

BIO634: Next generation Sequencing 2

Transcriptome and Biological Interpretation

8th-9th May, 2023

Deepak Kumar Tanwar

URPP Evolution in Action
Embedded bioinformatician



@d_k_tanwar



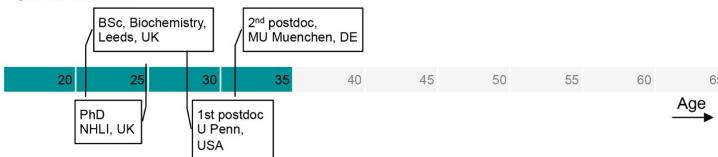
deepak.tanwar@evolution.uzh.ch

Leon (bioinformatics user)

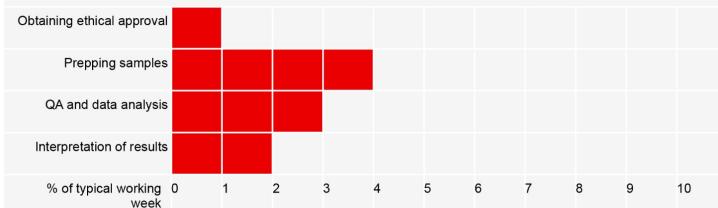
Leon is on his second postdoctoral fellowship, working on quorum sensing in bacteria. "I'm using a combination of transcriptomics, proteomics and metabolomics to understand these pathogenic changes better" he explains. "I end up with big spreadsheets of protein or gene IDs and I'm trying to piece together which signalling pathways are involved in flipping to the pathogenic state". He has been on an introductory Unix course but is much more comfortable with GUIs than with the command line. "I just have a visual brain", he says.



Career timeline



Typical activities



Distribution of time between bench-work and computational work



Preference for using GUI vs command line



Drivers

- Understanding what makes a usually harmless bacterium pathogenic in the lungs of people with cystic fibrosis

Goals

- QA of -omics data
- Statistical analysis of data
- Data integration and pathway analysis

Pain points

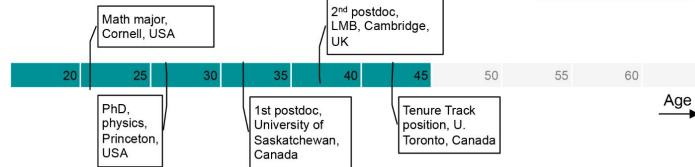
- Lack of access to departmental compute farm
- Sporadic to non-existent access to bioinformatics support

Martha (bioinformatics scientist)

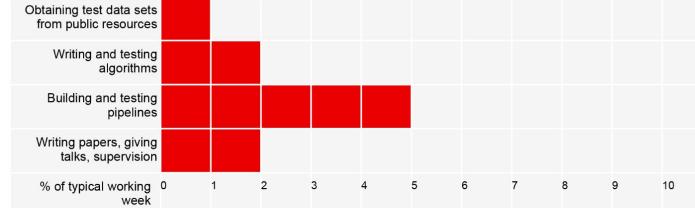
Martha is a senior bioinformatician in an international structural genomics consortium. Her biggest project is on predicting the functions of proteins whose structures have just been solved; she's building a structure-to-function prediction pipeline for the project. This is funded partly by the NIH and partly through industrial funding. She also has a fascination for predicting structure and usually has a student or two working on structural prediction projects.



Career timeline



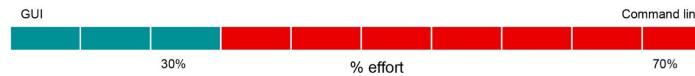
Typical activities



Distribution of time between bench work and computational work



Preference using for GUI vs command line



Drivers

- Understanding the relationship between sequence, structure and function
- Application to target discovery and validation

Goals

- Create a structure-to-function pipeline for molecular biologists
- Predict structures de novo from models of similar, solved structures

Pain points

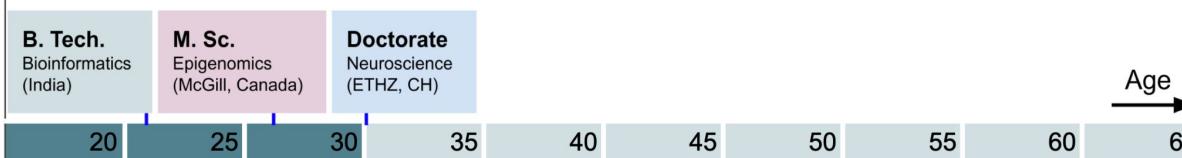
- Sometimes the guys in the lab expect her to fix their computers for them
- Finding students and more senior staff with adequate math

Deepak Tanwar (embedded bioinformatics scientist)

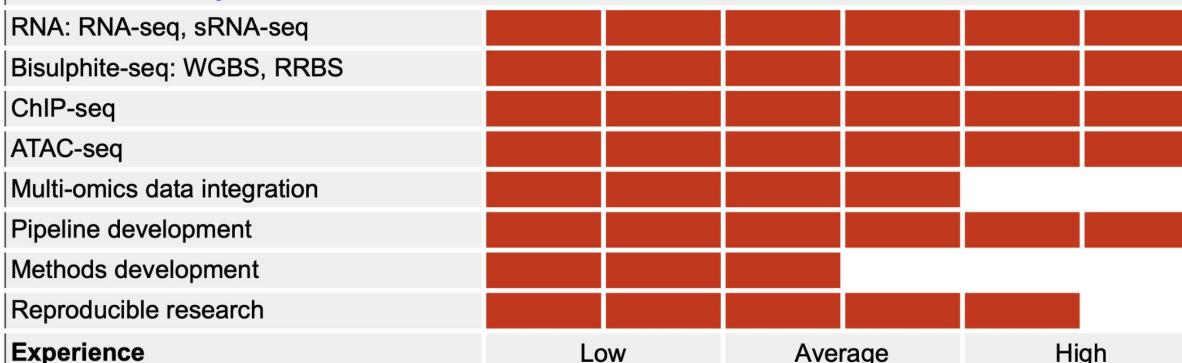
Deepak is a bioinformatician at the URPP Evolution in Action. Deepak has a strong background in multi-omics research and has expertise in reproducible data analysis, benchmarking, and methods development.



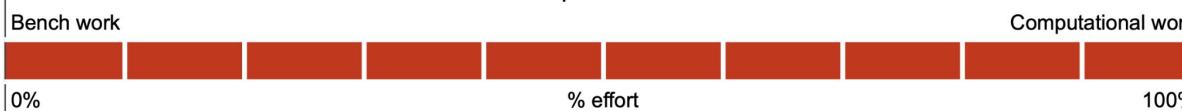
Education timeline



Multi-omics experience



Distribution of time between bench work and computational work



```
cd ~/Desktop
```

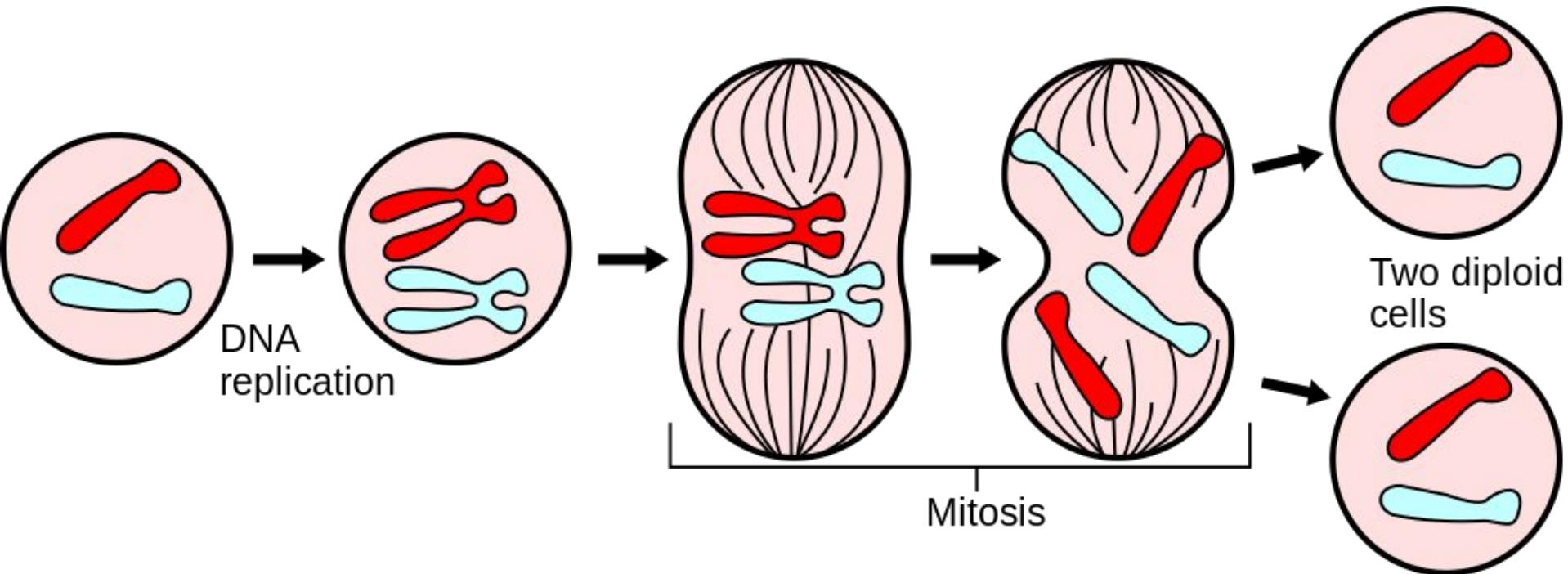
```
git clone https://github.com/urppeia/bio634.git
```

```
cd bio634
```

```
./login.sh      |      ./win_login.sh
```

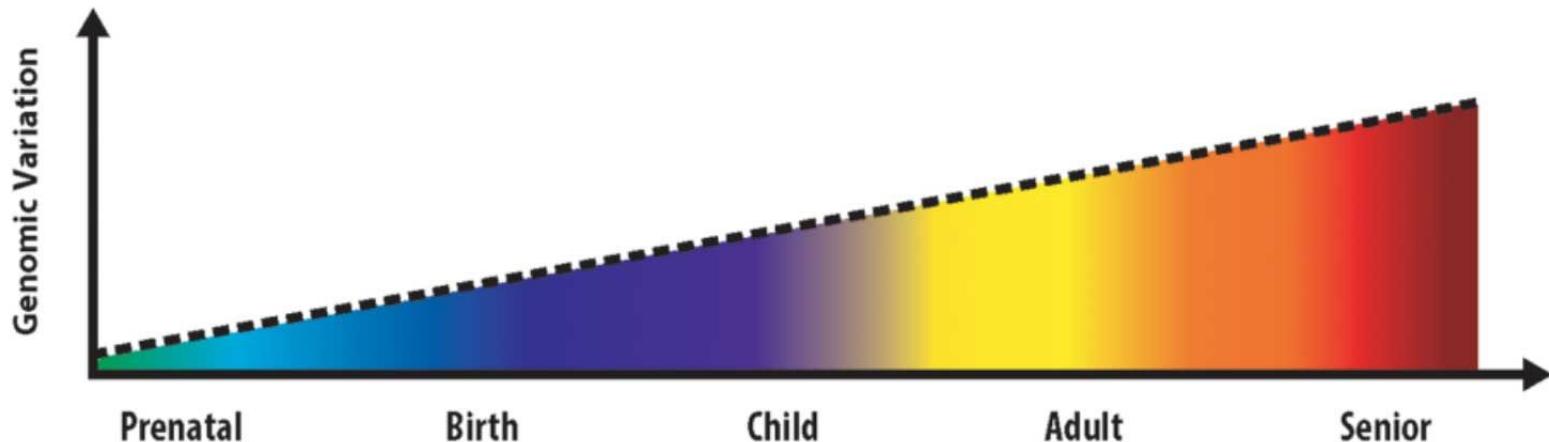
Introduction

All cells have same DNA

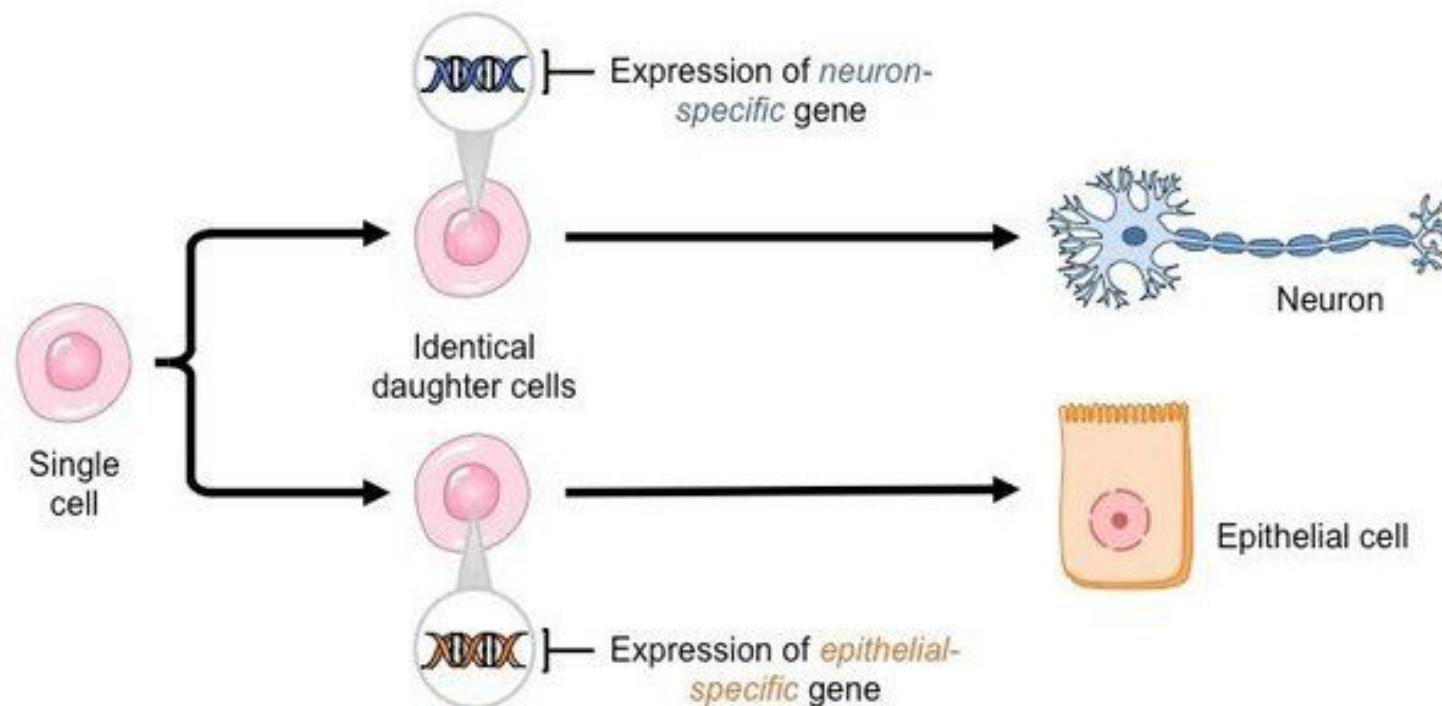


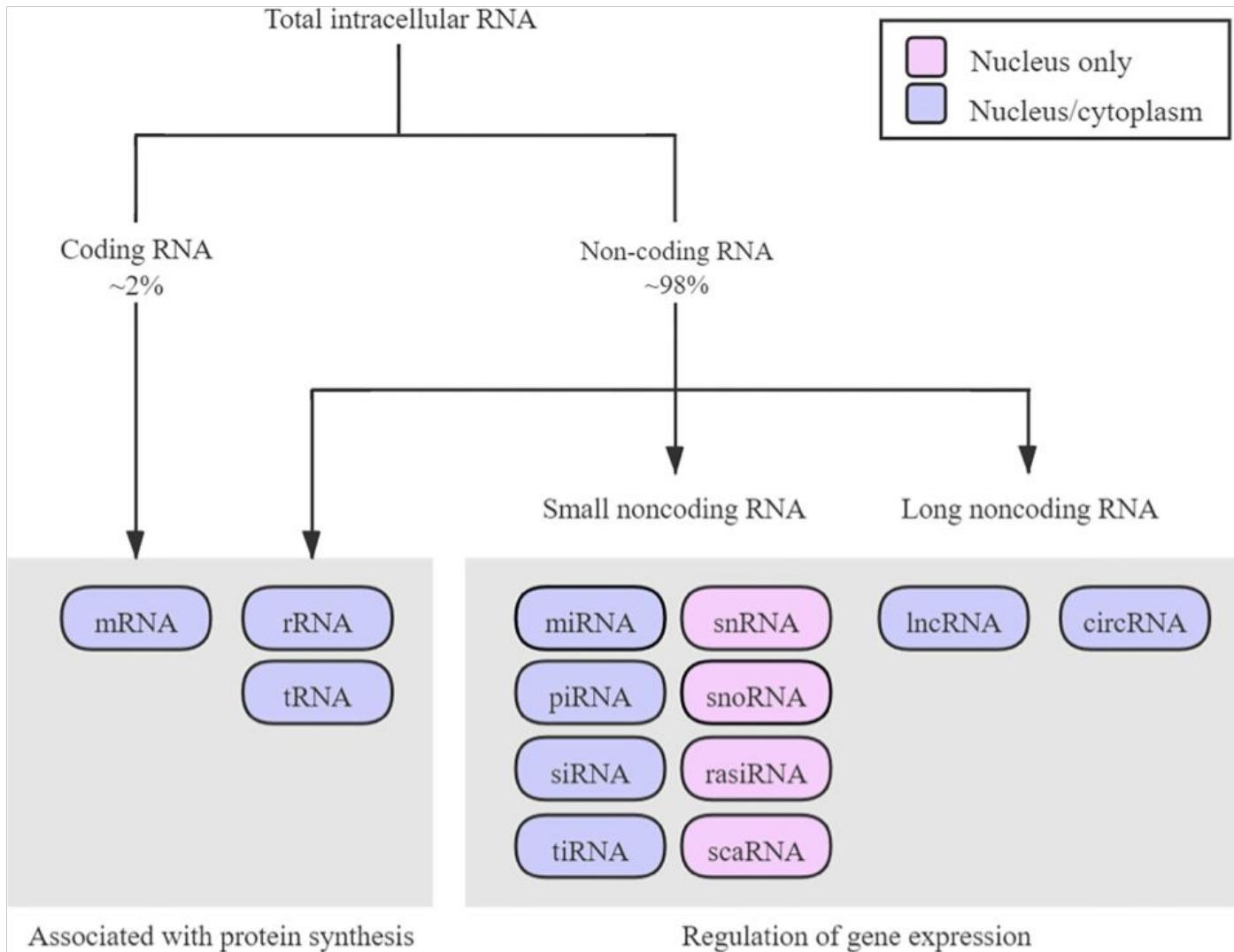
Not all our cells have the same DNA

Genomic Mosaicism Throughout Lifespan

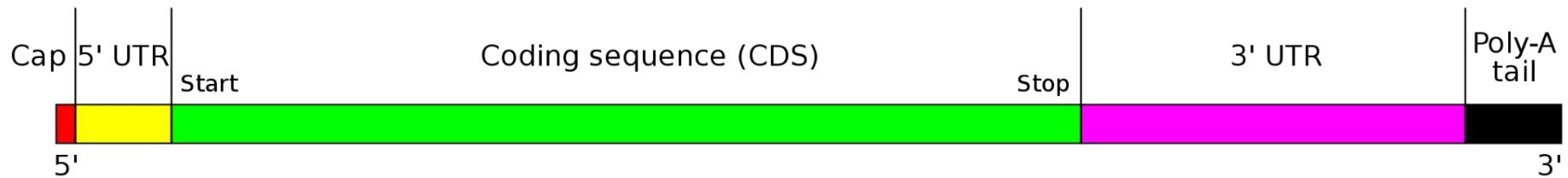


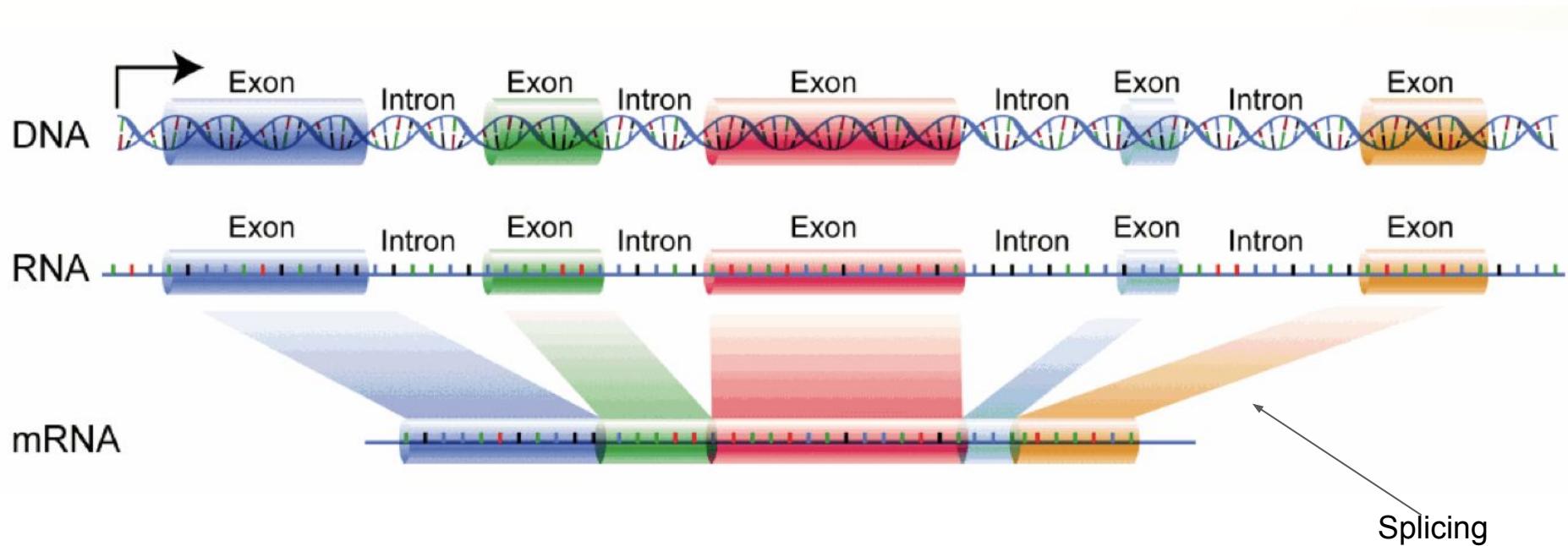
Specific RNA expression make specific cells

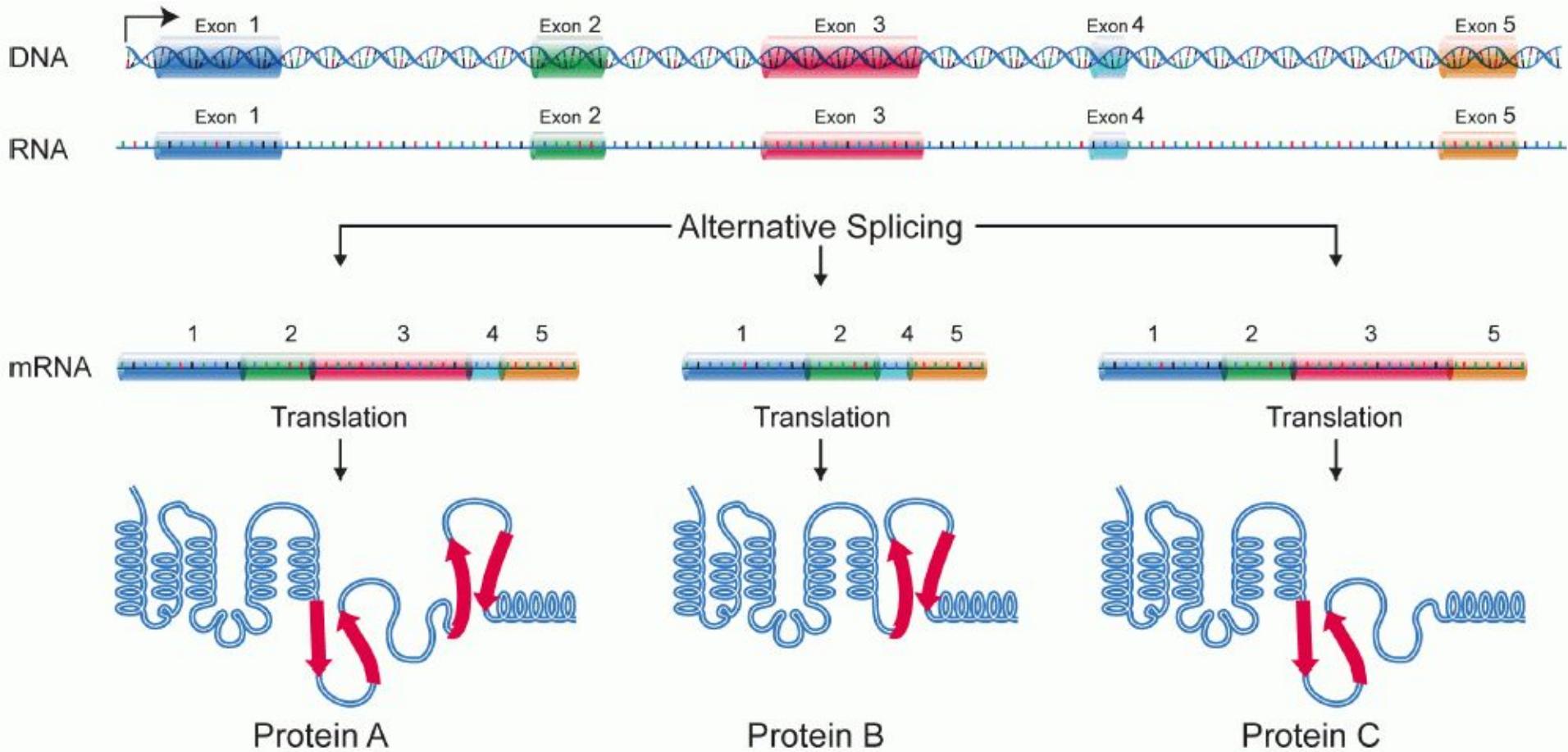




The structure of a typical human protein coding mRNA including the untranslated regions (UTRs)





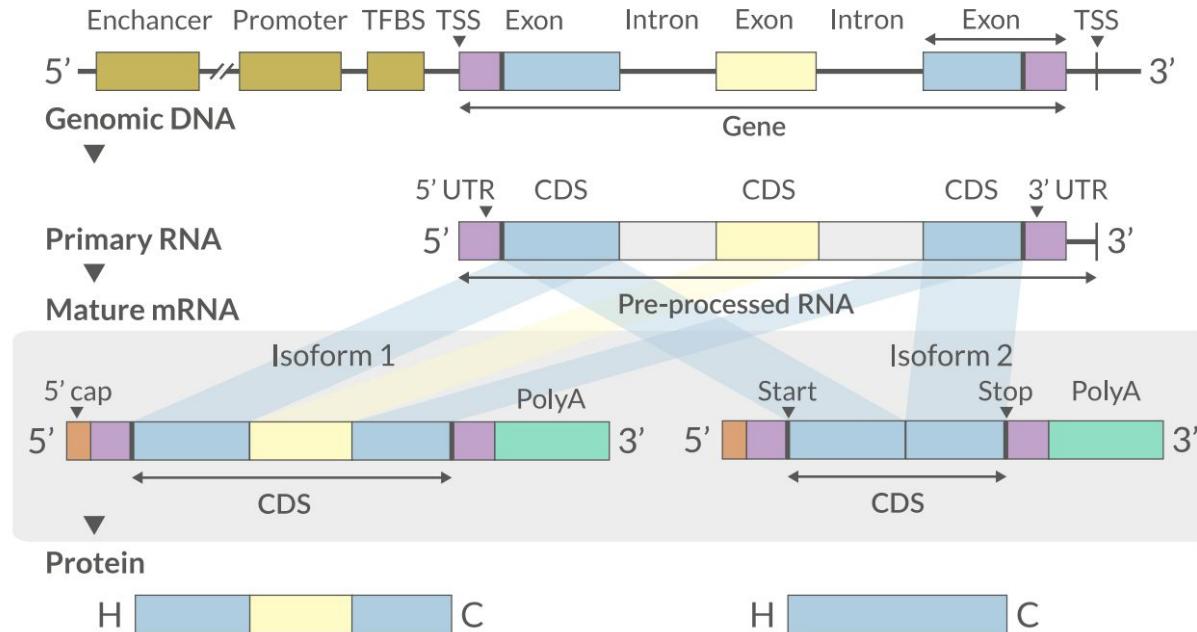


Perspective

Not all exons are protein coding: Addressing a common misconception

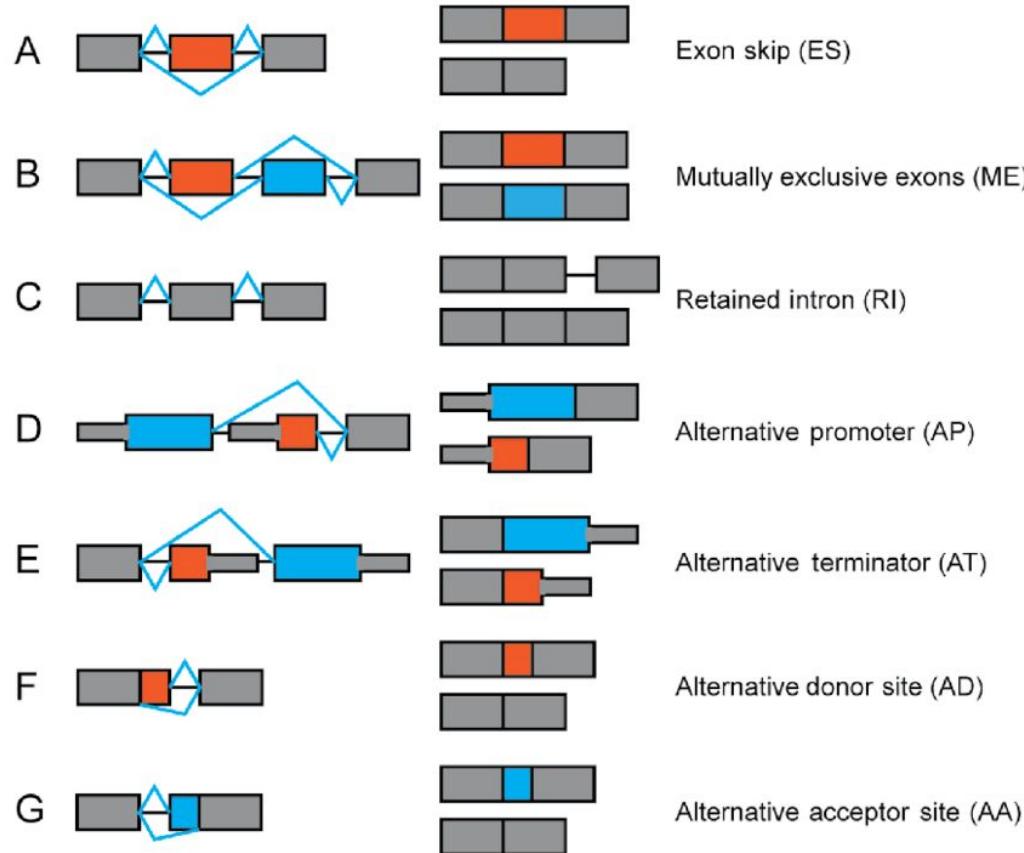
Julie L. Aspden,^{1,2,3} Edward W.J. Wallace,⁴ and Nicola Whiffin^{5,6,*}

Summary so far



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

One gene, many different mRNAs



Why should we sequence RNA?

Designing the right experiment

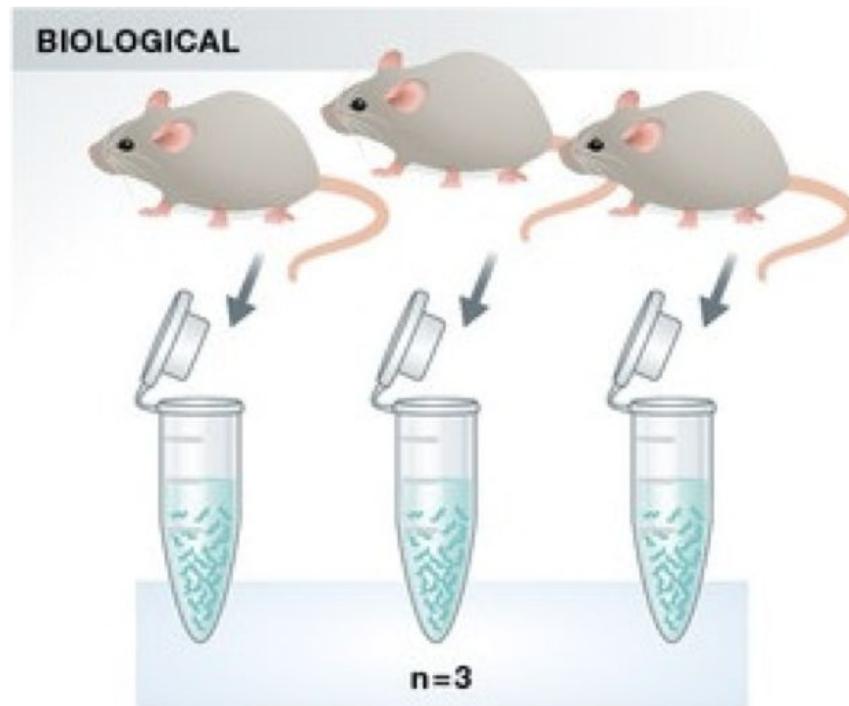
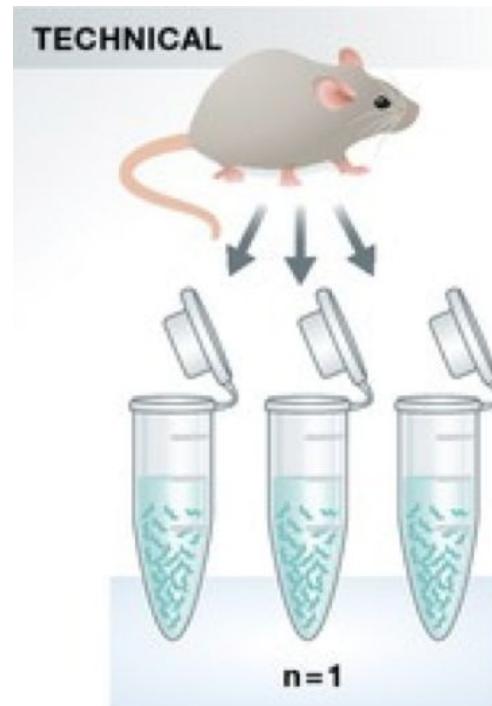
Clear objectives

Amenable to statistical analysis

Reproducible

Experiment design

Importance of replicates



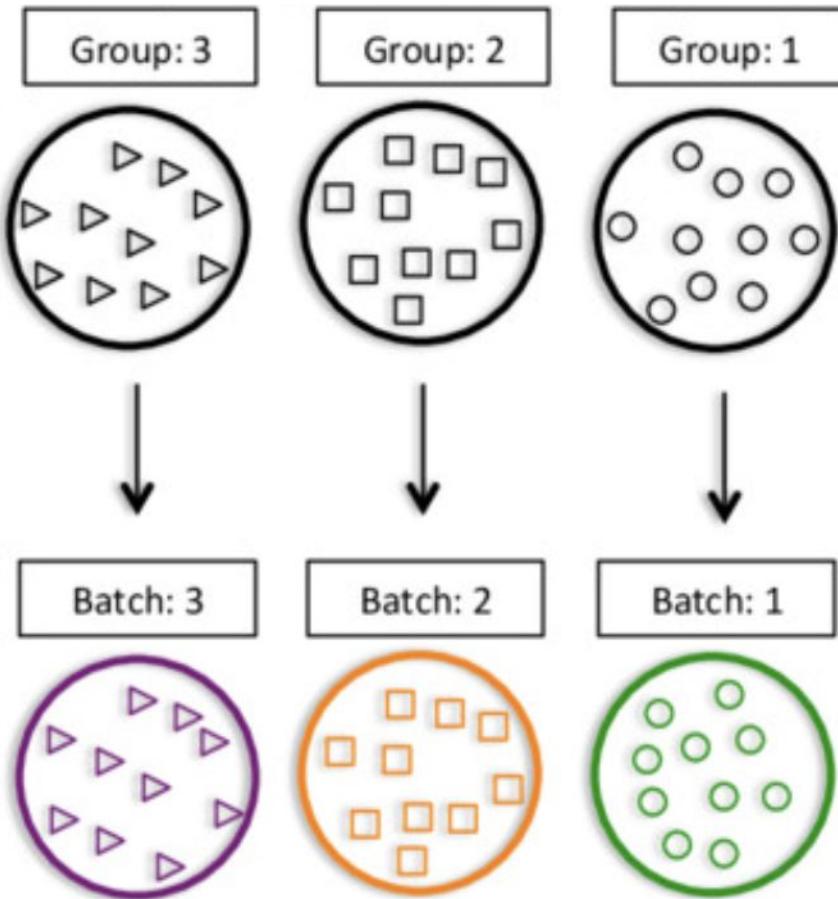
How much should one sequence?

The coverage is defined as:

$$\frac{\text{Read Length} \times \text{Number of Reads}}{\text{Length of Target Sequence}}$$

The amount of sequencing needed for a given sample is determined by the goals of the experiment and the nature of the RNA sample.

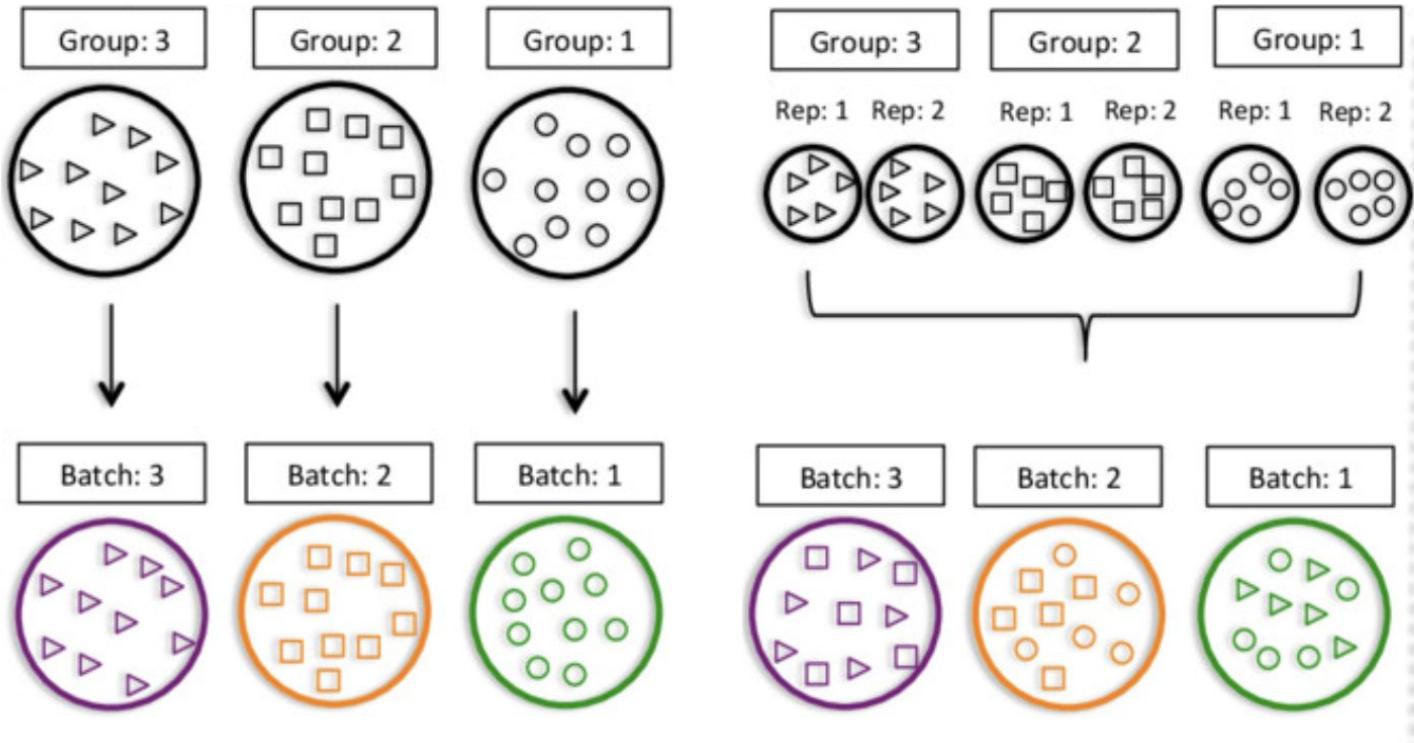
- For a general view of differential expression: 5–25 million reads per sample
- For alternative splicing and lowly expressed genes: 30–60 million reads per sample.
- In-depth view of the transcriptome/assemble new transcripts: 100–200 million reads
- Targeted RNA expression requires fewer reads.
- miRNA-Seq or Small RNA Analysis require even fewer reads.



Biological Group Processing Batch

What changes would you make here to make the experimental design more optimal?

Batch effects?



Write down everything.

High-throughput sequencing

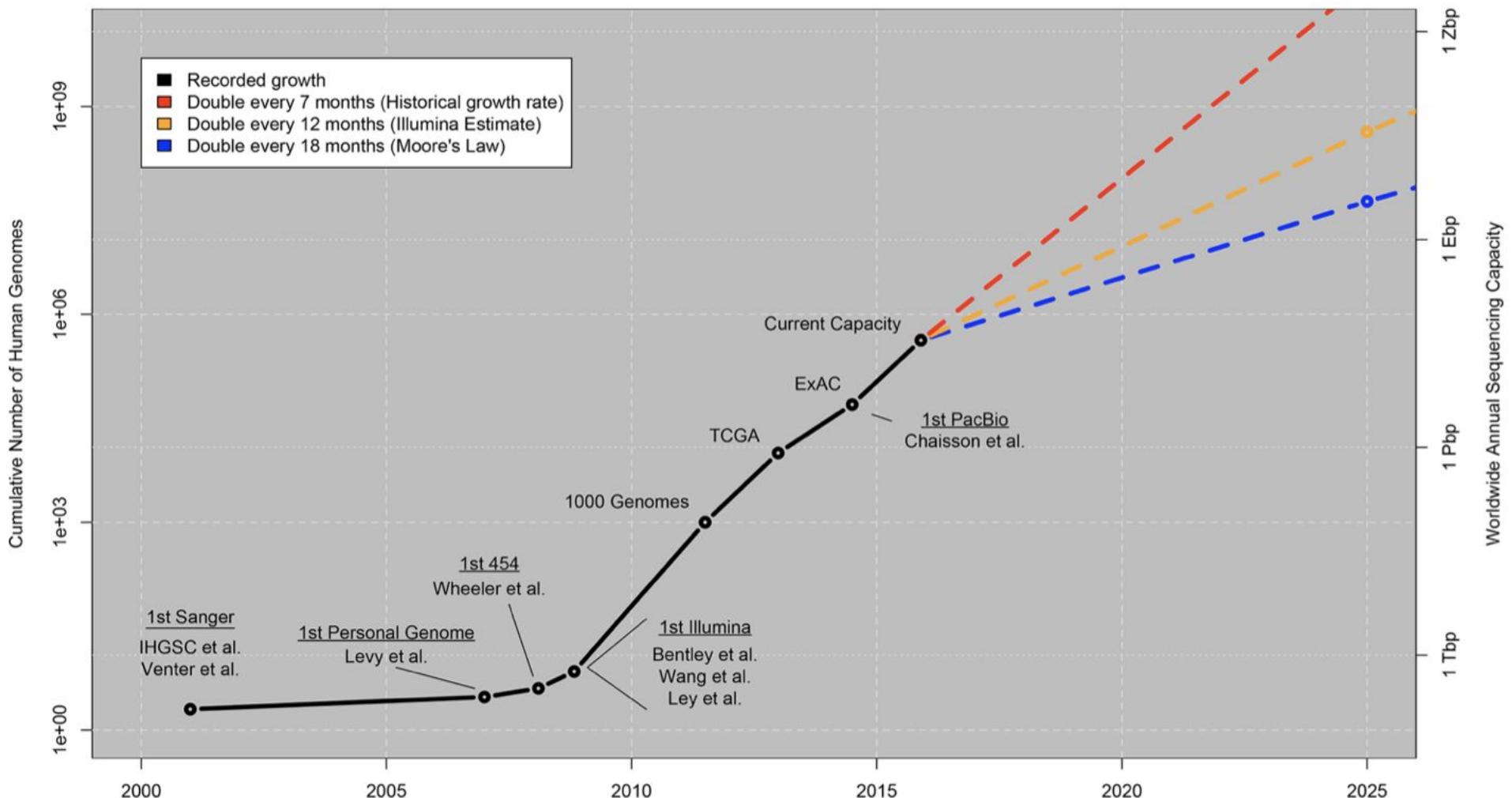
The Importance of Experimental Design



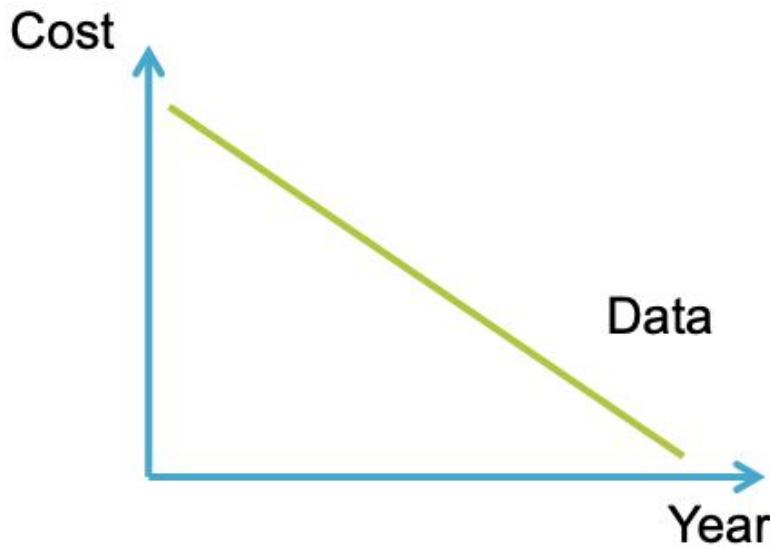
Let's see if the subject
responds to magnetic
stimuli... ADMINISTER
THE MAGNET!

Interesting...there seems
to be a significant
decrease in heart rate.
The fish must sense the
magnetic field.

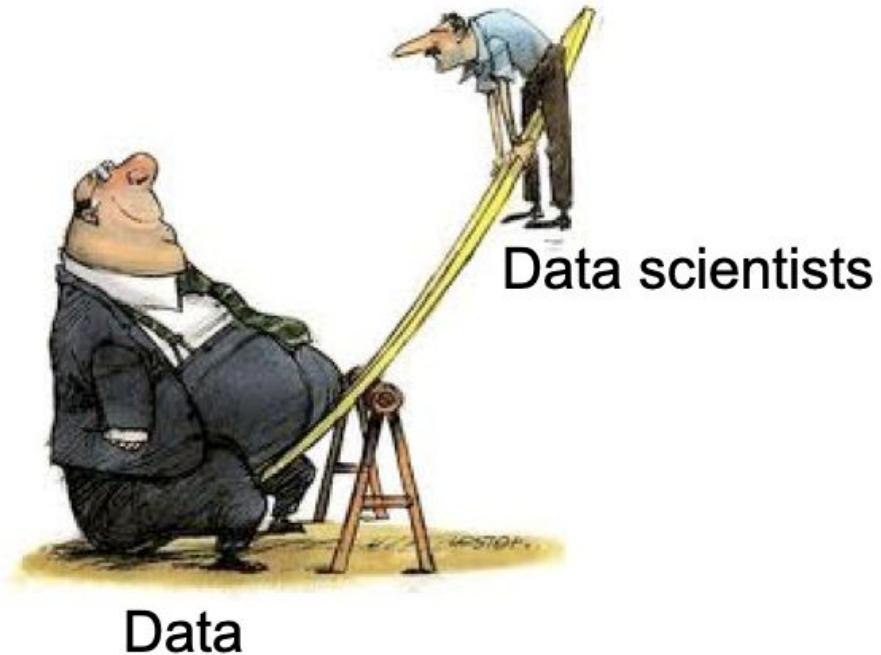
Growth of DNA Sequencing

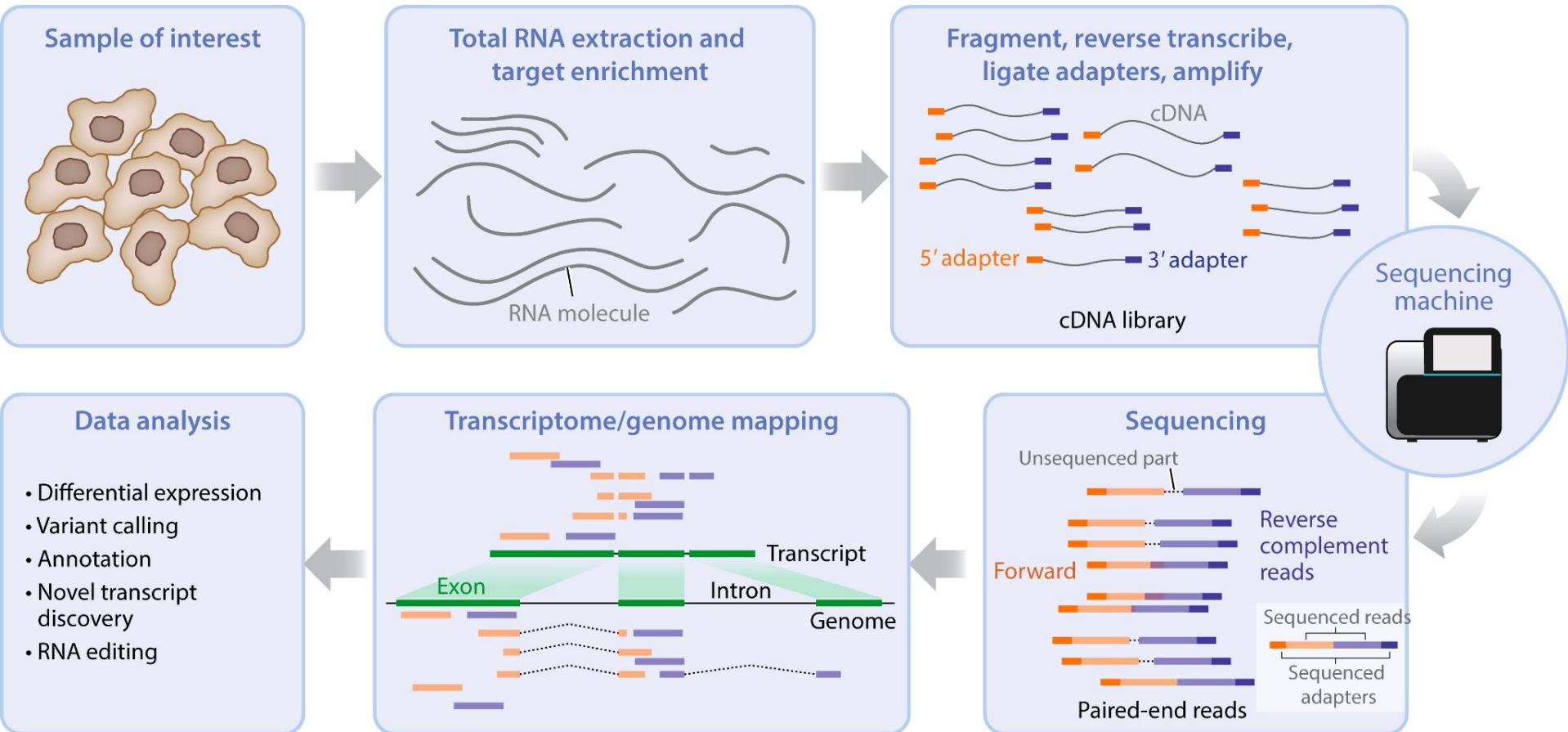


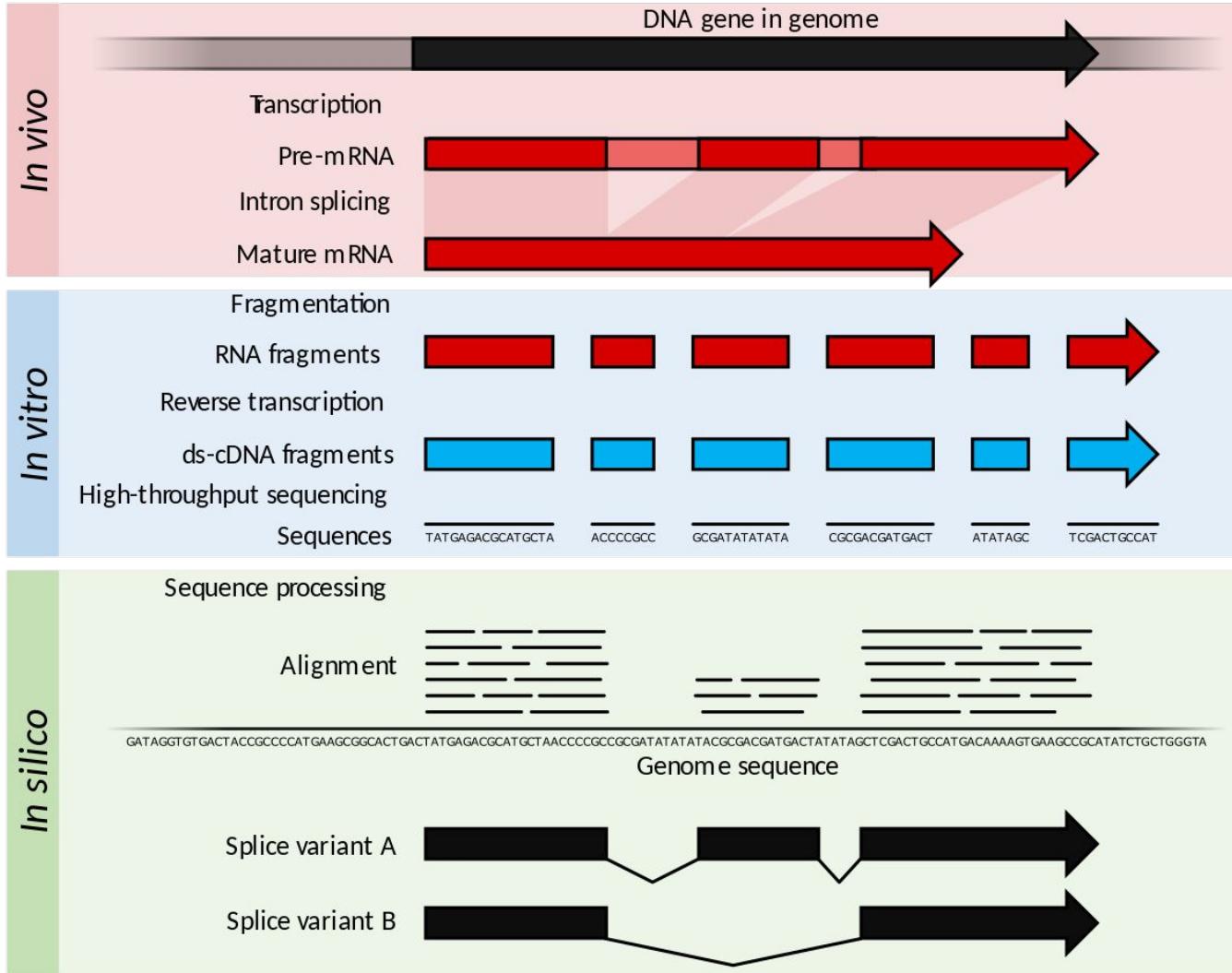
Growing burden on Bioinformaticians



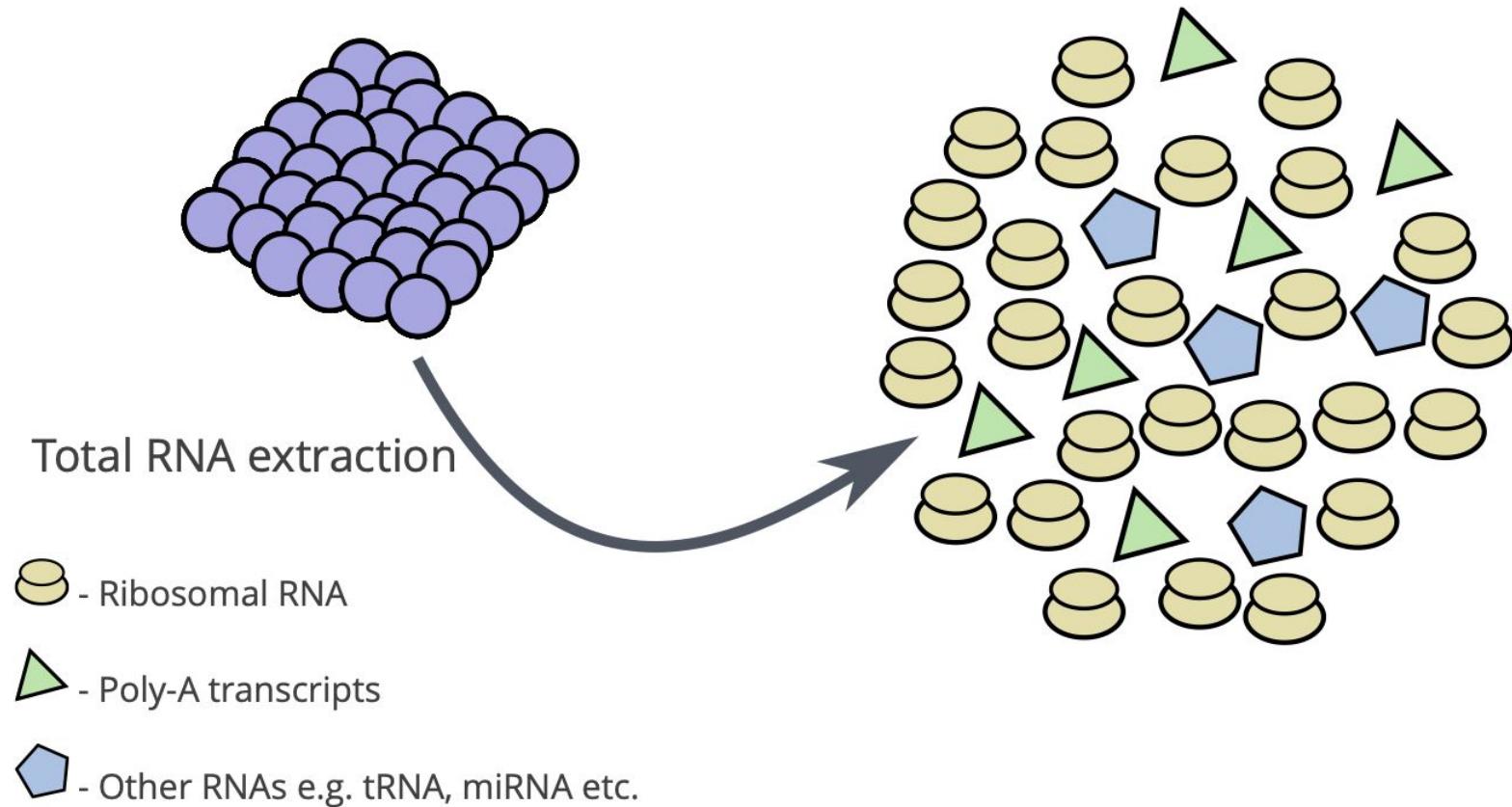
“Per base”

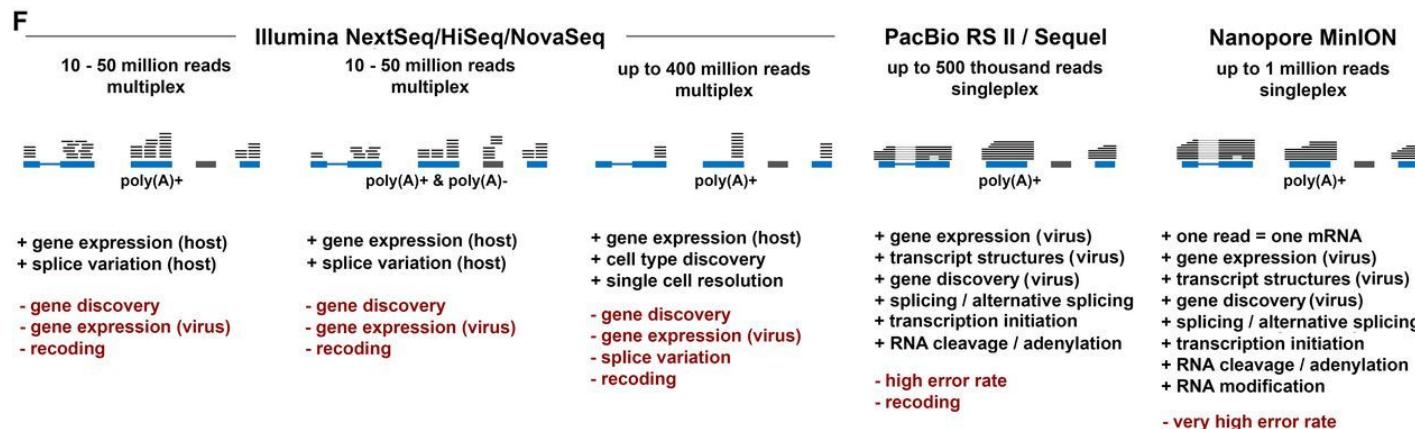
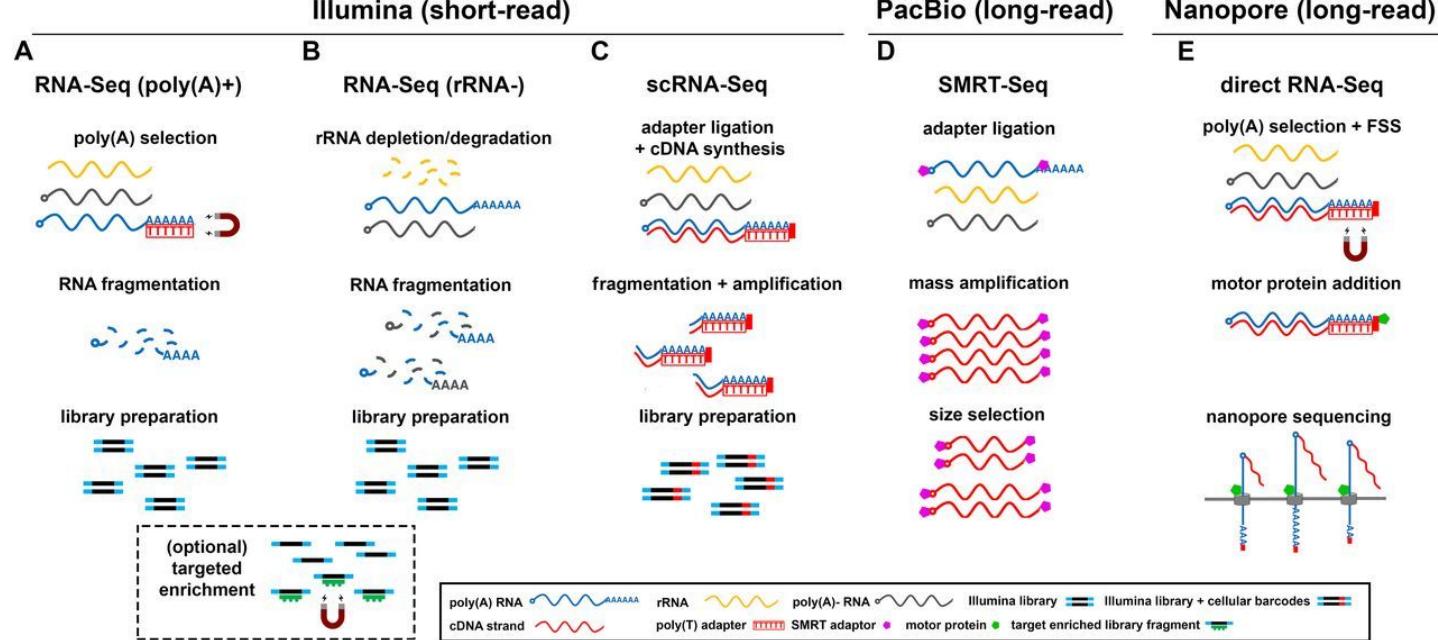






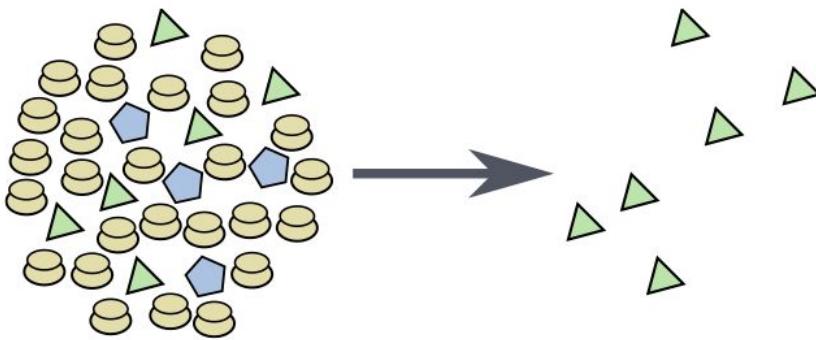
Library preparation



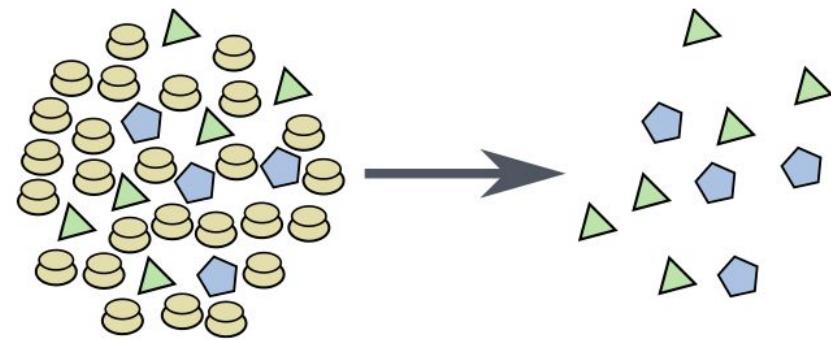


Library preparation

Poly-A Selection



Ribominus selection

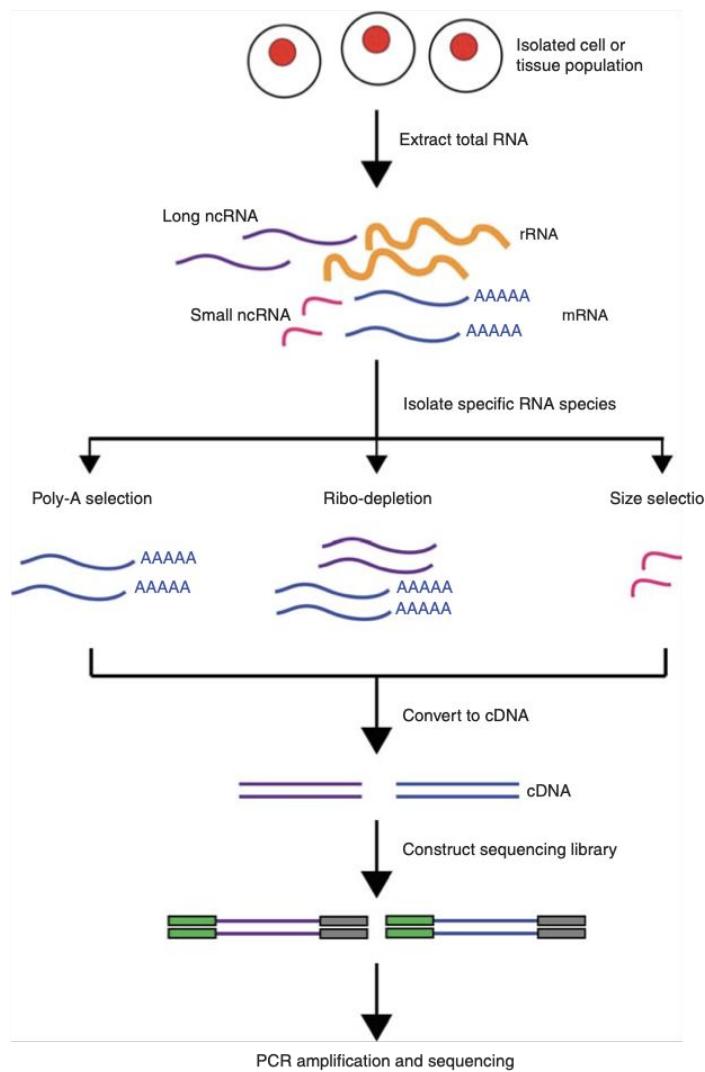


Poly-A transcripts e.g.:

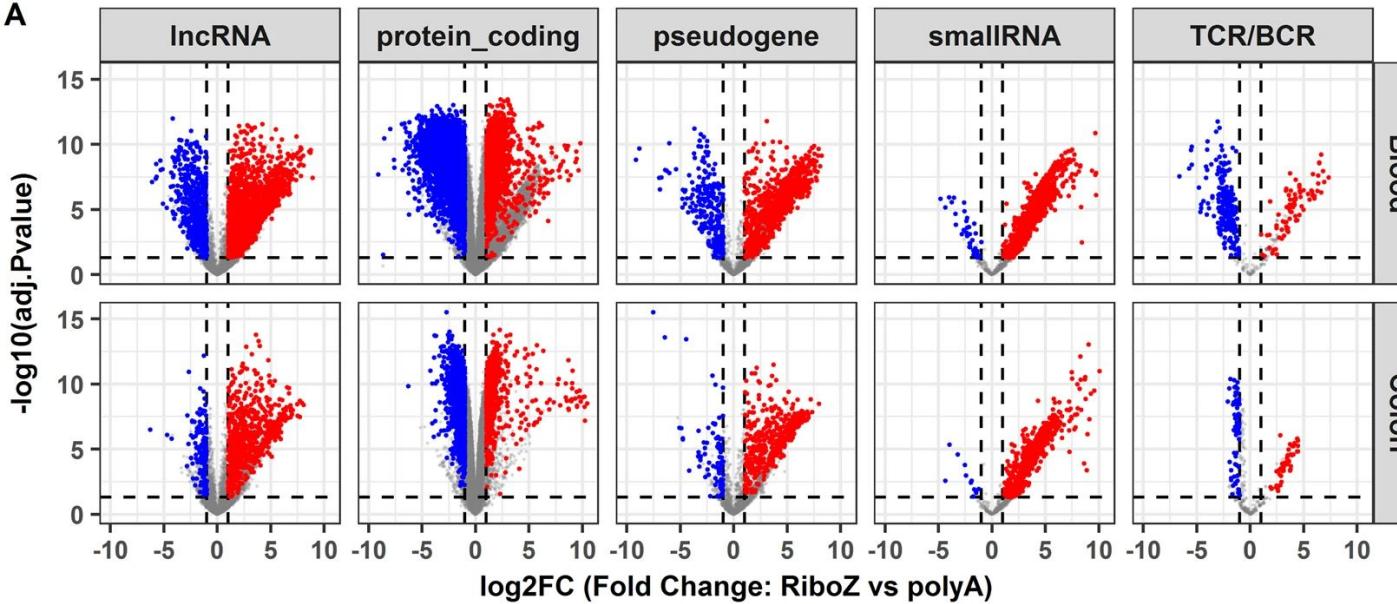
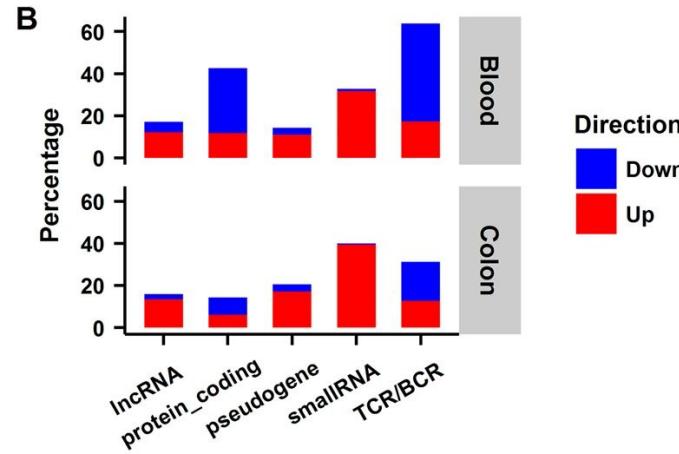
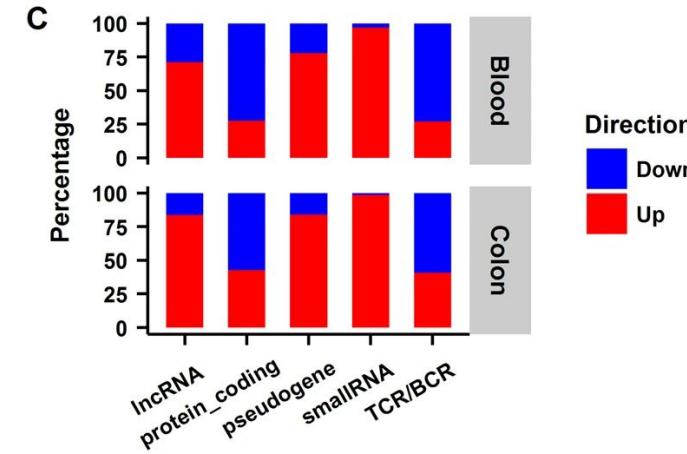
- mRNAs
- immature miRNAs
- snoRNA

Poly-A transcripts + Other mRNAs e.g.:

- tRNAs
- mature miRNAs
- piRNAs



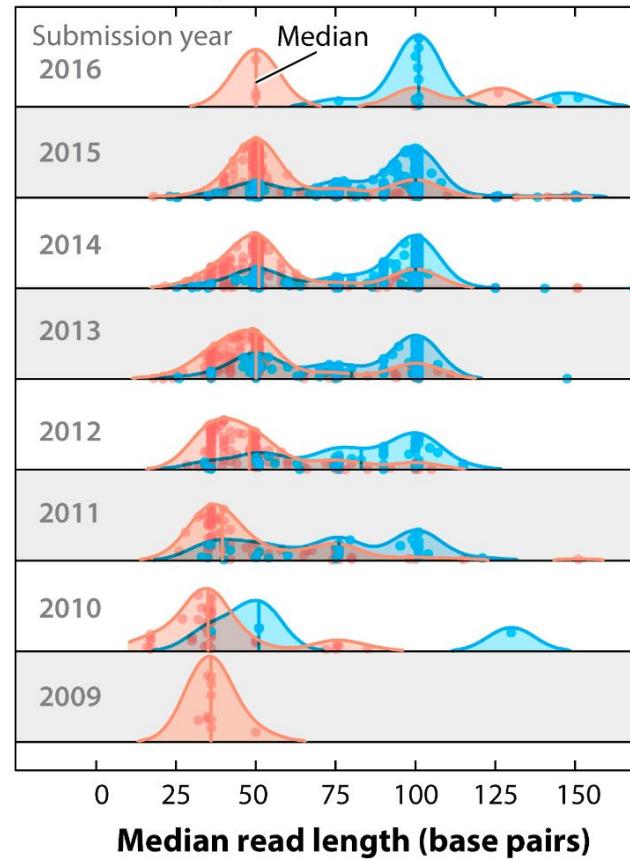
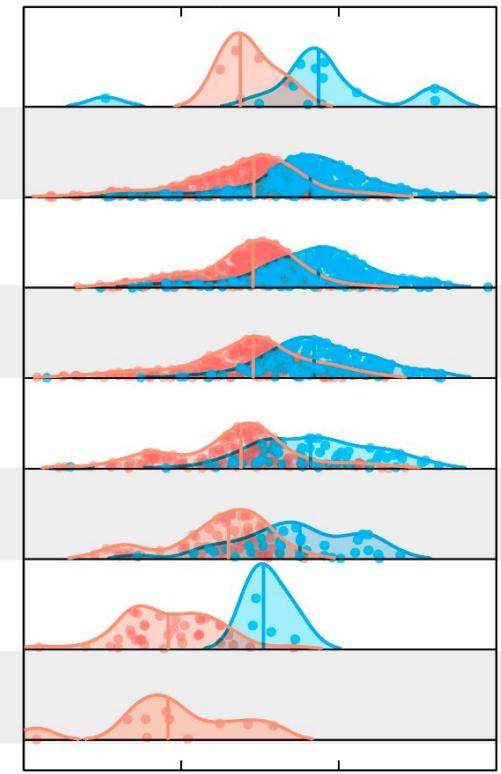
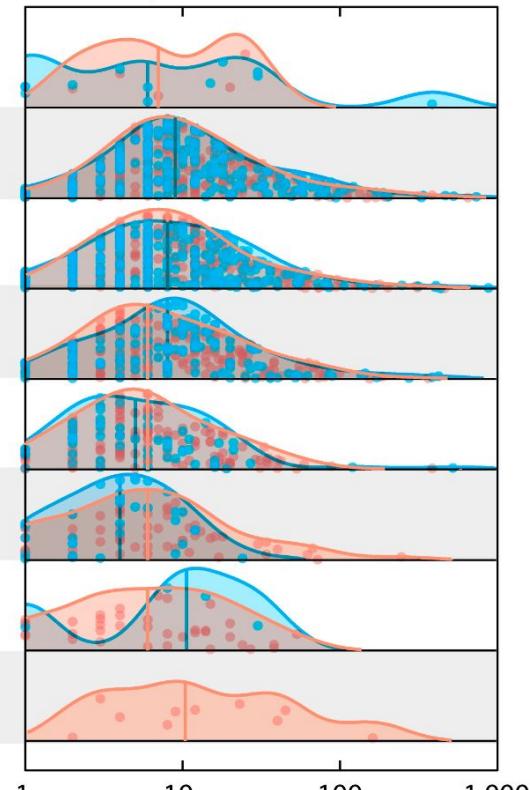
When would you do a size selection?

A**B****C**

"polyA+ selection and rRNA depletion both selectively omit a distinct set of RNAs, so different fractions of the transcriptome are sequenced; thus, generating incompatible datasets."

a Read length

Cumulative Number of Human Genomes

**b** Read depth**c** Sample size

Median read length (base pairs)

Single-end project ($n = 787$)

Median number of reads

Paired-end project ($n = 1,008$)

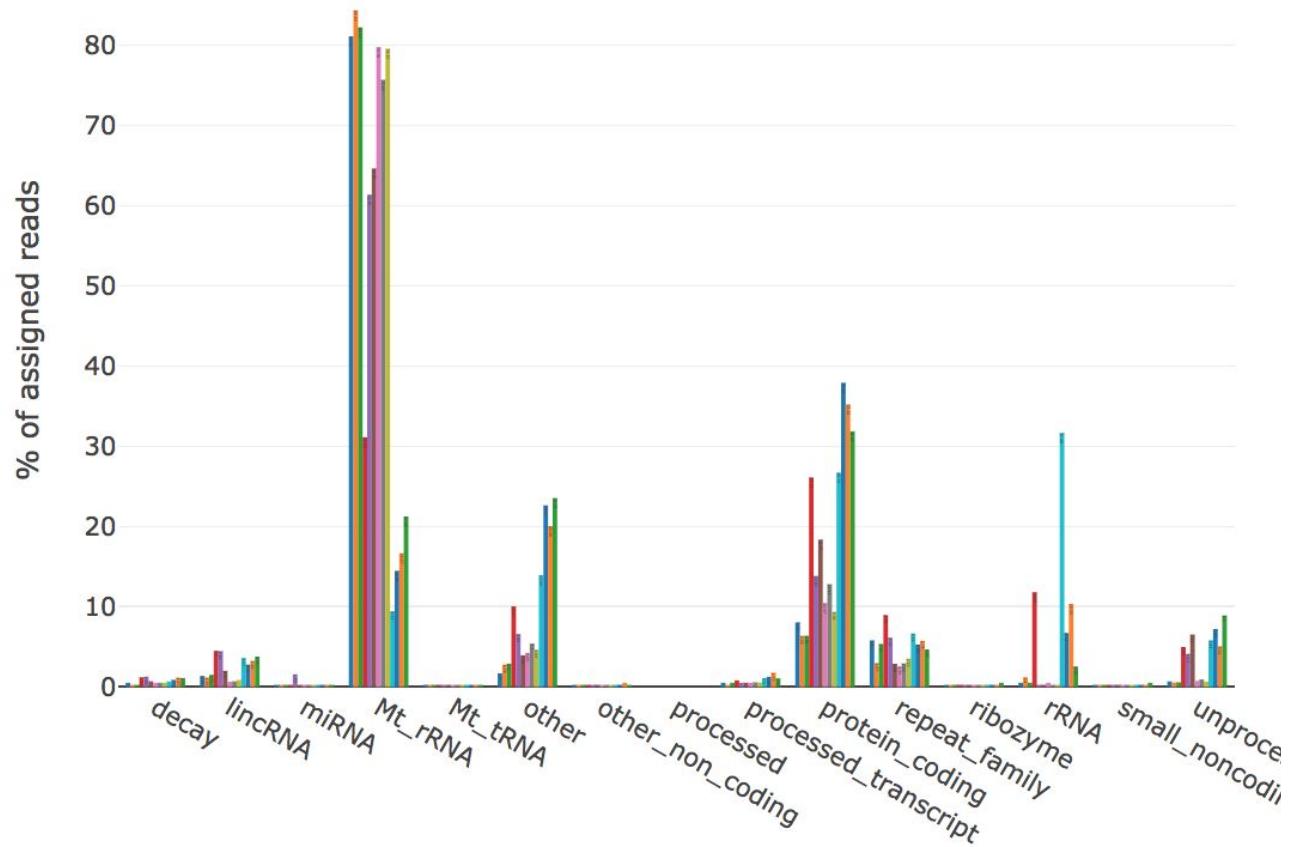
Number of samples

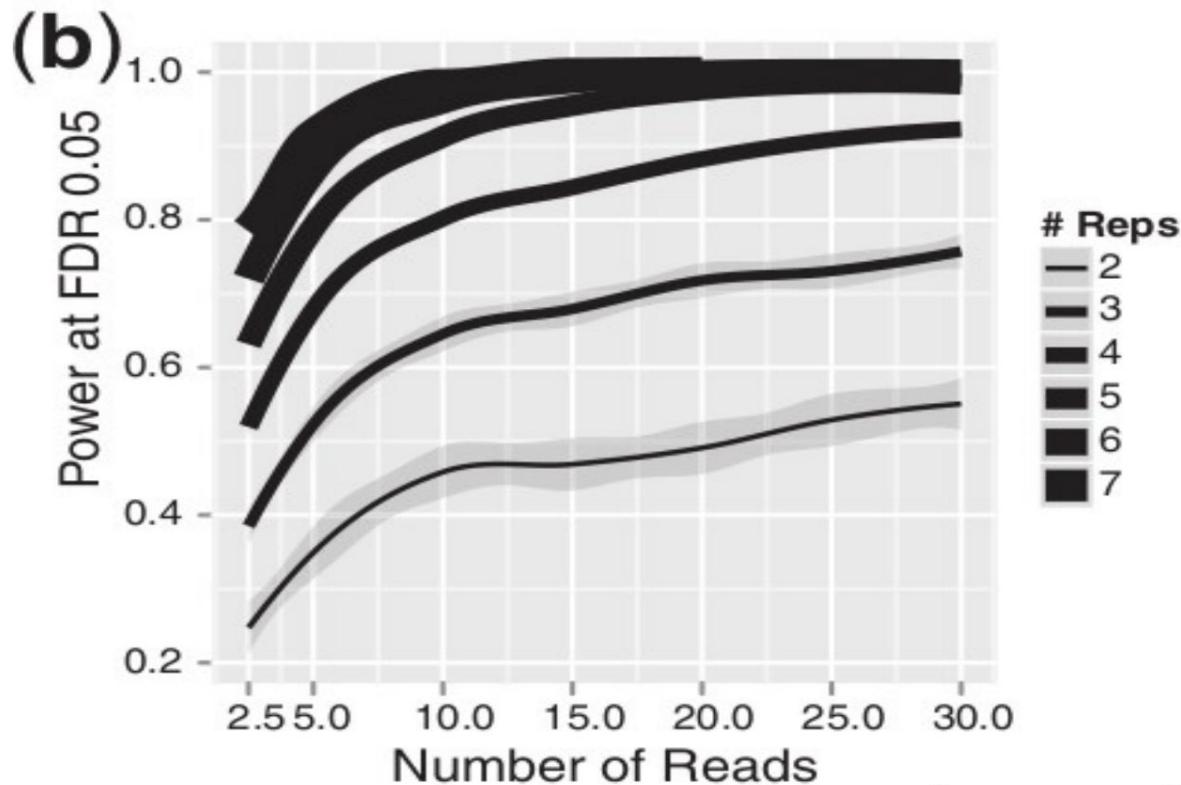
Single project

Should one have more replicates or more sequencing depth?

**Find problems in datasets generated, based on
previous figure.**

What is the problem for RNA-seq data for this graph?





Generally more replicates should be preferred if they can be obtained.

From Liu et al. 2014. RNA-seq differential expression studies: more sequence or more replication?

Raw data formats

FASTQ format

starting symbol → @HWI-EAS3X_10102_2_120_19829_1823#0/2 sequence identifier

sequence end → TCTAACTCTTACTTAGCATAGCTGTTAAAATTTTGAGTT sequence

start QS → +(optionally the same identifier) DEAEE:B:BE5EEEED=:DEA:-AE5DDBDFFEDEEDFAE quality score

Command to go to home directory

`cd`

`cd ~`

`cd $HOME`

Data for the course

```
wget https://bioinfo.evolution.uzh.ch/project.zip
```

Bash commands fresh up

Command to make a directory

Command to make a directory

mkdir

Command to change directory

Command to change directory

cd

Command to copy

Command to copy

cp

Command to rename

Command to rename

mv

Command to move

Command to move

mv

Command to go to home directory

Command to show your username

Command to show your username

whoami

Command to list files

Command to list files

ls

Command to show content of a file

Command to check if a command exists

Command to show content of a file

cat file

less file

more file

Command to make an empty file

Command to make an empty file

touch filename

Command to get help for a tool/ command

Command to get help for a tool/ command

help command

command -h

command --help

Command to softlink

Command to softlink

ln

Command to check if a command exists

which

Command to do 2nd task only after 1st task is done

Command to do 2nd task only after 1st task is done

task1 && task2

Quality control

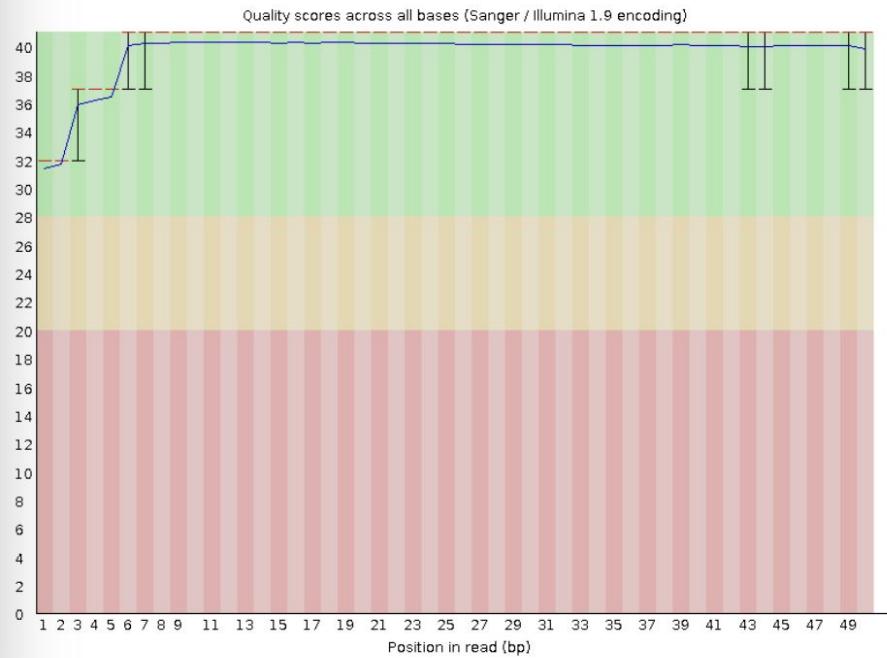
Basic statistics

FastQC

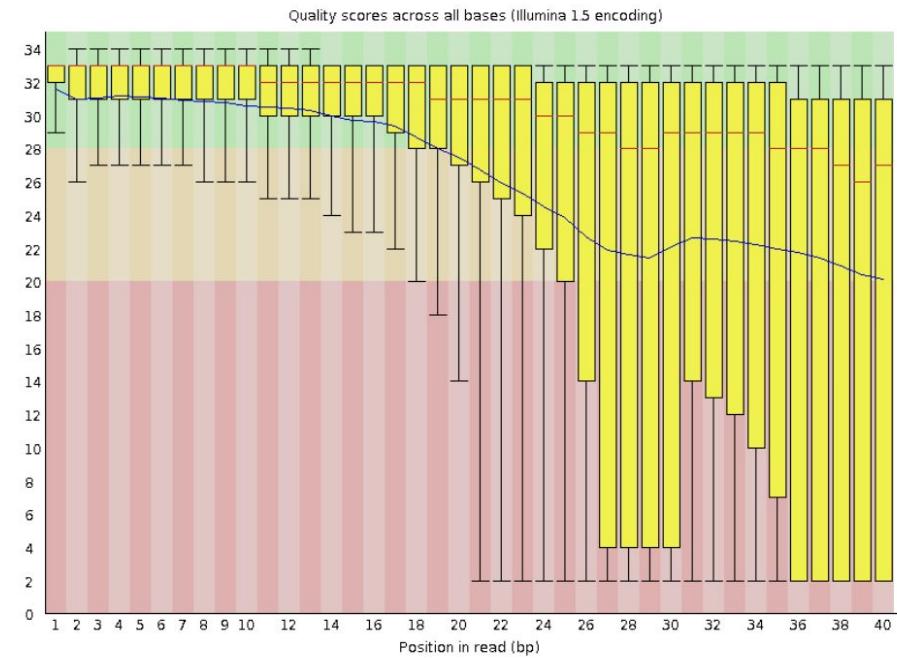
Measure	Value
Filename	Reference_sample
File type	Conventional base calls
Encoding	Sanger/ Illumina 1.9
Total sequences	143,077,301
Sequences flagged as poor quality	0
Sequence length	100
%GC	47
Duplicates	37,625,112
Sequencing type	Single end

Per base sequence quality

Good Data



Bad Data



+SEQ_ID

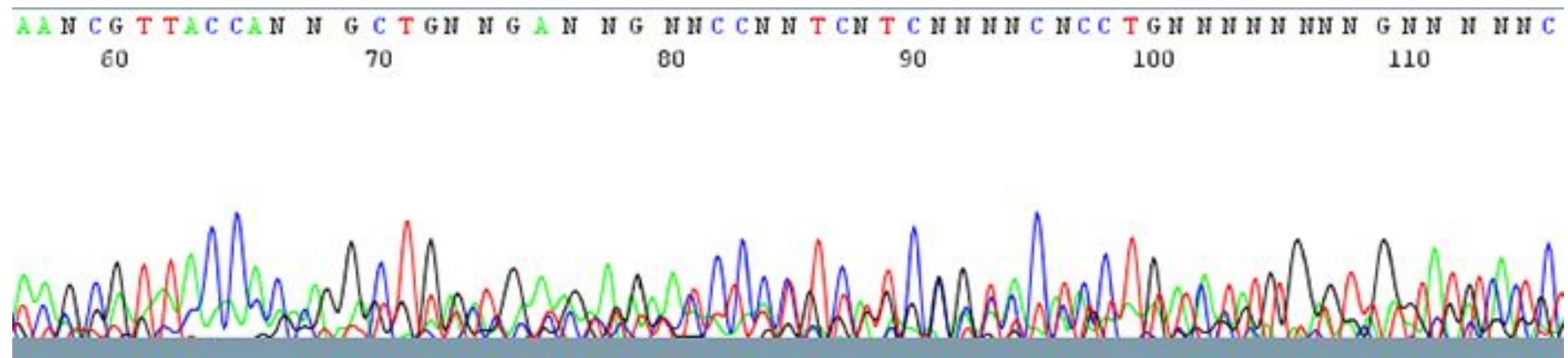
! ' ' * (((***+)) % % % ++) (% % % %) . 1 * *

A quality value Q is an integer representation of the probability p that the corresponding base call is incorrect.

$$Q = -10 \log_{10} P \quad \longrightarrow \quad P = 10^{\frac{-Q}{10}}$$

Phred Quality Score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%
40	1 in 10000	99.99%
50	1 in 100000	99.999%

Base cannot be determined

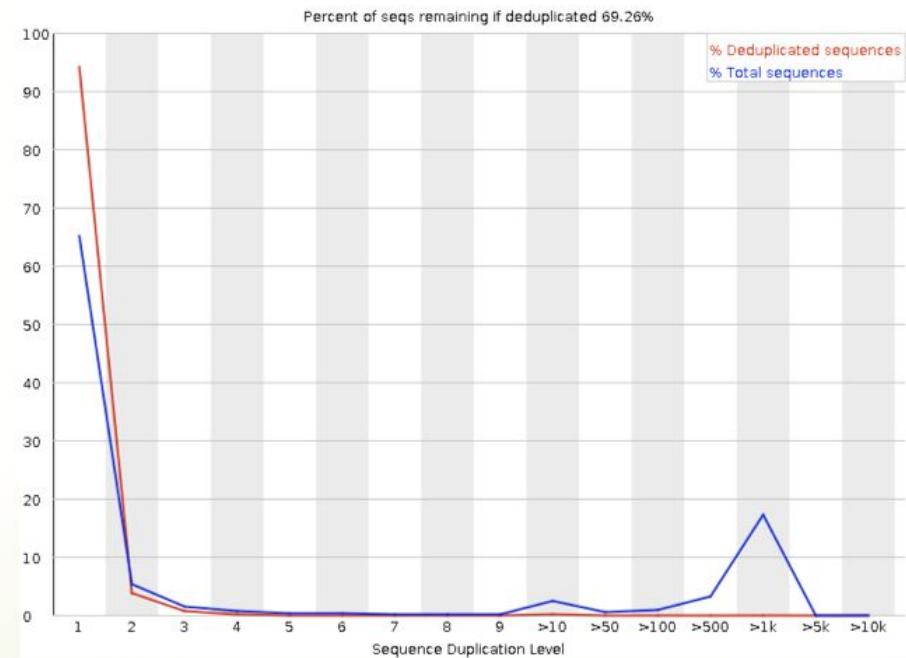
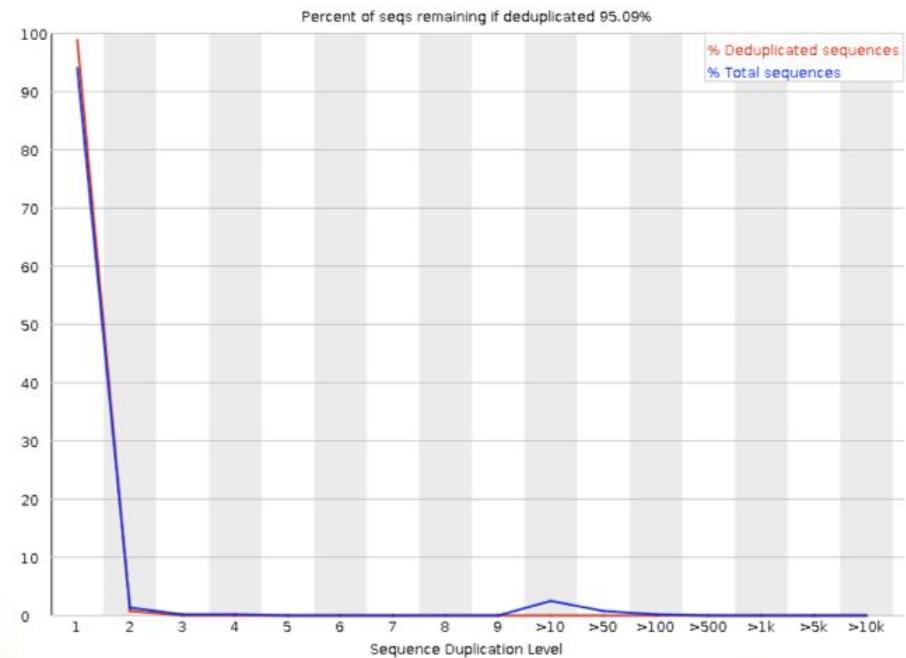


Removing bad data

NNNNNTGGCATATACGCGTGGCATATAADAPTER

- Remove N
- Remove Adapter
- Remove bases with bad quality
- Remove reads get shorter than 30 bp

Duplication problem



Performing quality check with fastqc

1. Check if fastqc is installed
2. Make a directory: analysis in the project directory
3. Go to analysis directory and make a directory: 01_qc_raw
4. Go to 01_qc_raw and make directories: log and output

Shell script vs Makefile

Shell script to run analysis

```
#!/usr/bin/bash or #!/bin/bash  
  
for i in ../../raw/*.gz  
  
do  
  
    fastqc $i -o ./output/ --noextract -f fastq -t 2  
  
done
```

What changes can you make to this script?

Shell script to run analysis

```
#!/usr/bin/bash or #!/bin/bash

for i in ../../raw/*.gz
do
    fastqc -v >./log/${i}.log && fastqc $i -o ./output/ --noextract -f fastq -t 2 2>>
./log/${i}.log
done
```

Shell script to run analysis

```
#!/usr/bin/bash  or #!/bin/bash

for i in ../../raw/*.gz
do
    fastqc -v >./log/${i}.stdout && fastqc $i -o ./output/ --noextract -f fastq -t 2 2>
    ./log/${i}.stderr
done
```

Power of Makefile

```
# This makefile will run the FastQC software to check the quality of FastQ files

SHELL:=/bin/bash
source_dir=../..raw
target_dir=./output
files := $(wildcard $(source_dir)/*.fq.gz)
targets := $(patsubst $(source_dir)/%.fq.gz, $(target_dir)/%_fastqc.zip, $(files))

all: $(targets)
$(target_dir)/%_fastqc.zip: $(source_dir)/%.fq.gz
    fastqc -v > ./log/$(basename $(notdir $@)).stdout && fastqc $< -o ./output/
--noextract -f fastq -t 2 2>./log/$(basename $(notdir $@)).stderr
```

Power of Makefile

```
# This makefile will run the FastQC software to check the quality of FastQ files

SHELL:=/bin/bash
source_dir=../data/raw
target_dir=./output
files := $(wildcard $(source_dir)/*.fq.gz)
targets := $(patsubst $(source_dir)/%.fq.gz, $(target_dir)/%_fastqc.zip, $(files))

all: $(targets)
$(target_dir)/%_fastqc.zip: $(source_dir)/%.fq.gz
    fastqc -v > ./log/$(basename $(notdir $@)).stdout && fastqc $< -o ./output/
--noextract -f fastq -t 2 2>./log/$(basename $(notdir $@)).stderr
```

Makefile for simple task

```
# Makefile to count unique lines
```

```
all: task1
```

```
task1:
```

```
    cat file.txt | uniq -c
```

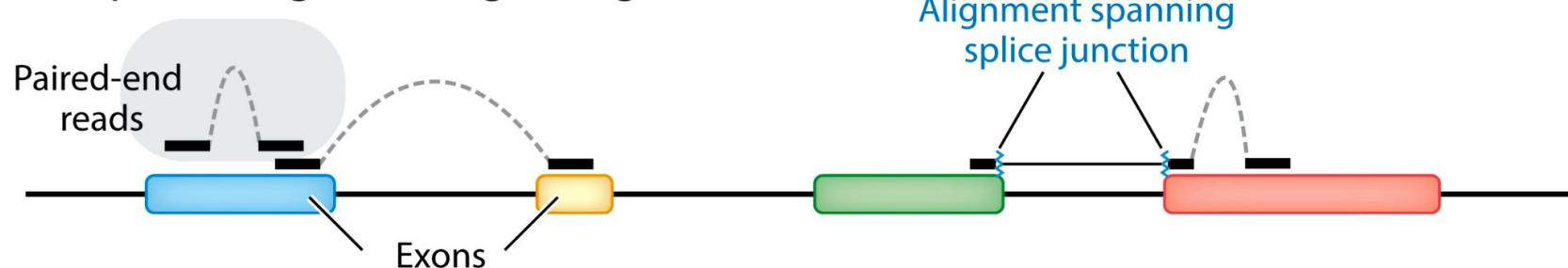
Removing bad data TrimGalore

NNNNNTGGCATATACGCGTGGCATATAADAPTER

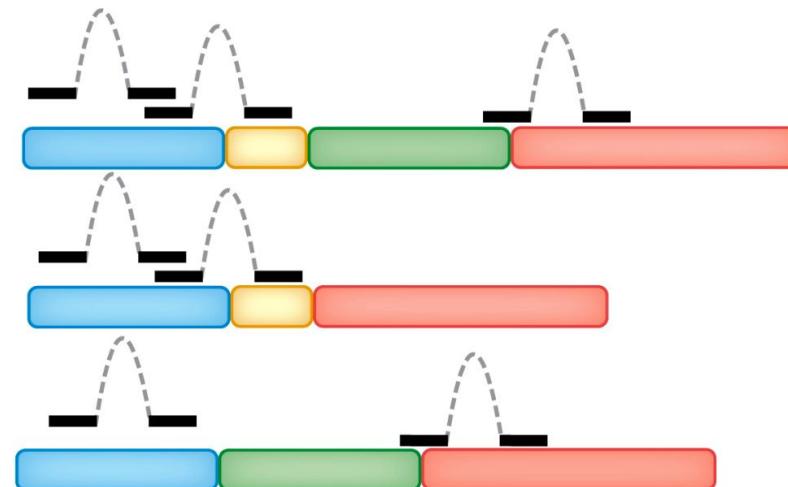
- Remove N
- Remove Adapter
- Remove bases with bad quality (phred scores less than 30)
- Remove reads get shorter than 30 bp

Alignment and read counts

a Spliced alignment against genome



b Unspliced alignment against transcriptome



Alignment

<https://youtu.be/4WRANhDiSHM>

<https://tinyheero.github.io/2015/09/02/pseudoalignments-kallisto.html>

Watch the video and read the link

Why not use BLAST?

Why not use BLAST?

Aligner	Human reference runtime (hrs)	Max mem used (GB)	Number of AMD 64 bit core processors
Bowtie2	0.62	9	17
BWA	0.66	9	17
BLAST	9.4	12	17

Questions

Question 1

What is the application of `BWT` in biology?

Question 2

Which tools are built on `BWT` for biology?

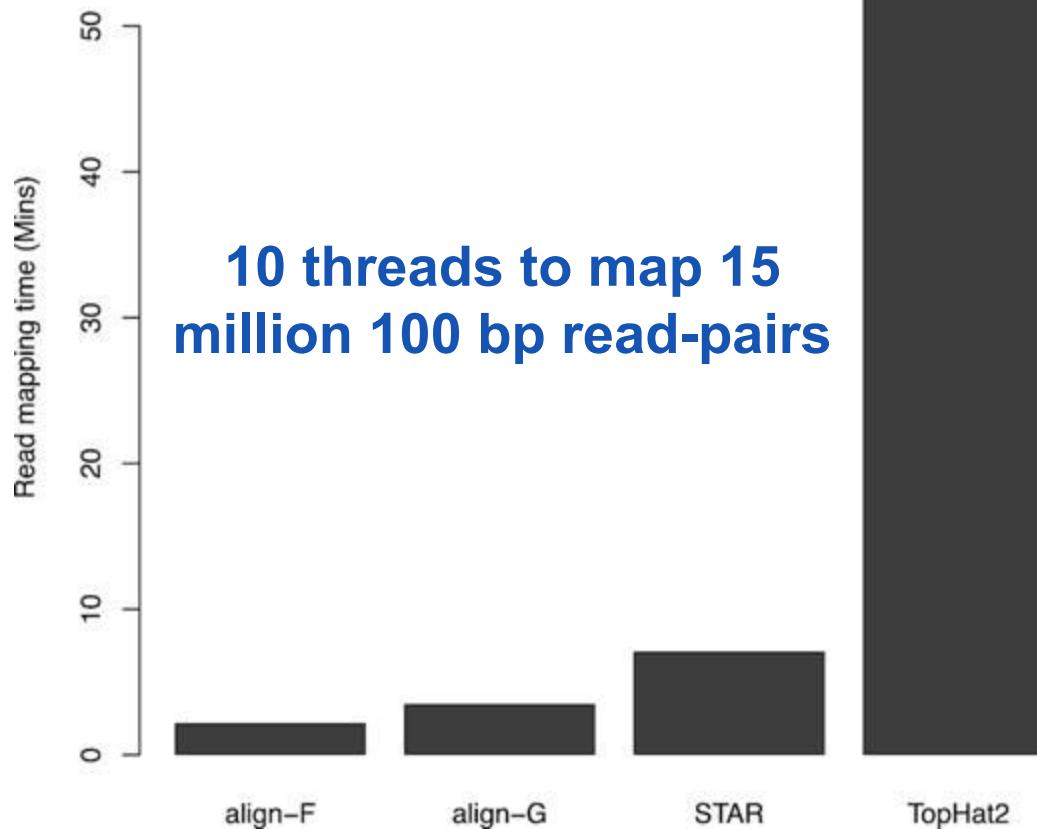
Question 3

Difference between alignment and `Pseudoalignment`?

Question 4

Which tools can perform `Pseudoalignment` and for which `*seq` data?

Why to use Rsubread?



Files needed for alignment

1. Trimmed FASTQ files
2. A Reference genome
3. Gene definition file
4. Tools to run alignment

<http://localhost:8780/>

rstudio

pass

SAM format

A

Coor	12345678901234	10	20	30	40
ref	AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGCGCCAT				
+r001/1	TTAGATAAAGGATA*CTG				
+r002	aaaAGATAA*GGATA				
+r003	gcctaAGCTAA				
+r004	ATAGCT.....TCAGC				
-r003	ttagctTAGGC				
-r001/2	CAGCGGCAT				

B

Header section @HD VN:1.5 SO:coordinate @SQ SN:ref LN:45											
Alignment section r001 99 ref 7 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATAACTG *											
r002 0 ref 9 30 3S6M1P1I4M * 0 0 AAAAGATAAGGATA *											
r003 0 ref 9 30 5S6M * 0 0 GCCTAACGCTAA * 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											

QNAME	FLAG	RNAME	POS	MAPQ	CIGAR	RNEXT	PNEXT	TLEN	SEQ	Optional fields in the format of TAG:TYPE:VALUE
(query template name, aka. read ID)	(indicates alignment information about the read, e.g. paired, aligned, etc.)	(reference sequence name, e.g. chromosome /transcript id)	(1-based position)	(mapping quality)	(summary of alignment, e.g. insertion, deletion)	(reference sequence name of the primary alignment of the NEXT read; for paired-end sequencing, NEXT read is the paired read; corresponding to the RNAME column)	(Position of the primary alignment of the NEXT read; for paired-end sequencing, corresponding to the POS column)	(the number of bases covered by the reads from the same fragment. In this particular case, it's 45 - 7 + 1 = 39 as highlighted in Panel A). Sign: plus for leftmost read, and minus for rightmost read	(read sequence)	

SAM Flags

Bit	Description
1	0x1 template having multiple segments in sequencing
2	0x2 each segment properly aligned according to the aligner
4	0x4 segment unmapped
8	0x8 next segment in the template unmapped
16	0x10 SEQ being reverse complemented
32	0x20 SEQ of the next segment in the template being reverse complemented
64	0x40 the first segment in the template
128	0x80 the last segment in the template
256	0x100 secondary alignment
512	0x200 not passing quality controls
1024	0x400 PCR or optical duplicate
2048	0x800 supplementary alignment

Example, flag 83 = 64+16+2+1 means it's first read (0x40) of pair-end reads (0x1) and it's mapped on minus strand (0x10) and both reads mapped (0x2).

<https://broadinstitute.github.io/picard/explain-flags.html>

Questions

1. Count the number of aligned reads in the BAM file

samtools view -c sample.bam

2. Store the unaligned reads from bam file in a new file

samtools view -bh -f 4 sample.bam > unmapped.bam

3. Display only header of the BAM file

samtools view -H sample.bam

4. Display first few lines of a BAM file with header

samtools view -h sample.bam | head -n 10

5. Display first few lines of a BAM file without header

samtools view sample.bam | head -n 10

Duplicated reads: which one is technical/ biological?

Reference ATCGGACTACTAACCTCGCGCGATATAC

Read 1
ATCGGACT

ATCGGACT

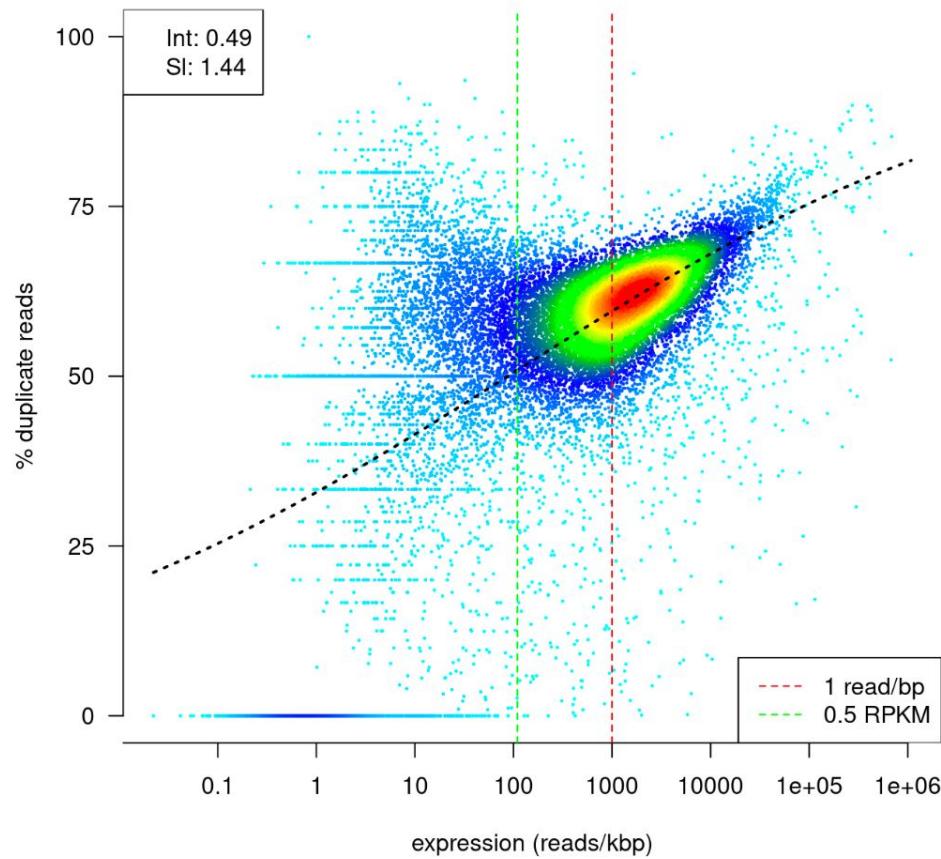
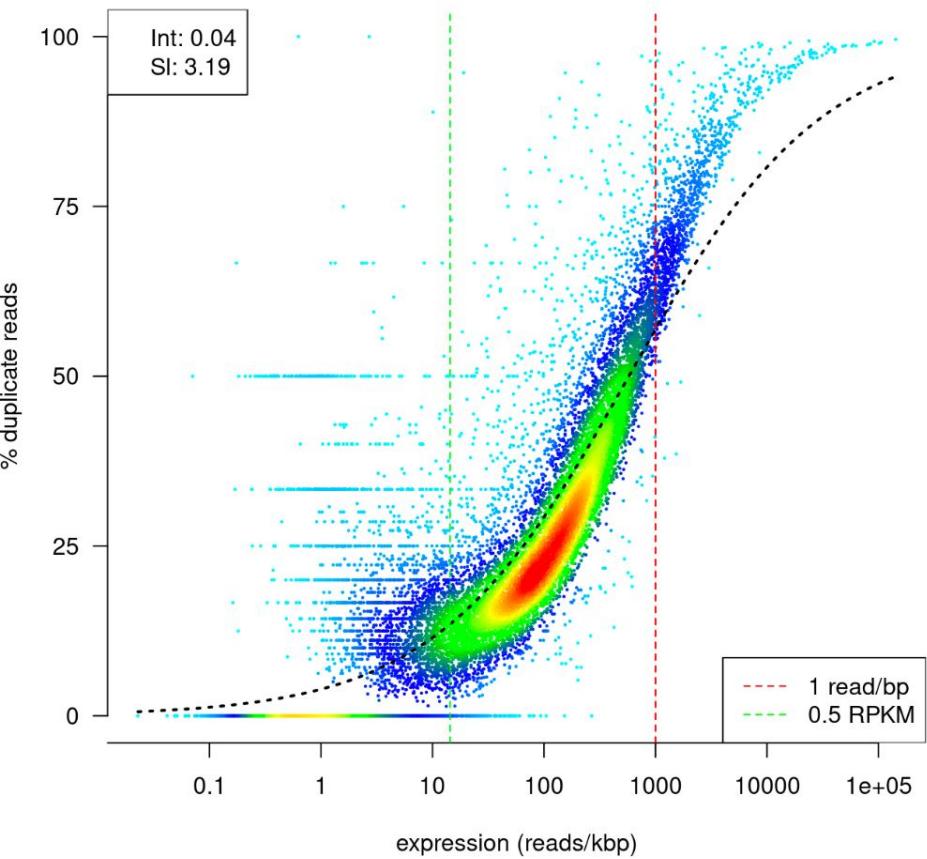
ACTACTAACCTCGCG

ACTACTAACCTCGCG

Reads duplication problem in RNA-seq

1. Should we remove duplicates, why or why not?
2. Which library preparation method would be best to distinguish between technical and biological duplicated reads?
3. How many PCR cycles should one use for optimal technical duplicates?
4. Is paired-end sequencing better than single-end sequencing for duplication problem?

What does this plot tell you?



**Make a flowchart of RNA-seq
pipeline so far**

Psudo-aligners

k -mers are nucleotides of length k

- Oct4 is 1574 nt long ($L = 1574$)
- $k = 7$ ($k= 7$)
- Oct4 will contain 1568 Kmers ($L-k+1$)

.....CTTGGAACAAAT.....

CTTGGAA

TTGGAAC

TGGAACA

GGAACAA

GAACAAAT

Splits up a read into the same k-mer size

Read1 = CTTGGAAACAAT

Kmer Read1
CTTGGAA
TTGGAAC
TGGAACA
GGAACAA
GAACAAT
AACAATA

Checks in which transcripts the Kmers exist and sums them up

Kmer Read1	Kmer	Oct4	Oct3	Oct2	Sox2	Sox3
CTTGGAA	CTTGGAA	TRUE	FALSE	TRUE	FALSE	FALSE
TTGGAAC	TTGGAAC	TRUE	FALSE	TRUE	FALSE	FALSE
TGGAACA	TGGAACA	TRUE	FALSE	FALSE	FALSE	FALSE
GGAACAA	GGAACAA	TRUE	TRUE	FALSE	FALSE	FALSE
GAACAAAT	GAACAAAT	TRUE	TRUE	FALSE	FALSE	FALSE
AACAATA	AACAATA	TRUE	FALSE	FALSE	FALSE	FALSE



Read	Oct4	Oct3	Oct2	Sox2	Sox3
Read 1	6	2	2	0	0

Assign the read to one or many transcript

Checks which of the transcripts the number of kmers matched is least likely to happen by chance and assign it to those transcript

Read	Oct4	Oct3	Oct2	Sox2	Sox3
Read 1	6	2	2	0	0



Assign read to transcripts

Read	Oct4
Read 1	1



Add read counts to transcripts in sample

Read	Oct4	Oct3	Oct2	Sox2	Sox3
Sample 1	+1	0	0	0	0

Redo the procedure for all reads

Read2 = GATACAGATAC
6 kmers of length 7

Read	Oct4	Oct3	Oct2	Sox2	Sox3
Read 2	0	0	0	6	6



Assign read to transcripts

Read	Sox2	Sox3
Read 2	1	1

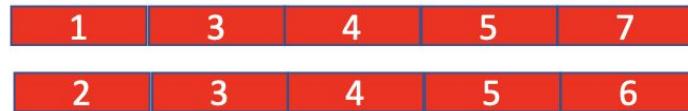
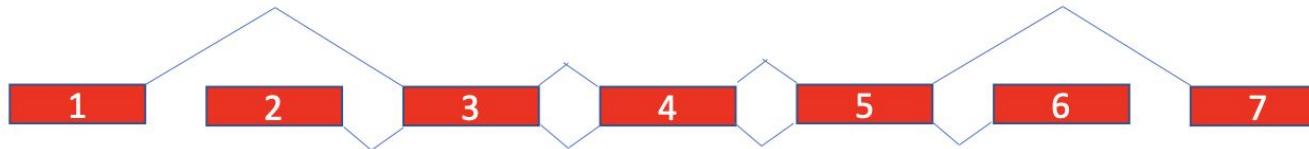


Add read counts to transcripts in sample

Read	Oct4	Oct3	Oct2	Sox2	Sox3
Sample 1	1	0	0	+1	+1

Kalisto builds a de Bruijn graph with all the k-mers

Red gene

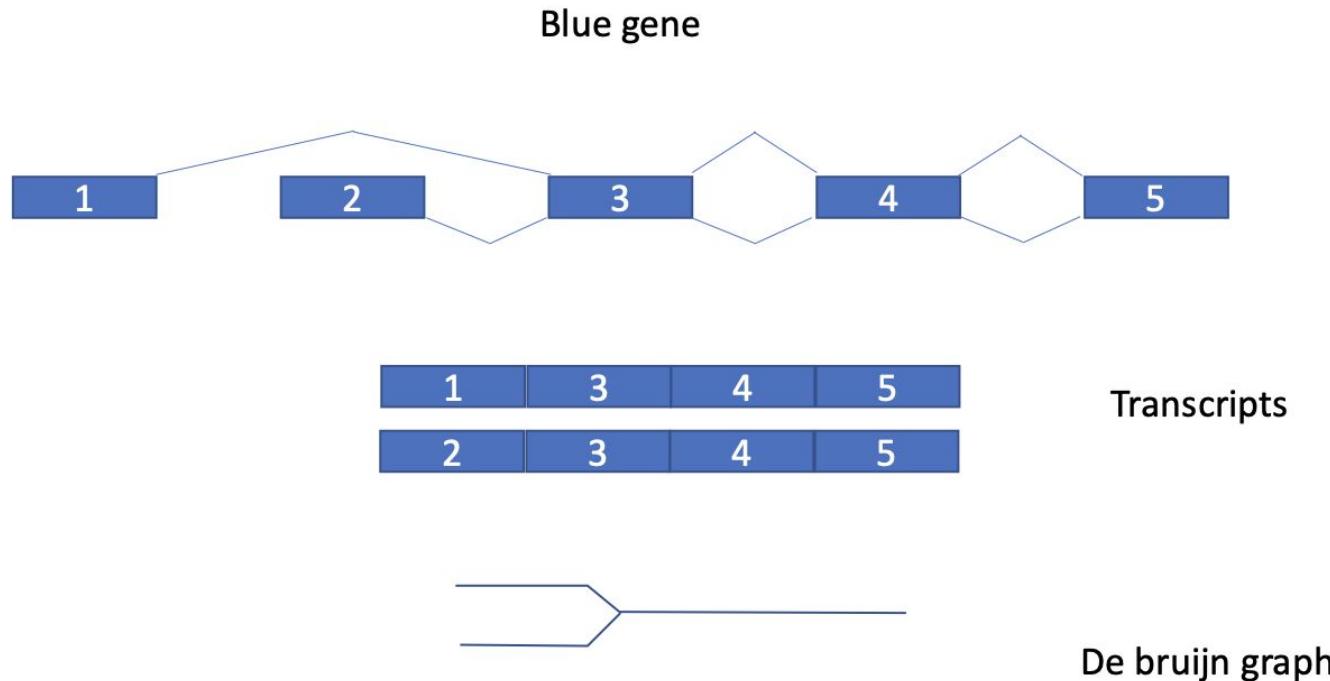


Transcripts

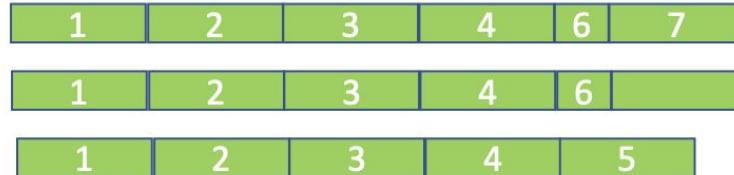


De bruijn graph

Kalisto builds a graph with all the k-mers



Kalisto builds a graph with all the k-mers

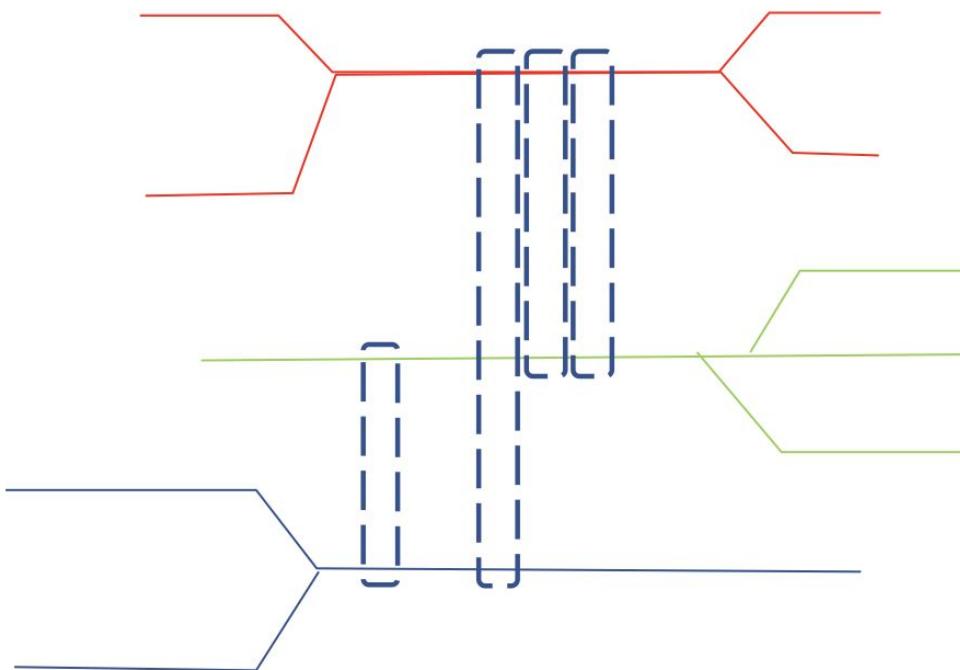


Transcripts



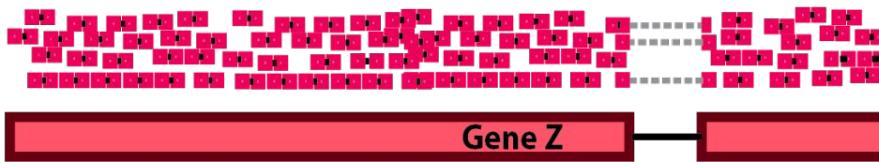
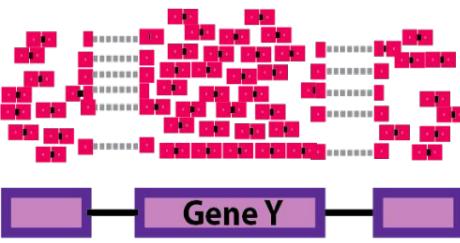
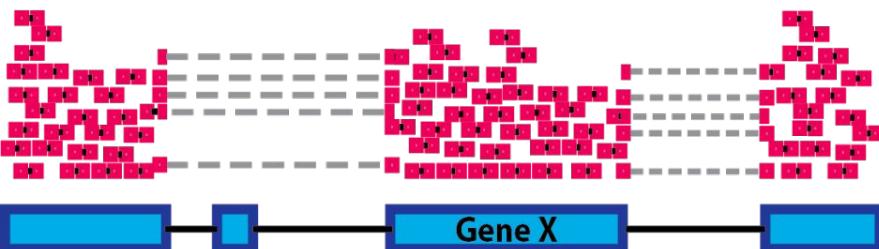
De bruijn graph

3 gene de-bruijn graph with parts in common.

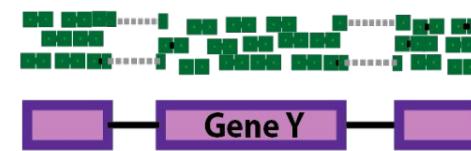
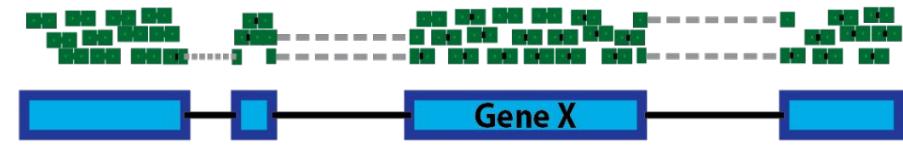


Normalisation

Sample A Reads



Sample B Reads



Aims of normalizing your data

1. Make the data comparable across features (compare genes)
2. Make data comparable across libraries (different samples)

Normalization/ scaling/ transformation: different goals

- **R/FPKM:** (Mortazavi et al. 2008)
 - **Correct for:** differences in sequencing depth and transcript length
 - **Aiming to:** compare a gene across samples and diff genes within sample
- **TMM:** (Robinson and Oshlack 2010)
 - **Correct for:** differences in transcript pool composition; extreme outliers
 - **Aiming to:** provide better across-sample comparability
- **TPM:** (Li et al 2010, Wagner et al 2012)
 - **Correct for:** transcript length distribution in RNA pool
 - **Aiming to:** provide better across-sample comparability
- **Limma voom (logCPM):** (Law et al 2013)
 - **Aiming to:** stabilize variance; remove dependence of variance on the mean

RPKM: Reads Per Kilobase Million

Gene name	Rep1 counts	Rep2 counts	Rep3 counts
A (2kb)	10	12	30
B (4kb)	20	25	60
C (1kb)	5	8	15
D (10kb)	0	0	1

RPKM: Reads Per Kilobase Million

Gene name	Rep1 counts	Rep2 counts	Rep3 counts
A (2kb)	10	12	30
B (4kb)	20	25	60
C (1kb)	5	8	15
D (10kb)	0	0	1
Total	35	45	106

RPKM: Reads Per Kilobase Million

Gene name	Rep1 counts	Rep2 counts	Rep3 counts
A (2kb)	10	12	30
B (4kb)	20	25	60
C (1kb)	5	8	15
D (10kb)	0	0	1
Total	35	45	106
Scaling factor of 1M (we use 10)	3.5	4.5	20.6

RPKM: Reads Per Kilobase Million

Gene name	Rep1 RPM	Rep2 RPM	Rep3 RPM
A (2kb)	10/3.5	12/4.5	30/20.6
B (4kb)	20/3.5	25/4.5	60/20.6
C (1kb)	5/3.5	8/4.5	15/20.6
D (10kb)	0/3.5	0/4.5	1/20.6
Scaling factor of 1M (we use 10)	3.5	4.5	20.6

RPKM: Reads Per Kilobase Million

Gene name	Rep1 RPM	Rep2 RPM	Rep3 RPM
A (2kb)	2.86	2.67	2.83
B (4kb)	5.71	5.56	5.66
C (1kb)	1.43	1.78	1.42
D (10kb)	0	0	0.09

RPKM: Reads Per Kilobase Million

Gene name	Rep1 RPKM	Rep2 RPKM	Rep3 RPKM
A (2kb)	2.86/2	2.67/2	2.83/2
B (4kb)	5.71/4	5.56/4	5.66/4
C (1kb)	1.43/1	1.78/1	1.42/1
D (10kb)	0/10	0/10	0.09/10

RPKM: Reads Per Kilobase Million

Gene name	Rep1 RPKM	Rep2 RPKM	Rep3 RPKM
A (2kb)	1.43	1.33	1.42
B (4kb)	1.43	1.39	1.42
C (1kb)	1.43	1.78	1.42
D (10kb)	0	0	0.009

TPM: Transcript Per Million

Gene name	Rep1 counts	Rep2 counts	Rep3 counts
A (2kb)	10	12	30
B (4kb)	20	25	60
C (1kb)	5	8	15
D (10kb)	0	0	1

TPM: Transcript Per Million

Gene name	Rep1 RPK	Rep2 RPK	Rep3 RPK
A (2kb)	10/2	12/2	30/2
B (4kb)	20/4	25/4	60/4
C (1kb)	5/1	8/1	15/1
D (10kb)	0/10	0/10	1/10

TPM: Transcript Per Million

Gene name	Rep1 RPK	Rep2 RPK	Rep3 RPK
A (2kb)	5	6	15
B (4kb)	5	6.25	15
C (1kb)	5	8	15
D (10kb)	0	0	0.1

TPM: Transcript Per Million

Gene name	Rep1 RPK	Rep2 RPK	Rep3 RPK
A (2kb)	5	6	15
B (4kb)	5	6.25	15
C (1kb)	5	8	15
D (10kb)	0	0	0.1
Total RPK	15	20.25	45.1

TPM: Transcript Per Million

Gene name	Rep1 RPK	Rep2 RPK	Rep3 RPK
A (2kb)	5	6	15
B (4kb)	5	6.25	15
C (1kb)	5	8	15
D (10kb)	0	0	0.1
Total RPK	15	20.25	45.1
Scaling factor of 1M (we use 10)	1.5	2.025	4.51

TPM: Transcript Per Million

Gene name	Rep1 TPM	Rep2 TPM	Rep3 TPM
A (2kb)	5/1.5	6/2.025	15/4.51
B (4kb)	5/1.5	6.25/2.025	15/4.51
C (1kb)	5/1.5	8/2.025	15/4.51
D (10kb)	0/1.5	0/2.025	0.1/4.51
Scaling factor of 1M (we use 10)	1.5	2.025	4.51

TPM: Transcript Per Million

Gene name	Rep1 TPM	Rep2 TPM	Rep3 TPM
A (2kb)	3.33	2.96	3.326
B (4kb)	3.33	3.09	3.326
C (1kb)	3.33	3.95	3.326
D (10kb)	0	0	0.02

RPKM vs TPM

What difference
do you see?

Gene name	Rep1 RPKM	Rep2 RPKM	Rep3 RPKM
A (2kb)	1.43	1.33	1.42
B (4kb)	1.43	1.39	1.42
C (1kb)	1.43	1.78	1.42
D (10kb)	0	0	0.009
Total	4.29	4.5	4.25

Gene name	Rep1 TPM	Rep2 TPM	Rep3 TPM
A (2kb)	3.33	2.96	3.326
B (4kb)	3.33	3.09	3.326
C (1kb)	3.33	3.95	3.326
D (10kb)	0	0	0.02
Total	10	10	10

How does TMM work?

<https://youtu.be/Wdt6jdi-NQo>

**Is TMM a intra sample or inter sample
normalisation?**

Reading method papers as a **bioinformatics user**

- Why do you want to read this paper?
- Where is the method available from?
- What is the niche?
- What should you look for in results?
- Would this method be useful for your research?

How can normalization be improved?

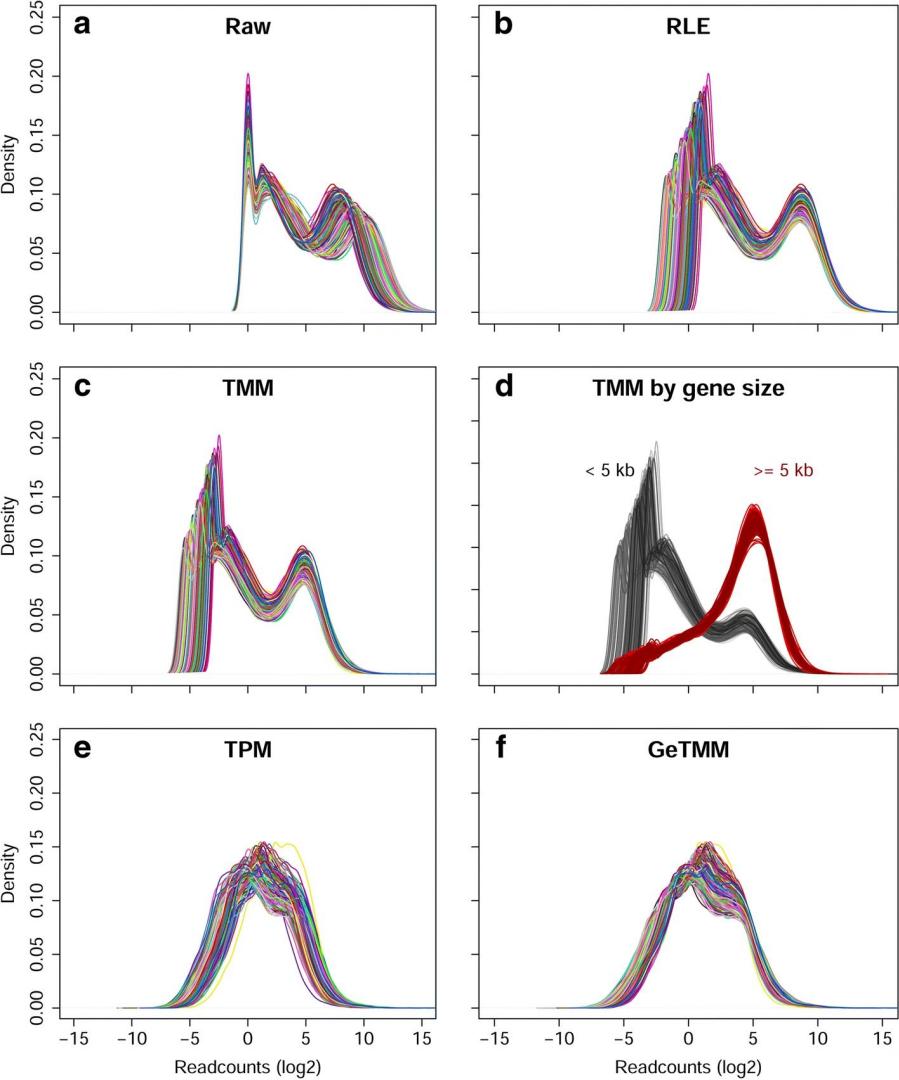
Methodology article | [Open Access](#) | Published: 22 June 2018

Gene length corrected trimmed mean of M-values (GeTMM) processing of RNA-seq data performs similarly in intersample analyses while improving intrasample comparisons

Marcel Smid , Robert R. J. Coebergh van den Braak, Harmen J. G. van de Werken, Job van Riet, Anne van Galen, Vanja de Weerd, Michelle van der Vlugt-Daane, Sandra I. Bril, Zarina S. Lalmahomed, Wigard P. Kloosterman, Saskia M. Wilting, John A. Foekens & Jan N. M. IJzermans, on behalf of the MATCH study group, John W. M. Martens & Anieta M. Sieuwerts

[BMC Bioinformatics](#) **19**, Article number: 236 (2018) | [Cite this article](#)

19k Accesses | **59** Citations | **3** Altmetric | [Metrics](#)



SVA

Published online 07 October 2014

Nucleic Acids Research, 2014, Vol. 42, No. 21 e161
doi: 10.1093/nar/gku864

svaseq: removing batch effects and other unwanted noise from sequencing data

Jeffrey T. Leek*

Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health Baltimore, MD 21212, US

Received June 24, 2014; Revised August 20, 2014; Accepted September 08, 2014

Article | [Open Access](#) | [Published: 15 September 2022](#)

Removing unwanted variation from large-scale RNA sequencing data with PRPS

[Ramyar Molania](#) , [Momeneh Foroutan](#), [Johann A. Gagnon-Bartsch](#), [Luke C. Gandolfo](#), [Aryan Jain](#),
[Abhishek Sinha](#), [Gavriel Olshansky](#), [Alexander Dobrovic](#), [Anthony T. Papenfuss](#)  & [Terence P. Speed](#) 

[Nature Biotechnology](#) **41**, 82–95 (2023) | [Cite this article](#)

21k Accesses | **3** Citations | **92** Altmetric | [Metrics](#)

Group exercise

<https://www.nature.com/articles/s41587-022-01440-w>

“Removing unwanted variation from large-scale RNA sequencing data with PRPS“

How to use their method?

Do you want to use their method, why?

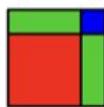
PCA

What is PCA?

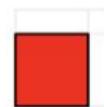
Google: “what is a PCA: stack exchange”

NMF

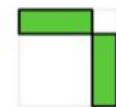
$$(a+b)^2$$



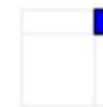
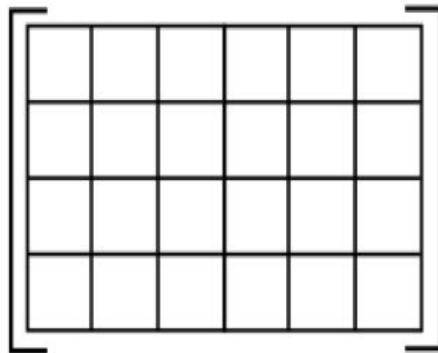
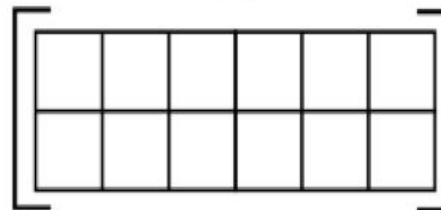
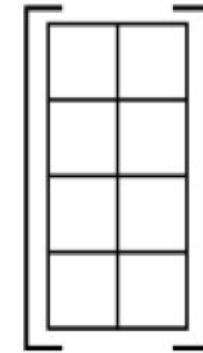
$$= a^2$$


$$+$$

$$2ab$$

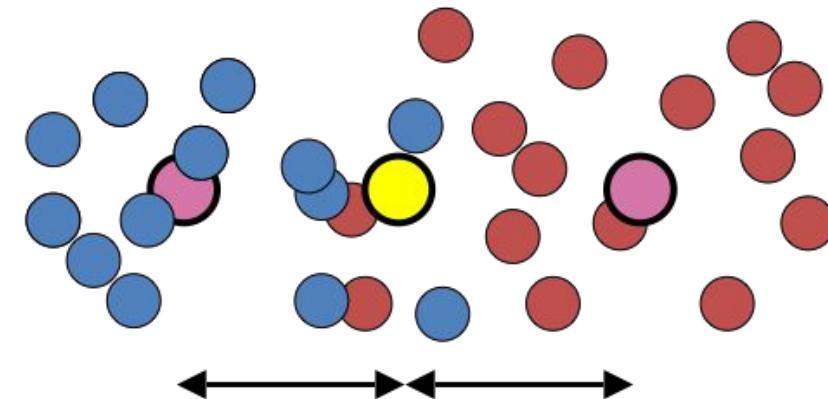
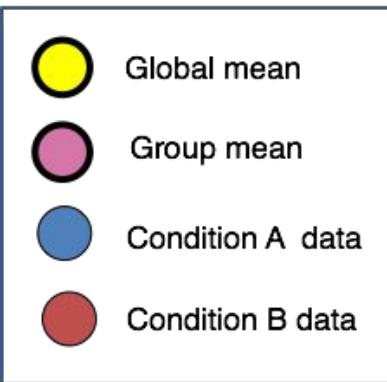

$$+$$

$$b^2$$


$$V$$

$$\approx$$
$$H$$

$$\times$$
$$W$$


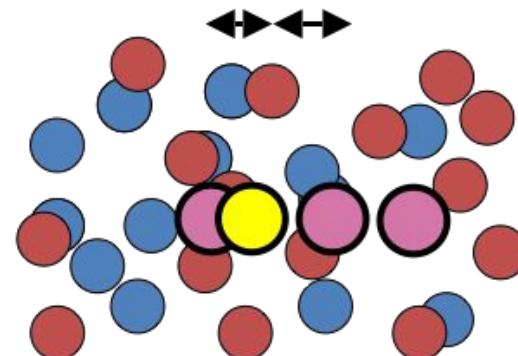
Differential analysis

Expression level



Significant
difference

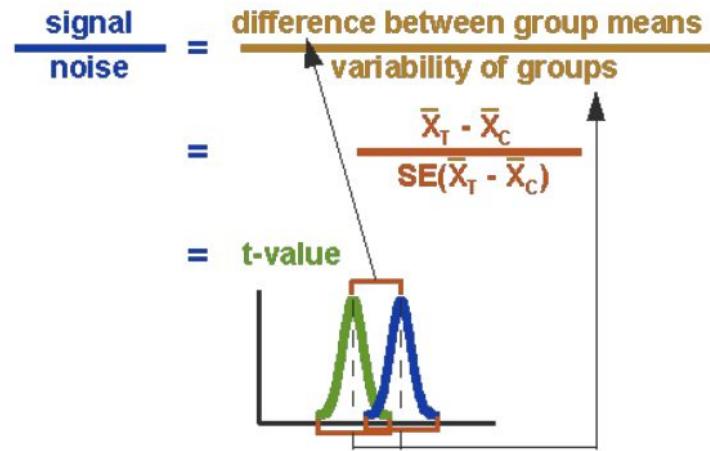
Deviations from global mean



No significant difference

Differential expression analysis?

Couldn't we just use a Student's t test for each gene?

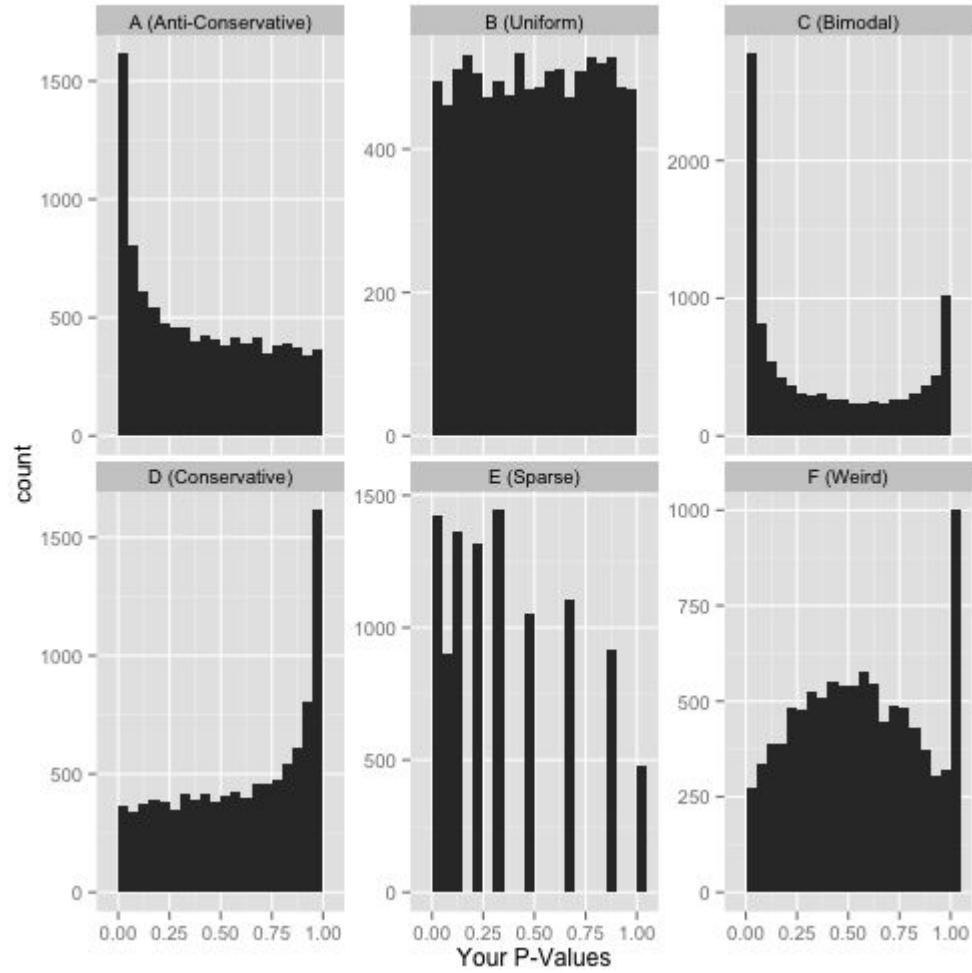


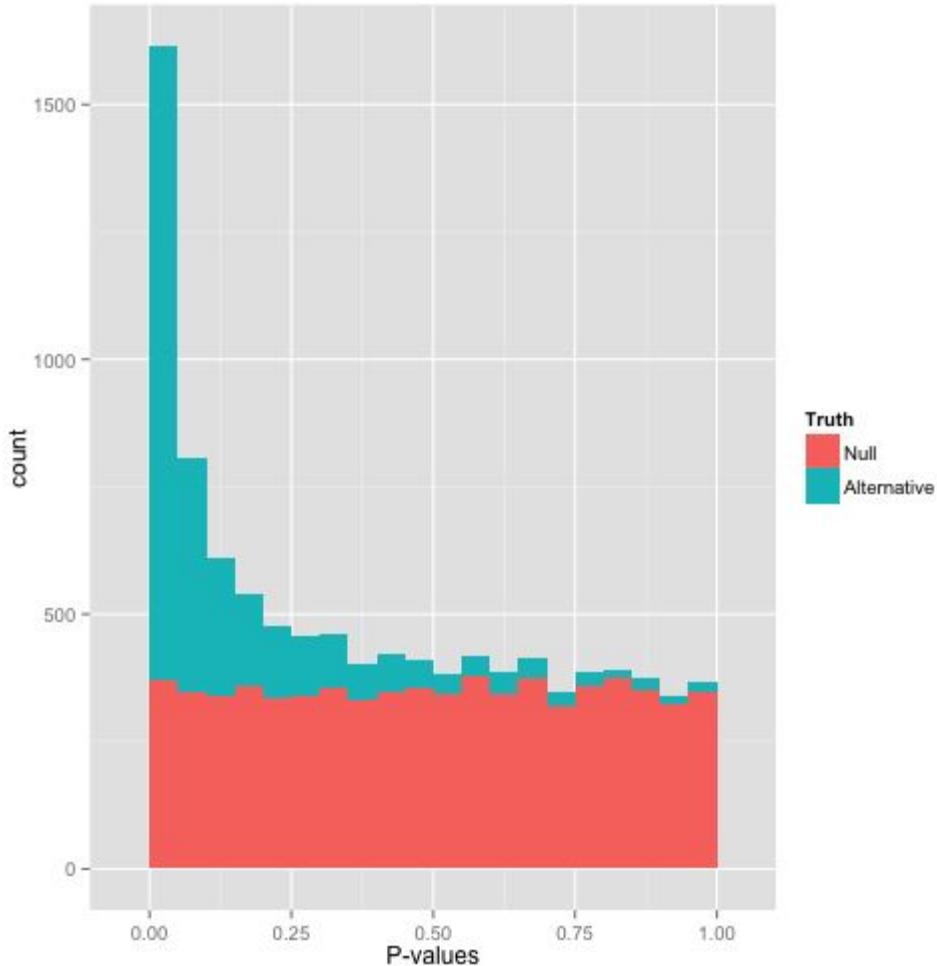
Problems with this approach:

http://www.socialresearchmethods.net/kb/stat_t.php

- May have **few replicates**
- **Multiple testing issues**
- Distribution is **not normal**

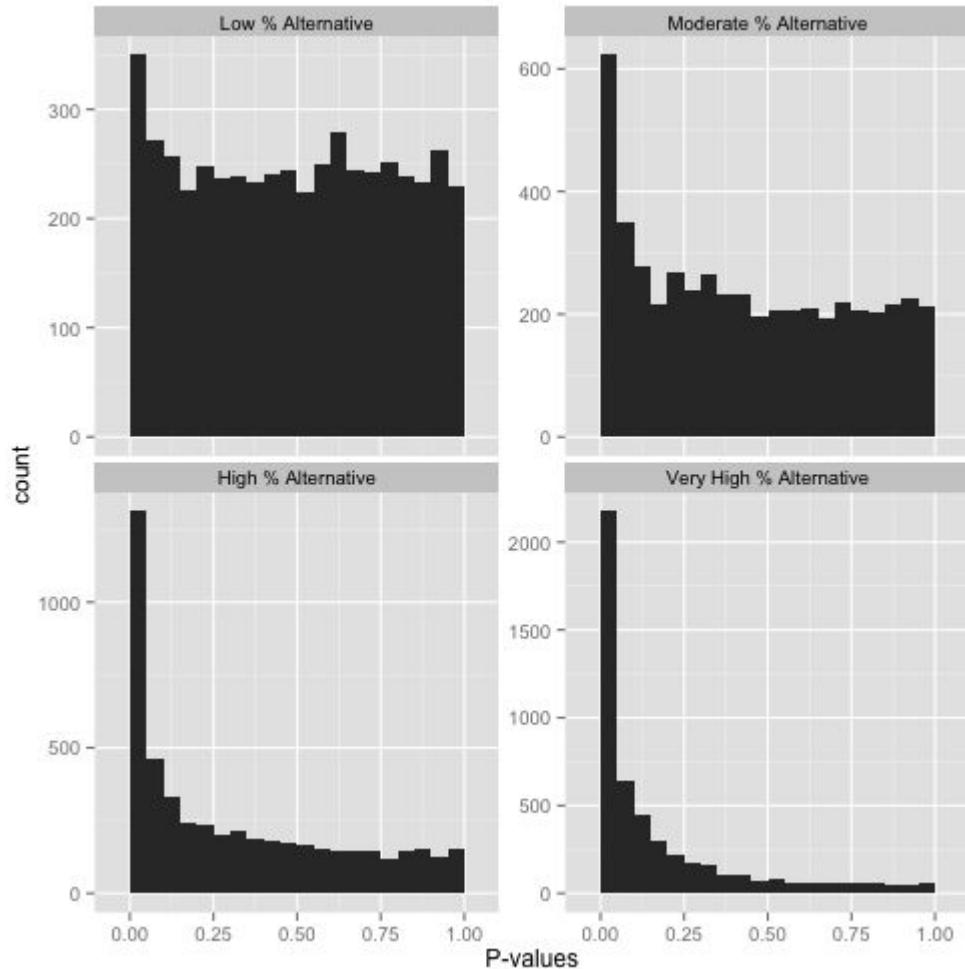
P-value distributions





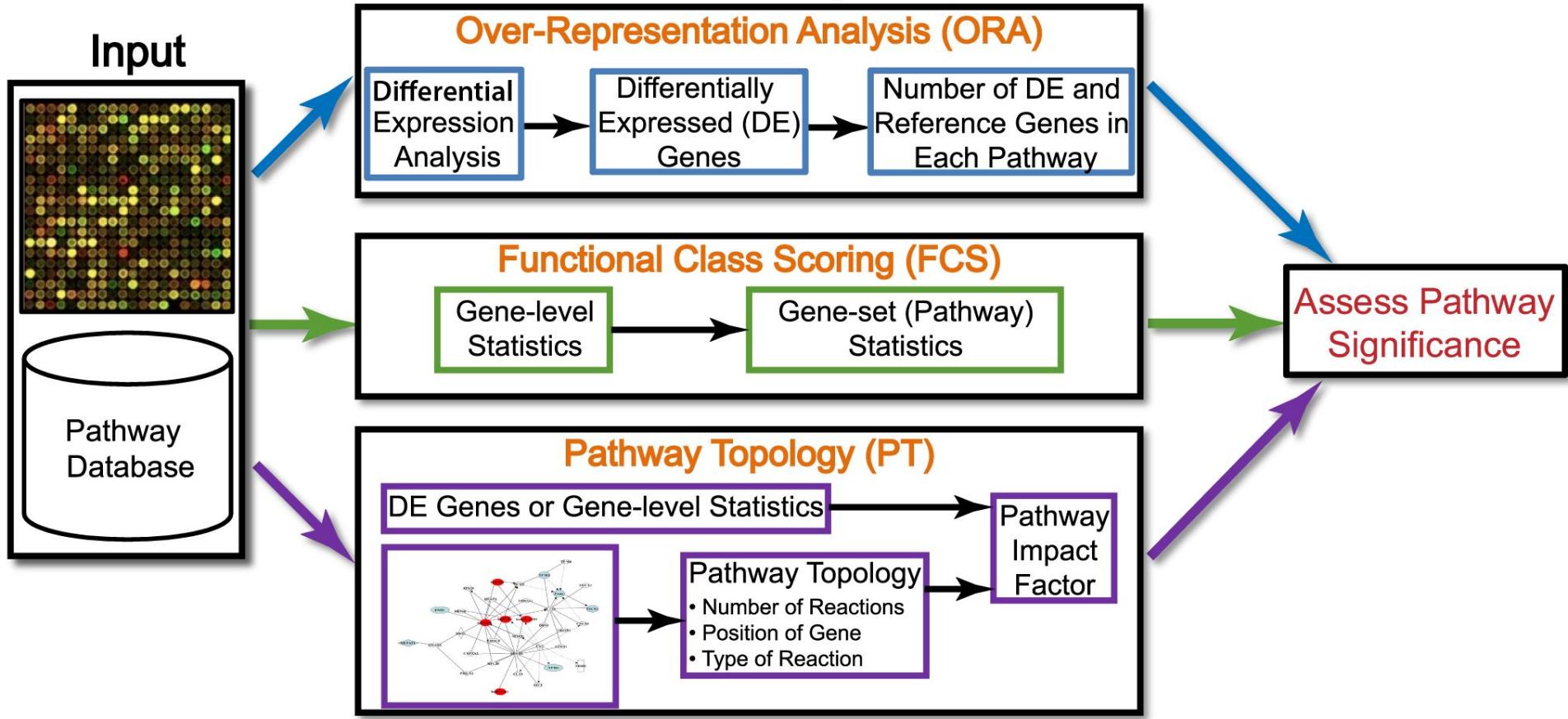
Anti-conservation: as per definition of a *p*-value:
under the null, it has a 5% chance of being less than .05, a 10% chance of being less than .1, etc.

Different anti-conservative distributions



Functional analysis

Functional Pathway Analysis

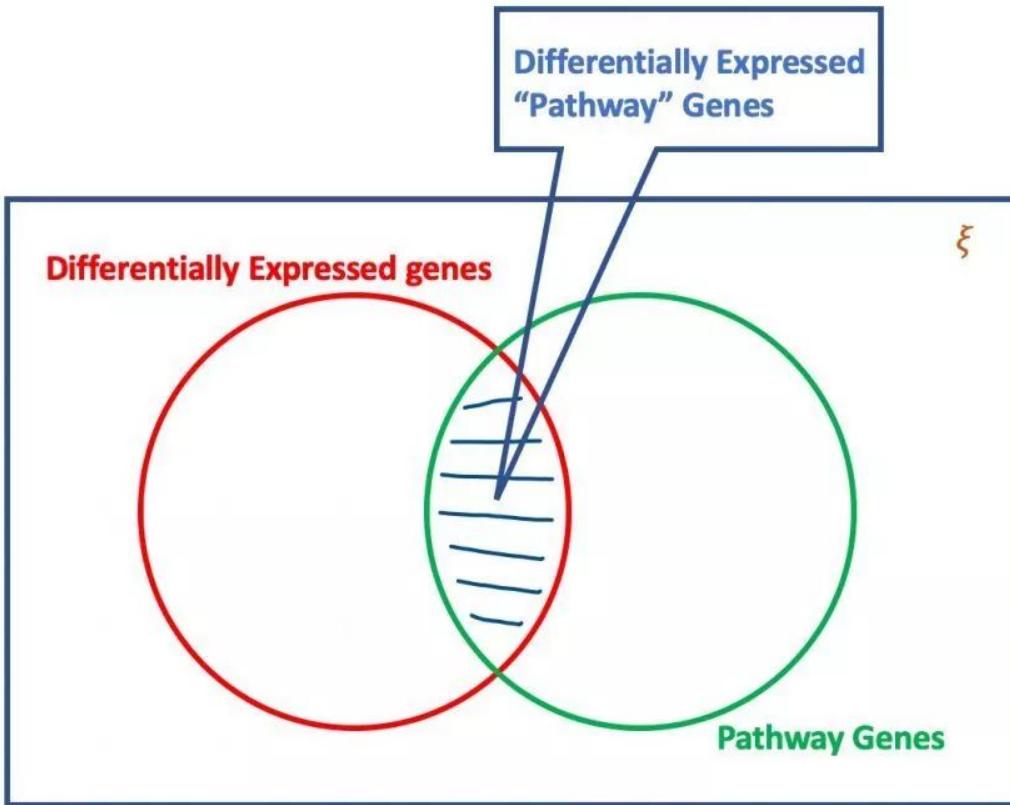


ORA

Statistically evaluates the fraction of genes in a particular pathway found among the set of genes showing changes in expression.

1. Select a list of genes with certain threshold (FDR ≤ 0.05)
2. For each pathway, count input genes that are part of the pathway
3. Repeat for an appropriate background list of genes
4. Every pathway is tested for over- or under-representation in the list of input genes

The most commonly used tests are based on the hypergeometric, chi-square, or binomial distribution



Problems with ORA

Cutoff? 0.051?

Treat all genes equally

Each gene is independent of other

Each pathway is independent of each other

FCS

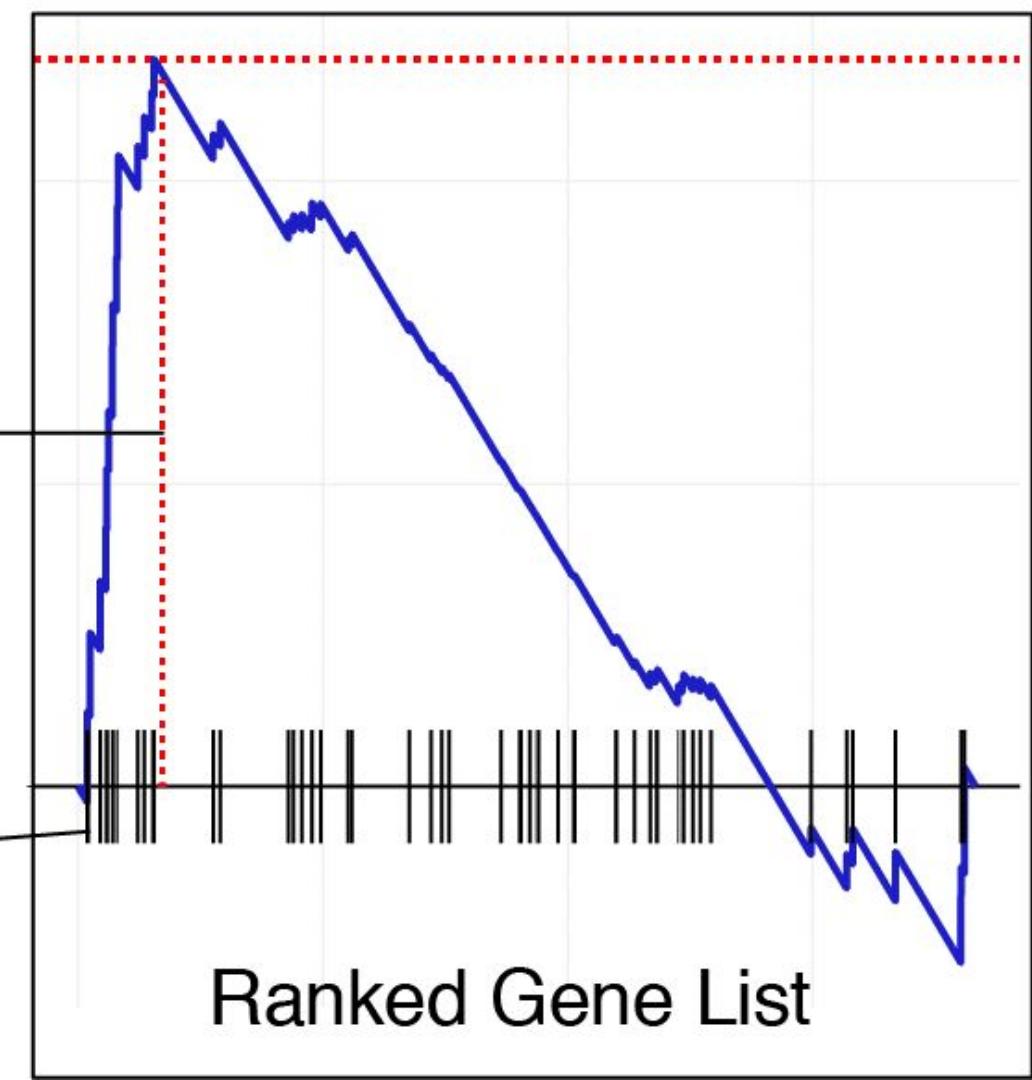
The hypothesis of functional class scoring (FCS) is that although large changes in individual genes can have significant effects on pathways, weaker but coordinated changes in sets of functionally related genes (i.e., pathways) can also have significant effects

1. Rank the genes
2. Perform gene-level statistics in a pathway
3. Calculate pathway level-statistics: – Kolmogorov-Smirnov statistic

Running Score

Enrichment Score

Genes in GO term



Problems with FCS

Each gene is independent of other

Each pathway is independent of each other

R package: fGSEA

Databases

- GO: BP, MF, CC
- KEGG
- Reactome
- DOSE
- DisGeNET
- MSigDb
- KEGG module
- WikiPathways
- TF
- miRNA
- "user input"
- PathGuide

Methods

- ORA
- GSEA
- SAFE
- PADOG
- ROAST
- CAMERA
- GSA
- GSVA/ssGSEA
- GlobalTest
- EBM
- MGSA
- GOSeq
- QUSAGE
- Pathview
- GOSemSim
- GGEA
- SPIA
- PathNet
- DEGraph
- TopologyGSA
- GANPA
- CePa
- NetGSA
- WGCNA

<https://yulab-smu.top/biomedical-knowledge-mining-book/index.html>

Problems with the databases?

Low resolution

Reactome pathways

\$`R-ATH-9675824.1: DNA replication Initiation`

"AT1G80190" "AT3G25100" "AT5G44635"
"AT5G49010" "AT2G07690" "AT3G12530"
"AT1G44900" "AT5G46280" "AT5G67100"
"AT5G22110" "AT3G55490" "AT1G19080"
"AT4G02060" "AT1G67630" "AT2G16440"
"AT1G67320"

\$`R-ATH-9675885.1: Lagging strand synthesis`

"AT1G07370" "AT4G18590" "AT5G67100"
"AT3G52630" "AT5G63960" "AT1G78650"
"AT2G29570" "AT1G67630" "AT2G42120"
"AT1G67320"

Make your own database?

database_seeds

\$paper1_day1

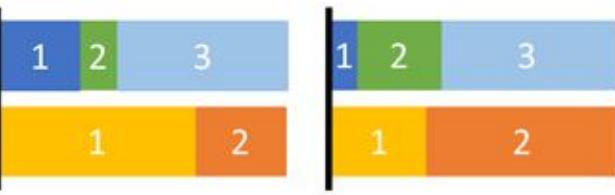
Gene1, Gene2, Gene3, Gene4

\$paper2_day2

Gene3, Gene4, Gene5, Gene6

RNA-seq data analysis pipeline

Additional analysis

	<u>Condition 1</u>	<u>Condition 2</u>	<u>Entity</u>	<u>Differential Expression</u>
[i]			Gene A	DGE ✓
			Gene B	DGE ✗
[ii]			Transcript A.1	DTE ✗
			Transcript A.2	DTE ✓
[iii]			Transcript A.3	DTE ✓
			Transcript B.1	DTE ✓
			Transcript B.2	DTE ✓
			Gene A	DTU ✓
			Transcript A.1	DTU ✓
			Transcript A.2	DTU ✓
			Transcript A.3	DTU ✗
			Gene B	DTU ✓
			Transcript B.1	DTU ✓
			Transcript B.2	DTU ✓

Gene expression

IsoformSwitchAnalyzeR: analysis of changes in genome-wide patterns of alternative splicing and its functional consequences

Kristoffer Vitting-Seerup  * and Albin Sandelin  *

The Bioinformatics Centre, Department of Biology and Biotech Research & Innovation Centre, University of Copenhagen, 2200 Copenhagen N, Denmark

*To whom correspondence should be addressed.

Associate Editor: Bonnie Berger

Received on August 28, 2018; revised on December 11, 2018; editorial decision on March 16, 2019; accepted on April 9, 2019

METHOD

Open Access

BANDITS: Bayesian differential splicing accounting for sample-to-sample variability and mapping uncertainty



Simone Tiberi*  and Mark D. Robinson

METHODOLOGY ARTICLE

Open Access

Streamlining differential exon and 3' UTR usage with diffUTR



Stefan Gerber^{1,2}, Gerhard Schratt² and Pierre-Luc Germain^{1,3,4*} 

Software | [Open Access](#) | [Published: 29 December 2008](#)

WGCNA: an R package for weighted correlation network analysis

[Peter Langfelder](#) & [Steve Horvath](#) 

[BMC Bioinformatics](#) **9**, Article number: 559 (2008) | [Cite this article](#)

366k Accesses | **10634** Citations | **82** Altmetric | [Metrics](#)

Methodology article | [Open Access](#) | Published: 27 February 2019

A high-throughput SNP discovery strategy for RNA-seq data

[Yun Zhao](#), [Ke Wang](#), [Wen-li Wang](#), [Ting-ting Yin](#), [Wei-qi Dong](#) & [Chang-jie Xu](#) 

[BMC Genomics](#) **20**, Article number: 160 (2019) | [Cite this article](#)

17k Accesses | **33** Citations | **3** Altmetric | [Metrics](#)

Learning outcomes

- QC, alignment, normalisation, differential analysis, functional analysis
- Interpretation of data
- Installing tools and packages
- Interpretation of plots
- Skimming through method papers
- Makefile
- Figures: heatmap, PCA plot, dot-plot

**Thank you for your
attention!**

Feedback

<https://bit.ly/bio634>

