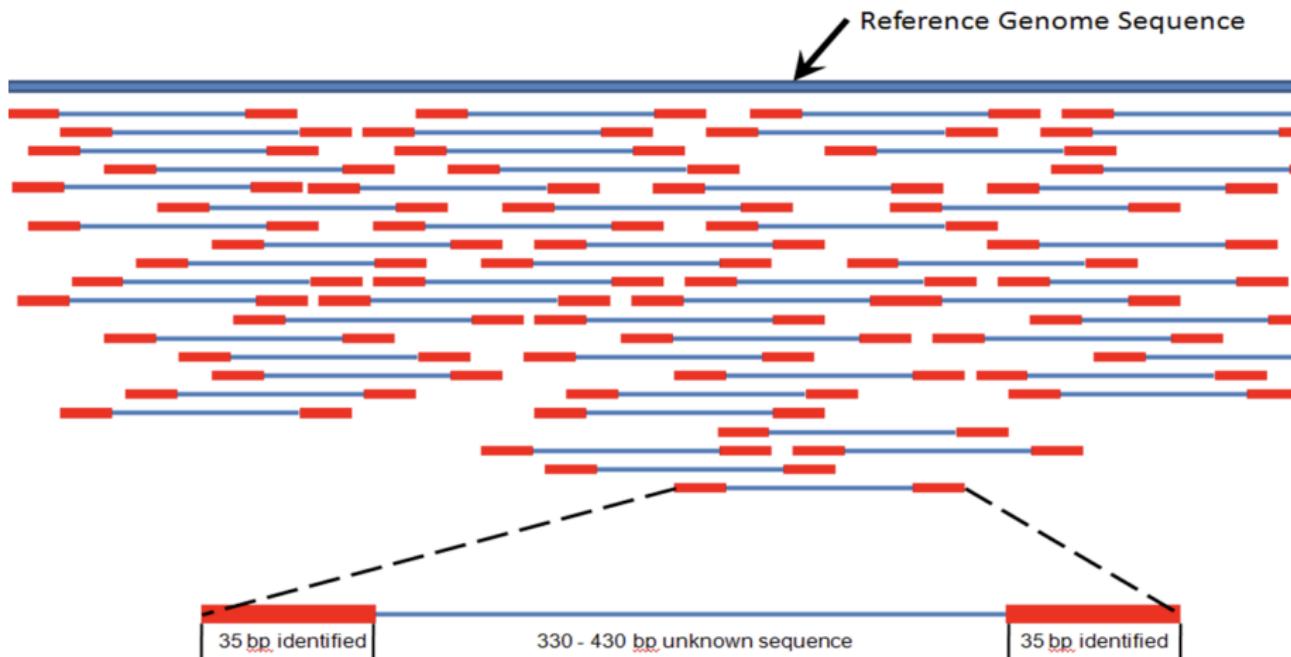
N content across all bases



Alignment to the Reference Genome

Goal: Find the genomic Location of origin for the sequencing read.

Software: Bowtie2, TopHat, STAR, Subread/Rsubread, many others!



Indexing your genome

Abraham Lincoln: “Give me six hours to chop down a tree and I will spend the first four sharpening the axe.” (4 minutes indexing the genome, 2 minutes aligning the reads)

Indexing your genome

Note that you will rarely do this for human alignment. You will usually download an existing index given to you by others who have already done this work. You will do this often if you are aligning microbial reads, e.g. MTB or some other organism for which others have not already made your index for you.

```
buildindex(basename="rna_seq/genome/ucsc.hg19.chr1_120-150M",
           reference="rna_seq/genome/ucsc.hg19.chr1_120-150M.fasta.gz")
```

Took me ~0.2 minutes!

Aligning your reads:

Note that this outputs results in a .bam file and not a .sam file

```
align(index="rna_seq/genome/ucsc.hg19.chr1_120-150M",
      readfile1="rna_seq/reads/R01_10_short500K.fq.gz",
      output_file="rna_seq/alignments/R01_10_short.bam",
      nthreads=4)
```

My laptop is an Apple M2, which has 8 cores (used 4 cores), 24GB RAM:

- Took 15.7 minutes to align ~60M reads to the 30M bases
- Took 0.7 minutes to align ~6.5M reads to the 30M bases
- Took 0.3 minutes to align ~500K reads to the 30M bases

Aligned Sequencing Data Formats (SAM and BAM)

Note that Rsubread outputs a .bam file (bam = binary alignment map) and not a .sam file (sam = sequence alignment map). Here is some information about a .sam file: [https://en.wikipedia.org/wiki/SAM_\(file_format\)](https://en.wikipedia.org/wiki/SAM_(file_format))

Col	Field	Type	Regexp/Range	Brief description
1	QNAME	String	[!-?A-~]{1,255}	Query template NAME
2	FLAG	Int	[0,2 ¹⁶ -1]	bitwise FLAG
3	RNAME	String	* [!-()+-<>-~] [!-~]*	Reference sequence NAME
4	POS	Int	[0,2 ²⁹ -1]	1-based leftmost mapping POSition
5	MAPQ	Int	[0,2 ⁸ -1]	MAPping Quality
6	CIGAR	String	* ([0-9]+[MIDNSHPX=])+	CIGAR string
7	RNEXT	String	* = [!-()+-<>-~] [!-~]*	Ref. name of the mate/next segment
8	PNEXT	Int	[0,2 ²⁹ -1]	Position of the mate/next segment
9	TLEN	Int	[-2 ²⁹ +1,2 ²⁹ -1]	observed Template LENGTH
10	SEQ	String	* [A-Za-z=.]+	segment SEQuence
11	QUAL	String	[!-~]+	ASCII of Phred-scaled base QUALity+33

Aligned Sequencing Data Formats (SAM and BAM)

To convert .sam to .bam or vice versa, a package called Rsamtools. Using Rsamtools, you can convert bam to sam as follows:

```
asSam("rna_seq/alignments/R01_10_short.sam",
      overwrite=T)
```

To convert to bam:

```
#asBam("rna_seq/alignments/R01_10_short.sam")
```

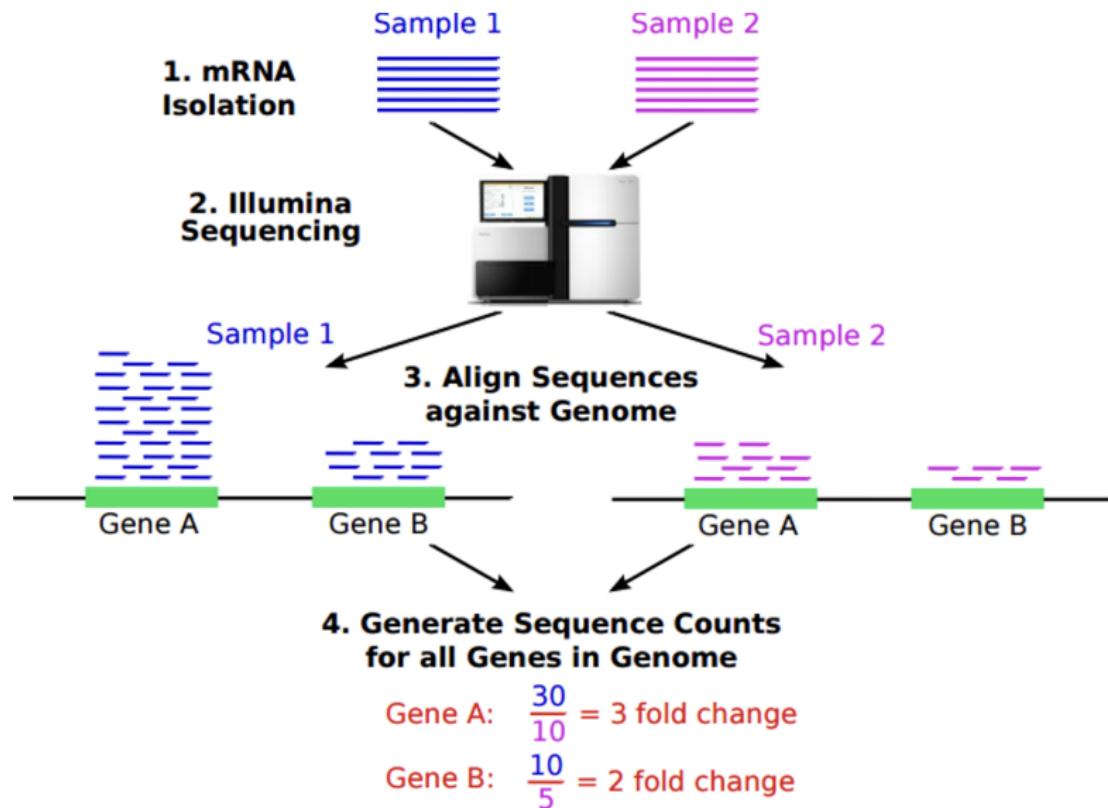
Feature counts

Now we can count reads hitting genes.

Approaches/software:

- HT-Seq
- STAR
- Cufflinks
- RPKM FPKM or CPM
- RSEM
- edgeR
- findOverlaps (GenomicRanges)
- featureCounts (Rsubread)

Feature counts



Feature counts

```
fCountsList = featureCounts(  
  "rna_seq/alignments/R01_10_short.bam",  
  annot.ext="rna_seq_files/genome/genes.chr1_120-150M.gtf",  
  isGTFAnnotationFile=TRUE)  
  
featureCounts = cbind(fCountsList$annotation[,1],  
                      fCountsList$counts)  
  
write.table(featureCounts,  
            "rna_seq/alignments/R01_10_short.features.txt",  
            sep="\t", col.names=FALSE, row.names=FALSE, quote=FALSE)
```

Use the Single Cell Toolkit (SCTK) to analyze your RNA-seq data!

- Inputs: RNA-seq, Nanostring, Proteomic, immunological assay data
- Interactive analyses and visualization of data
- Save results, figures, etc
- Sophisticated data structures
- R/Bioconductor package



Bioconductor
OPEN SOURCE SOFTWARE FOR BIOINFORMATICS

Search:

[Home](#) [Install](#) [Help](#) [Developers](#) [About](#)

[Home](#) » [Bioconductor 3.8](#) » [Software Packages](#) » [singleCellITK](#)

singleCellITK

platforms all rank 1102 / 1649 posts 0 in Bioc 1 year
build ok updated since release

DOI: [10.18129/B9.bioc.singleCellITK](https://doi.org/10.18129/B9.bioc.singleCellITK)

Interactive Analysis of Single Cell RNA-Seq Data

Documentation »

Bioconductor

- Package [vignettes](#) and manuals.
- [Workflows](#) for learning and use.
- [Course and conference](#) material.
- [Videos](#).
- Community [resources](#) and [tutorials](#).

R / [CRAN](#) packages and [documentation](#)

The screenshot shows the main landing page of the SCTK v1.3.7 web application. At the top, there is a navigation bar with links for "Upload", "Data Summary & Filtering", "Visualization & Clustering", "Batch Correction", "Differential Expression", "Enrichment Analysis", and "Sample Size". Below the navigation bar, the title "Single Cell Toolkit" is displayed in large, bold letters, followed by the subtitle "Filter, cluster, and analyze single cell RNA-Seq data". A link "Need help? Read the docs." is also present. The background is light gray.

Upload

(help)

Choose data source:

- Upload files
- Upload SCTKExperiment RDS File
- Use example data

Upload data in tab separated text format:

Example count file:

Gene	Cell1	Cell2	...	CellN
Gene1	0	0	...	0
Gene2	5	6	...	0
Gene3	4	3	...	8
...
CountM	10	10	...	0

Example sample annotation file:

Cell	Annot1	...
Cell1	a	..
Cell2	a	..
Cell3	b	..
...

Example feature file:

Gene	Annot2	...
Gene1	a	..
Gene2	a	..
Gene3	b	..
...

Installing and using the SCTK

```
install.packages("devtools")
devtools::install_github("wewanjohnson/singleCellTK")
library(singleCellTK)
singleCellTK()

### Example: open downstream_analysis/
### features_combined.txt and meta_data.txt
```

Batch effects

Batch Effect: Non-biological variation due to differences in batches of data that confound the relationships between covariates of interest.

Batch effects are caused by differences in:

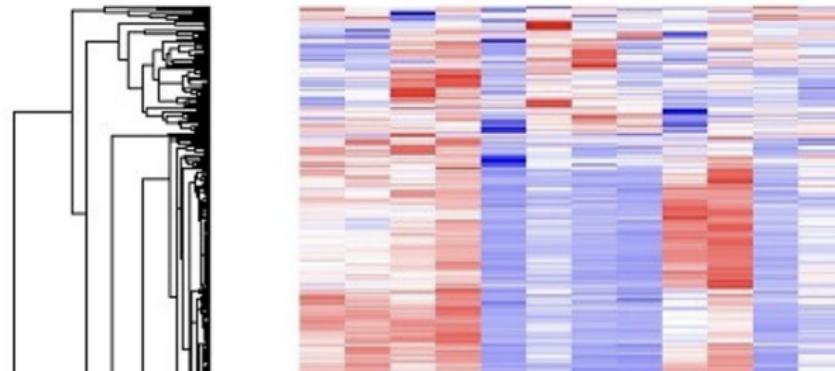
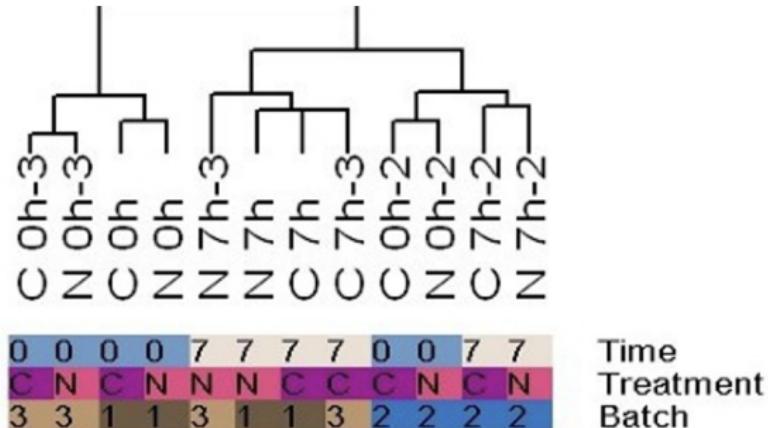
- Gene expression profiling platform
- Lab protocol or experimenter
- Time of day or processing
- Atmospheric ozone level (Rhodes et al. 2004)

Batch Effect Example #1: Nirtic Oxide

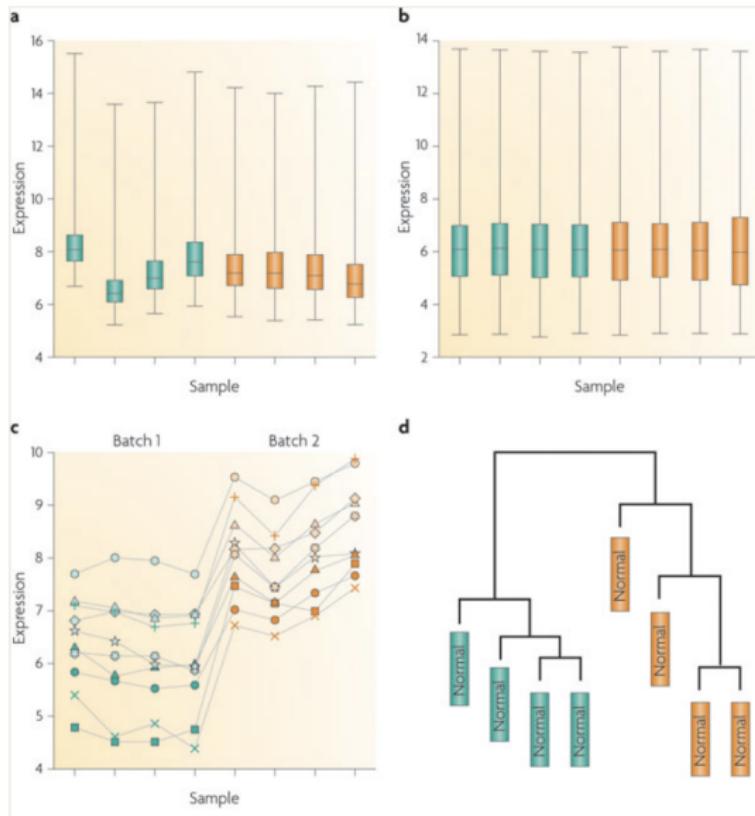
Example 1 resulted from an oligonucleotide microarray (Affymetrix HG-U133A) experiment on human lung fibroblast cells (IMR90) designed to reveal whether exposing mammalian cells to nitric oxide (NO) stabilizes mRNAs. Control samples and samples exposed to NO for 1 h were then transcription inhibited for 7.5 h.

Microarray data were collected at baseline (0 h, just before transcription inhibition) and at the end of the experiment (after 7.5 h) for both the control and the NO-treated group. It was hypothesized that NO will induce or inhibit the expression of some genes, but would also stabilize the mRNA of many genes, preventing them from being degraded after 7.5 h.

Batch Effect Example #1: Nirtic Oxide



Batch Effect Example #2: Control Gene Expression



ComBat Batch Adjustment

Consider the following model:

$$Y_{ijg} = \alpha_g + X\beta_g + \gamma_{ig} + \delta_{ig}\epsilon_{ijg}$$

where:

- α_g is the overall gene expression
- X is a design matrix
- β_g contains the regression coefficients
- The error terms $\epsilon_{ijg} \sim N(0, \sigma_g^2)$
- γ_{ig} and δ_{ig} are additive and multiplicative batch effects

ComBat Batch Adjustment

Adjust for batch effects:

$$Y_{ijg}^* = \frac{Y_{ig} - \hat{\alpha}_g - X\hat{\beta}_g - \hat{\gamma}_{ig}}{\hat{\delta}_{ig}} + \hat{\alpha}_g + X\hat{\beta}_g$$

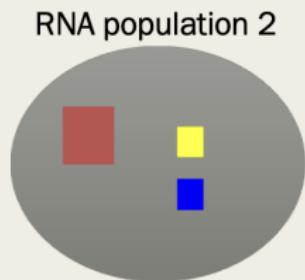
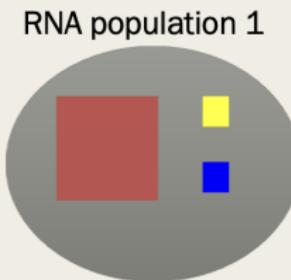
Answer: Empirical Bayes!

Normalization

Need to normalize data because of:

- Sequencing depth difference in each RNA sample
- RNA composition differences
- Highly expressed genes can consume a substantial proportion of RNA-Seq reads, causing other genes to be under-sampled
- Different methods
 - Log counts
 - Counts per million (CPM and logCPM; RPKM, FPKM)
 - Trimmed mean of M-values (edgeR/limma)
 - Median of Ratios method (DESeq)

Normalization



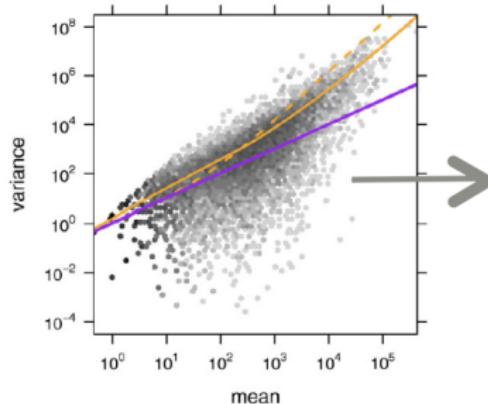
Assuming equal sequencing depth, yellow and blue will get lower RPKM in RNA population 1 although they have equal expression levels

Problem of overdispersion:

Alignment and feature counting result in discrete count data (i.e. the number of reads to each gene). A first thought might be to use a Poisson distribution to model the counts. However, the Poisson makes a strict mean-variance assumption (i.e. they are the same. Studies have demonstrated that a negative binomial fits data better.

Problem of overdispersion:

Many studies have shown that the variance grows faster than the mean in RNAseq data. This is known as **overdispersion**.



Software use:
Negative binomial distribution
Non-parametric methods
Voom transformation

- Mean count vs variance of RNA seq data. Orange line: the fitted observed curve. Purple: the variance implied by the Poisson distribution.

Data Structures

A data structure is a particular way of organizing data in a computer so that it can be used effectively. The idea is to reduce the space and time complexities of different tasks.

Data Structures

Data structures in R programming are tools for holding multiple values, variables, and sometimes functions

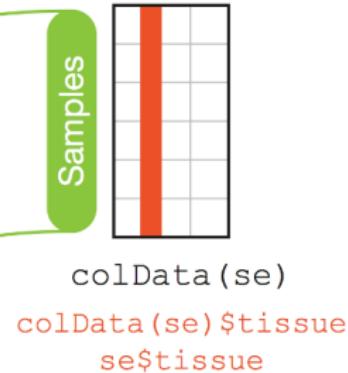
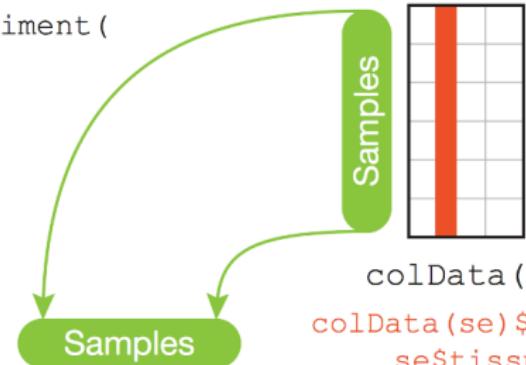
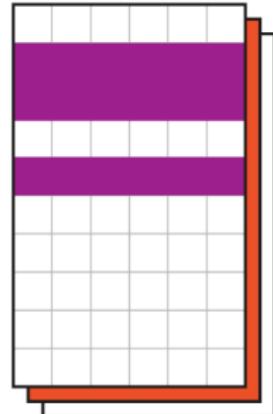
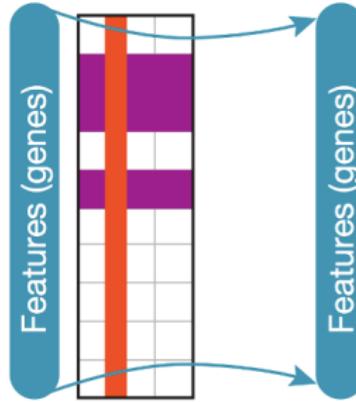
Please think very carefully about the way you manage and store your data! This can make your life much easier and make your code and data cleaner and more portable!

Data Structures

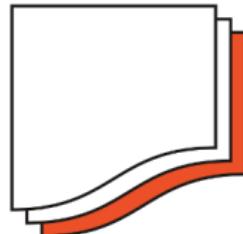
There are advanced R data structures, **S3** and **S4** class objects, that can facilitate object orientated programming. One useful example of an S4 class data structure is the **SummarizedExperiment** object.

Data Structures

```
se <- SummarizedExperiment(  
  assays,  
  rowData,  
  colData,  
  exptData  
)
```



se %in% CNVs



Visualization and Dimension reduction

Using an example dataset from: Verma, et al., 2018

```
## read in data
counts <- read.table(
  "rna_seq/downstream_analysis/features_combined.txt",
  sep="\t", header=T, row.names=1)
meta_data <- read.table(
  "rna_seq/downstream_analysis/meta_data.txt",
  sep="\t", header=T, row.names=1)
group <- meta_data$Disease
```

Visualization and Dimension reduction

```
## Make SummarizedExperiment
sce_hivtb <- SummarizedExperiment(assays=list(counts=counts),
                                     colData = meta_data)

## Make log counts, counts per million (cpm), logcpm
sce_hivtb <- mkAssay(sce_hivtb, log = TRUE,
                      counts_to_CPM = TRUE)
assays(sce_hivtb)

## List of length 4
## names(4): counts log_counts counts_cpm log_counts_cpm
```

Principal Components Analysis (PCA)

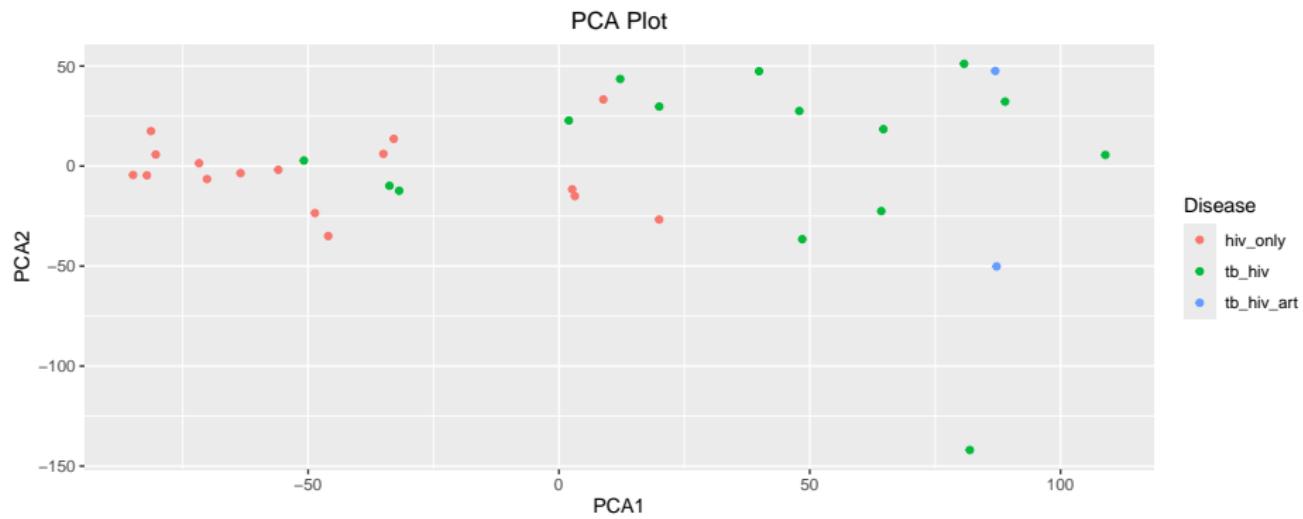
```
set.seed(1)
pca_out <- prcomp(t(assay(sce_hivtb, "log_counts_cpm")))

pca_plot <- as.data.frame(pca_out$x)
pca_plot$Disease <- as.factor(sce_hivtb$Disease)

g <- pca_plot %>% ggplot(aes(x=PC1, y=PC2, color=Disease)) +
  geom_point(size=1.5) + xlab("PCA1") + ylab("PCA2") +
  theme(plot.title = element_text(hjust = 0.5)) +
  ggtitle("PCA Plot")

plot(g)
```

Principal Components Analysis (PCA)



Uniform Manifold Approximation and Projection (UMAP)

For more on UMAP, please visit the following excellent tutorial: <https://pair-code.github.io/understanding-umap/>

Uniform Manifold Approximation and Projection (UMAP)

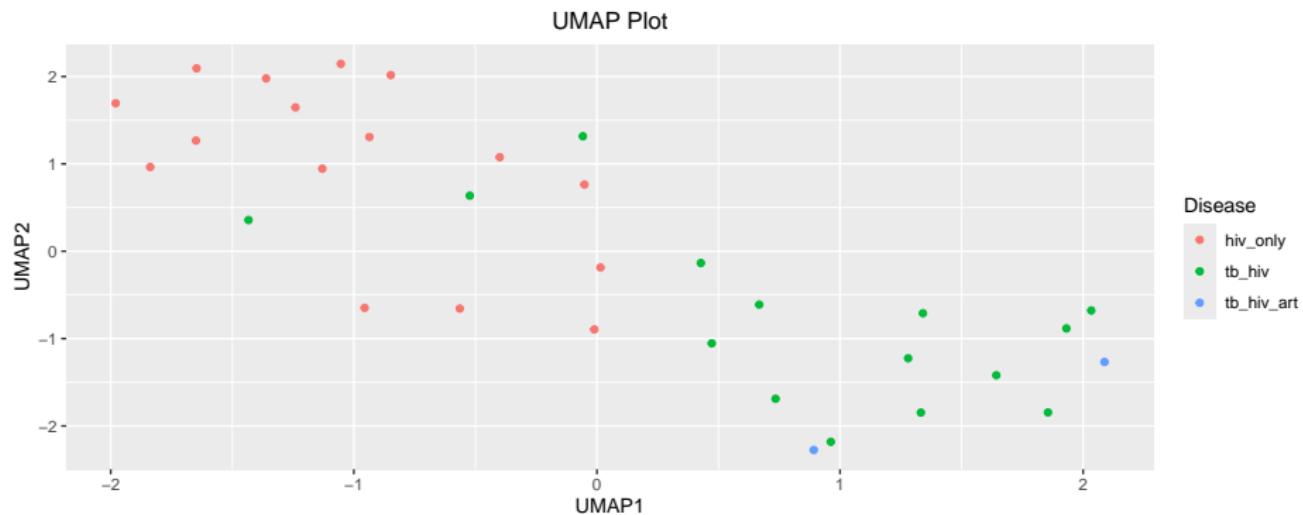
```
set.seed(1)
umap_out <- umap(t(assay(sce_hivtb,"log_counts_cpm")))

umap_plot <- as.data.frame(umap_out$layout)
umap_plot$Disease <- as.factor(sce_hivtb$Disease)

g <- umap_plot %>% ggplot(aes(x=V1, y=V2, color=Disease)) +
  geom_point(size=1.5) + xlab("UMAP1") + ylab("UMAP2") +
  theme(plot.title = element_text(hjust = 0.5)) +
  ggtitle("UMAP Plot")

plot(g)
```

Uniform Manifold Approximation and Projection (UMAP)



Differential Expression

Table I: Software packages for detecting differential expression

Method	Version	Reference	Normalization ^a	Read count distribution assumption	Differential expression test
edgeR	3.0.8	[4]	TMM/ <u>Upper quartile/RLE (DESeq-like)/None</u> (all scaling factors are set to be one)	Negative binomial distribution	Exact test
DESeq	1.10.1	[5]	DESeq sizeFactors	Negative binomial distribution	Exact test
baySeq	1.12.0	[6]	Scaling factors (<u>quantile</u> /TMM/total)	Negative binomial distribution	Assesses the posterior probabilities of models for differentially and non-differentially expressed genes via empirical Bayesian methods and then compares these posterior likelihoods
NOIseq	1.1.4	[7]	<u>RPKM</u> /TMM/ <u>Upper quartile</u>	Nonparametric method	Contrasts fold changes and absolute differences within a condition to determine the null distribution and then compares the observed differences to this null
SAMseq (samr)	2.0	[8]	SAMseq specialized method based on the mean read count over the null features of the data set	Nonparametric method	Wilcoxon rank statistic and a resampling strategy
Limma	3.14.4	[9]	TMM	voom transformation of counts	Empirical Bayes method
Cuffdiff 2 (Cufflinks)	2.0.2-beta	[10]	Geometric (DESeq-like)/quartile/classic-fpkm	Beta negative binomial distribution	t-test
EBSeq	1.1.7	[11]	DESeq median normalization	Negative binomial distribution	Evaluates the posterior probability of differentially and non-differentially expressed entities (genes or isoforms) via empirical Bayesian methods

^aIn case of availability of several normalization methods, the default one is underlined.

EdgeR Example

Implements statistical methods for DE analysis based on the negative binomial model:

```
#Gene Filtering  
counts<-counts [which(rowSums(cpm(counts))>1),]  
#Computes library size  
dge <- DGEList(counts=counts, group=group)  
#TMM normalization  
dge <- calcNormFactors(dge)  
# Design matrix  
design<-model.matrix(~group)  
#Estimates common, trended and tagwise dispersion  
dge<-estimateDisp(counts,design)
```

EdgeR Example

In negative binomial models, each gene is given a dispersion parameter. Dispersions control the variances of the gene counts and underestimation will lead to false discovery and overestimation may lead to a lower rate of true discovery

EdgeR Example

```
## Performs likelihood ratio tests
#fits a negative binomial GLM with the dispersion estimates
fit<-glmFit(counts,design, dispersion=dge$tagwise.dispersion)
# Performs likelihood ratio test
# Compares the goodness of the fit, full versus reduced model
lrt<-glmLRT(fit, coef=2)
# Prints the top results
topTags(lrt)
```

```
## Coefficient: grouptb_hiv
##          logFC    logCPM        LR      PValue
## IL1R2     4.334196  8.207931 100.01344 1.513665e-23 2.933634e
## AP3B2     5.758193  2.952770   71.58586 2.654527e-17 2.572370e
## FCGR1C    2.818498  4.536149   65.07230 7.220003e-16 4.664362e
## VNN1      3.150158  8.071776   64.39089 1.020287e-15 4.943545e
## CYP1B1    3.135471  6.873000   63.00698 2.059751e-15 7.984006e
## TL1BP1    2.706121  6.107474   60.75863 6.451068e-15 2.084003e
```

EdgeR Example

```
# Perform quasi-likelihood F-tests
## Replace the chisquare approximation to the likelihood
## ratio statistic with a quasi-likelihood F-test,
## more control of error rate
fit<-glmQLFit(counts, design,
                 dispersion=dge$tagwise.dispersion)
## use for small dataset, reflects uncertainty in estimating
## control when the number of replicates is small dispersion
## for each gene, more robust and reliable error rate
qlf<-glmQLFTest(fit, coef=2)
topTags(qlf)
```

```
## Coefficient: grouptb_hiv
##          logFC      logCPM          F      PValue
## IL1R2    4.334269  8.207931 104.44957 1.620179e-24 3.140069e
## AP3B2    5.765316  2.952770  74.47435 6.158177e-18 5.967581e
## VNN1    -2.150104  2.071776  67.10428 2.554151e-16 1.507305e
```

EdgeR Example

```
#For visualization, heatmaps/PCA  
Logcpm<-cpm(counts,log=TRUE)
```