# Overview of RNA-Seq

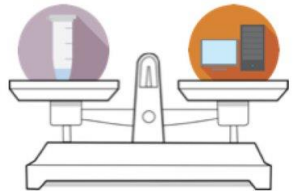Dr. Princess Rodriguez

2025-01-29

# Learning Objectives:

- Understand applications of RNA sequencing
- Introduce the overall differential expression workflow
- Understand experimental design concepts such as replicates and batch effects
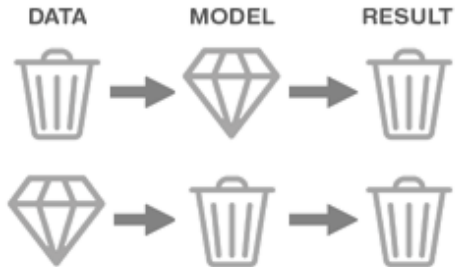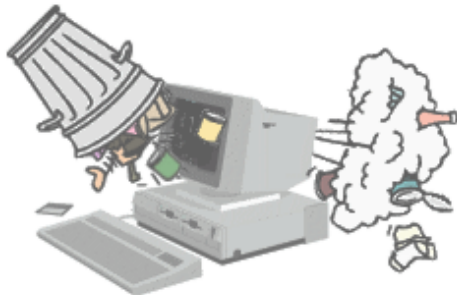- Understand different types of library preps, their requirements and uses.

"The quality of your data is at least directly proportional to the quality of your specimen."

**David B. Williams**

Transmission Electron Microscopy: A Textbook for Materials Science

# Garbage In, Garbage Out



DATA → MODEL → RESULT

Your input will define the quality of output you get!

# Overview of RNA-seq

RNA-seq is an exciting experimental technique that is utilized to explore and/or quantify gene expression within or between conditions.
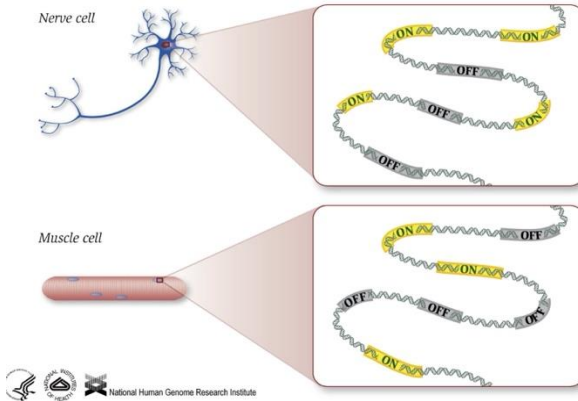


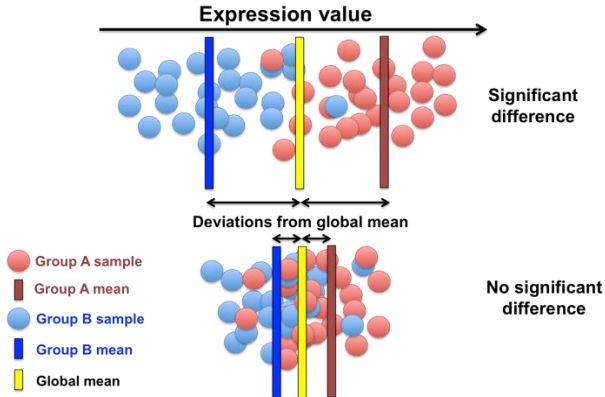**Figure 1:** Gene Expression in Cells

# The Transcriptome

The transcriptome is defined as a collection of all the transcript readouts present in a cell. RNA-seq data can be used to explore and/or <u>quantify</u> the transcriptome of an organism, which can be utilized for the following types of experiments:

- **Differential Gene Expression**: *quantitative* evaluation and comparison of transcript levels between conditions
- **Transcriptome assembly**: building the profile of transcribed regions of the genome, a *qualitative* evaluation
- **Refinement of gene models**: building better gene models and verifying them using transcriptome assembly
- **Metatranscriptomics**: community transcriptome analysis

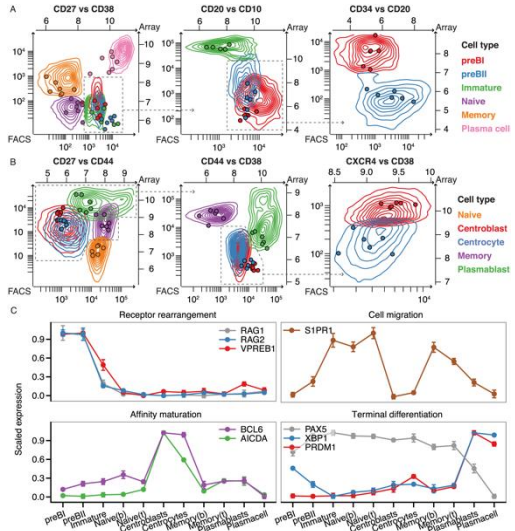# Basic types of questions answered:

What genes are differentially expressed between conditions?

## Other questions answered:

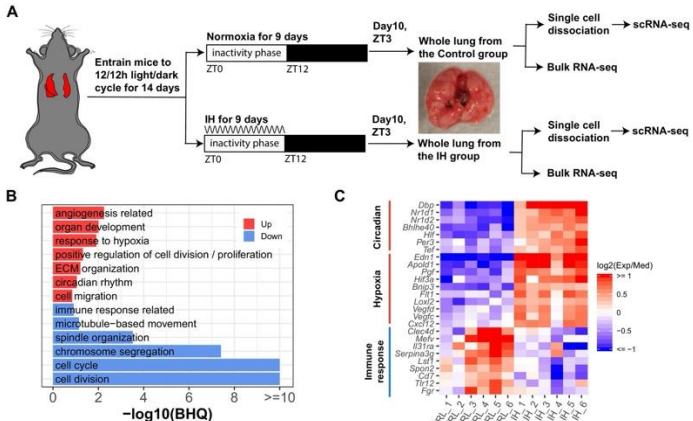Are there any trends in gene expression across development?

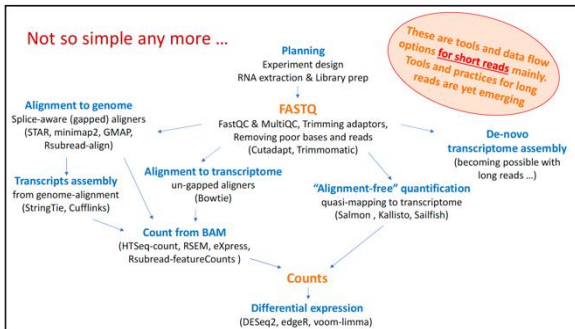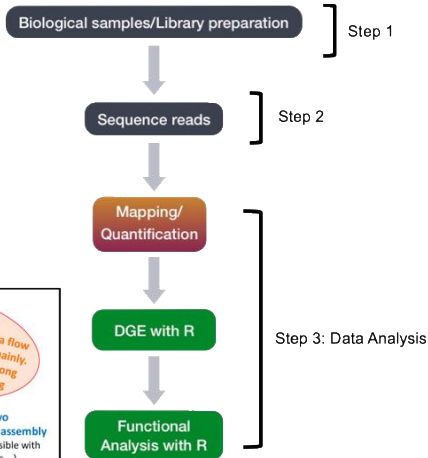Which groups of genes change similarly over time or across conditions?



https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0138236

# Basic types of questions answered:

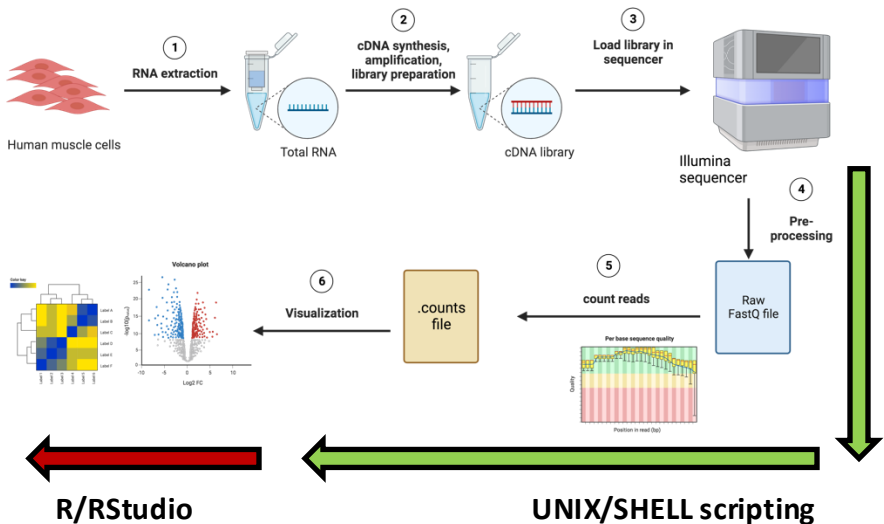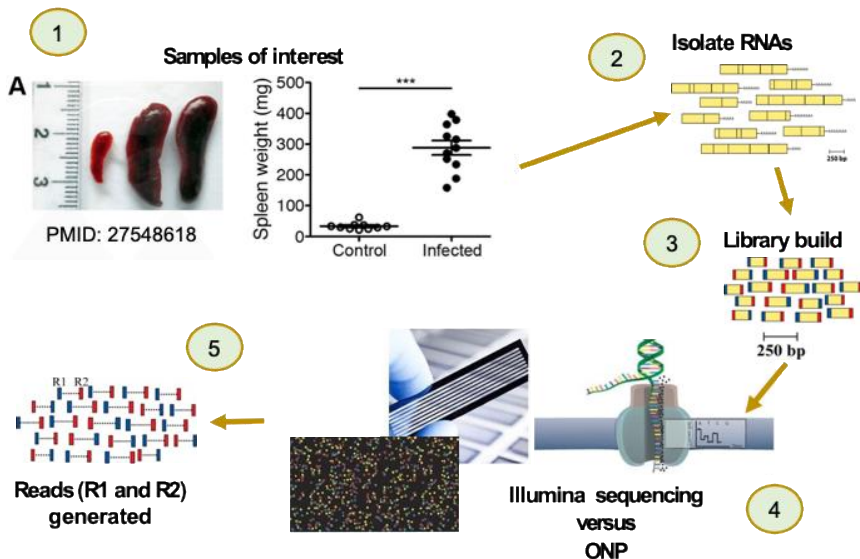What processes or pathways are enriched in condition of interest?

## Basic Principals

- Study Design
- Quality Assessment (UNIX)
- Trimming & Preprocessing (UNIX)
- Alignment (UNIX)
- Visualization of BAMs/counts (R)

# RNA-Seq Bioinformatic Pipeline

# Experimental workflow



**1** Samples of interest

PMID: 27548618

**2** Isolate RNAs

250 bp

**3** Library build

250 bp

**4** Illumina sequencing versus ONP

**5** Reads (R1 and R2) generated

R1   R2

## Biological Replicates

Experimental replicates can be performed as **technical replicates** or **biological replicates**.
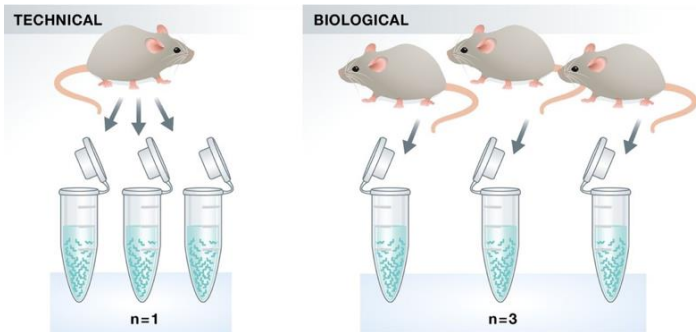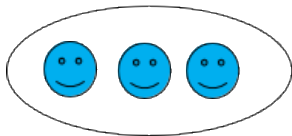
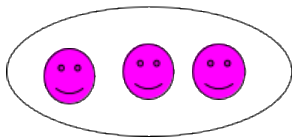

**Figure 16:** Biological Replicates

*Image credit: Klaus B., EMBO J (2015) **34**: 2727-2730*

- **Technical replicates:** use the same biological sample to repeat the technical or experimental steps in order to accurately measure technical variation and remove it during analysis.
- **Biological replicates** use different biological samples of the same condition to measure the biological variation between samples.

# Biological Replicates



Condition 1



Condition 2

❖ To detect Differentially Expressed Genes (DEGs) between groups we should have several samples, which are also known as biological replicates
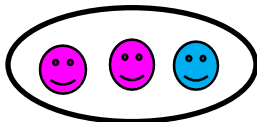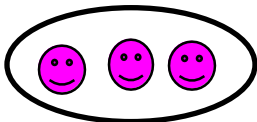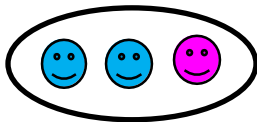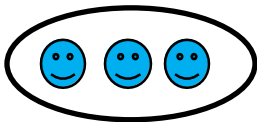
# Probability of detecting DEGs

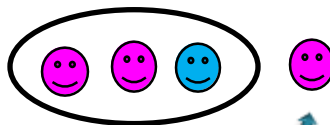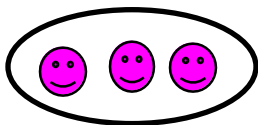|  | Replicates per group | | |
|---|---|---|---|
|  | 3 | 5 | 10 |
| Fold change | | | |
| 2 | 87% | 98% | 100% |

PMID: 26813401

# Grouping of Replicates
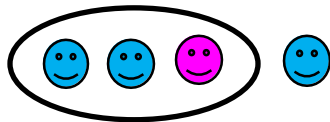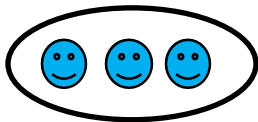


What you want

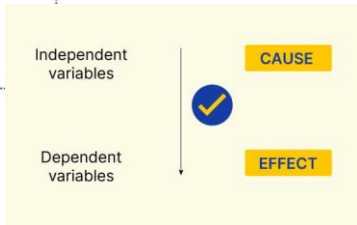What you get

# Grouping of Replicates



What you want

What you get

That spare comes in handy
Highly recommend especially
with mice!

# What causes this?
## Confounding variables

A variable that influences or *confounds* the relationship between an independent and dependent variable

# Examples of confounding variables



A new technician is
running the sequencer

# Examples of confounding variables



16S rRNA gene:

Sample collection

DNA (and RNA) isolation

16S rRNA gene PCR amplification of variable 4 region

Amplicon pooling and Illumina sequencing

Data analysis

5 counts of *L. iners*
3 counts of *G. vaginalis*
2 counts of *A. vaginae*

**Extracting DNA/RNA with two different kits!**

# Examples of confounding variables



16S rRNA gene:

V4

Sample collection → DNA (and RNA) isolation → 16S rRNA gene PCR amplification of variable 4 region → Amplicon pooling and Illumina sequencing → Data analysis

5 counts of *L. iners*
3 counts of *G. vaginalis*
2 counts of *A. vaginae*

**Sequencing on multiple different types of platforms**

# Examples of confounding variables



16S rRNA gene:

Sample collection → DNA (and RNA) isolation → 16S rRNA gene PCR amplification of variable 4 region → Amplicon pooling and Illumina sequencing → Data analysis

5 counts of *L. iners*
3 counts of *G. vaginalis*
2 counts of *A. vaginae*

**Inappropriate multiplexing strategy**

# Multiplexing



Generate & pool
indexed cDNA libraries

Sequence pooled
libraries on a single
lane

*in silico*: Demultiplex
the data on index

sample1    sample2    sample3    sample4    sample5    sample6
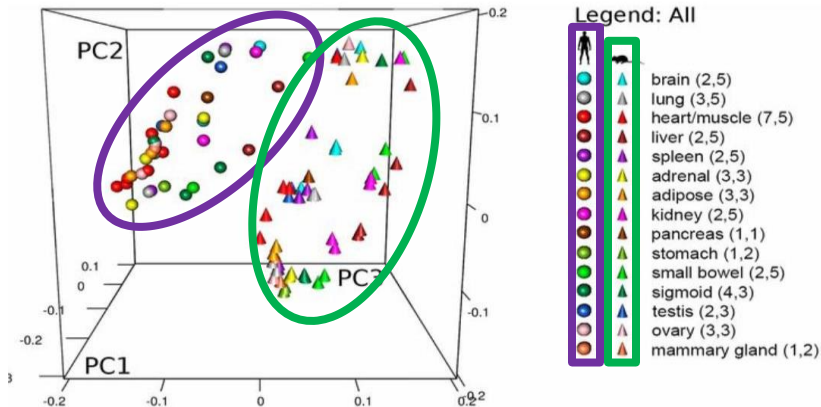
# ENCODE reported that gene expression was likely to follow a species-specific rather tissue-specific pattern



Legend: All

| | | |
|---|---|---|
| | | brain (2,5) |
| | | lung (3,5) |
| | | heart/muscle (7,5) |
| | | liver (2,5) |
| | | spleen (2,5) |
| | | adrenal (3,3) |
| | | adipose (3,3) |
| | | kidney (2,5) |
| | | pancreas (1,1) |
| | | stomach (1,2) |
| | | small bowel (2,5) |
| | | sigmoid (4,3) |
| | | testis (2,3) |
| | | ovary (3,3) |
| | | mammary gland (1,2) |

# ENCODE reported that gene expression was likely to follow a species-specific rather tissue-specific pattern

Reanalysis of Mouse ENCODE data suggests mouse and human genes are expressed in tissue-specific, rather than species-specific, patterns.

*May 19, 2015*
JYOTI MADHUSOODANAN

Late last year, members of the Mouse ENCODE consortium reported in *PNAS* that, across a wide range of tissues, gene expression was more likely to follow a species-specific rather than tissue-specific pattern. For example, genes in the mouse heart were expressed in a pattern more similar to that of other mouse tissues, such as the brain or liver, than the human heart.

WIKIMEDIA, RAMA

But earlier this month, Yoav Gilad of the University of Chicago called these results into question on Twitter. With a dozen or so 140-character dispatches (including three heat maps), Gilad suggested the results published in *PNAS* were an anomaly—a result of how the tissue samples were sequenced in different batches. If this "batch effect" was eliminated, he proposed, mouse and human tissues clustered in a tissue-specific manner, confirming previous results rather than supporting the conclusions reported by the Mouse ENCODE team.
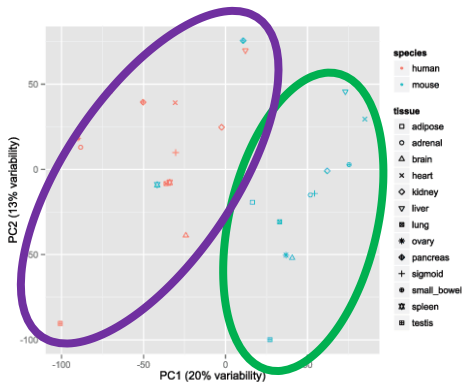
Sequence study design (sequencer ID, run ID, lane number):

| D87PMJN1 (run 253, lane 7) | D87PMJN1 (run 253, lane 8) | D4LHBFN1 (run 276, lane 4) | MONK (run 312, lane 6) | HWI-ST373 (run 375, lane 7) |
|---|---|---|---|---|
| heart | adipose | adipose | heart | brain |
| kidney | adrenal | adrenal | kidney | pancreas |
| liver | sigmoid colon | sigmoid colon | liver | brain |
| small bowel | lung | lung | small bowel | spleen |
| spleen | ovary | ovary | testis | ● human |
| testis | | pancreas | | ● mouse |

Sequencing lane (a batch effect) was almost completely confounded with species in the PNAS study. From @Y_Gilad
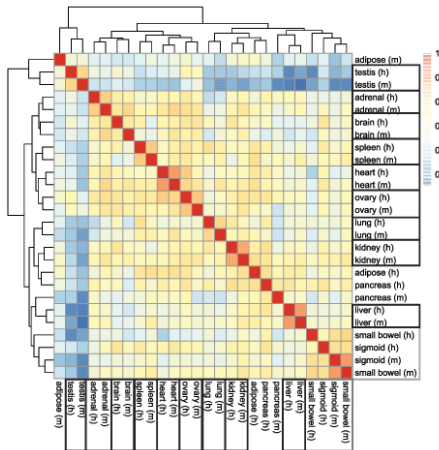
# Before accounting for batch effect
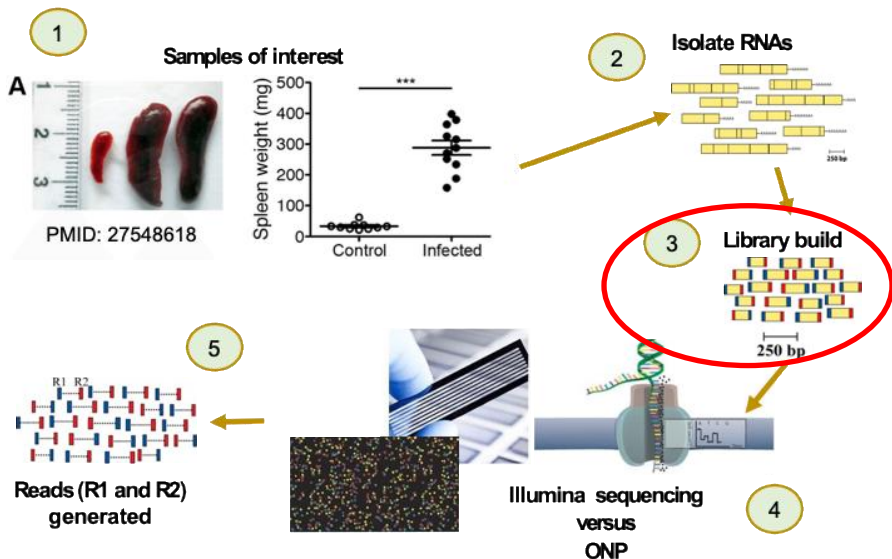


*Samples grouped by animal*

*Samples now grouped by tissue!*

# *What does this all means?*

- Its sometimes impossible for bioinformaticians to partition biological variation from technical variation, when these two sources of variation *are confounded*.

- No amount of statistical sophistication can separate confounded factors after data have been collected.

- *…these confounding variables may or may not be in your control!*

   **A well-planned experiment with an additional sample, does end up saving you time and money down the road. Its up to you to recognize this!**

# Experimental workflow



**1** Samples of interest

PMID: 27548618

**2** Isolate RNAs

250 bp

**3** Library build

250 bp

**4** Illumina sequencing
versus
ONP

**5** Reads (R1 and R2)
generated

# RNA composition

RNA comes in many
different flavors
- Ribosomal-related
  RNAs:
  - rRNA, tRNA,
    snoRNA (up to 90%
    of RNAs)
- Protein-coding RNAs:
  - mRNA
- Regulatory RNAs:
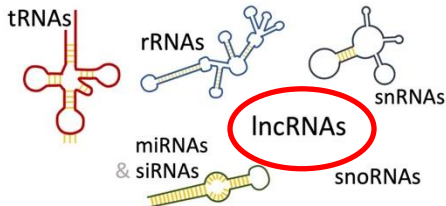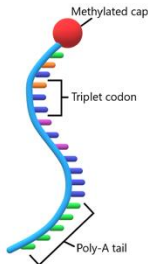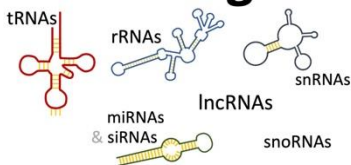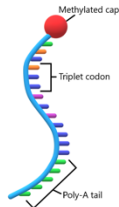  - microRNAs,
    lncRNAs

**"Comprehensive" transcriptome analysis**

*Two different protocols/kits!*

# RNA composition

RNA comes in many different flavors
- Ribosomal-related RNAs:
  - rRNA, tRNA, snoRNA (up to 90% of RNAs)
- Protein-coding RNAs:
  - mRNA
- Regulatory RNAs:
  - microRNAs, lncRNAs



tRNAs

rRNAs

snRNAs

lncRNAs

miRNAs & siRNAs

snoRNAs

messenger RNA

Methylated cap

Triplet codon

Poly-A tail

**The RNA sample undergoes either selection of the mRNA (polyA selection) or depletion of the rRNA. The resulting RNA is fragmented.**

# Poly-A versus rRNA depletion?

If you are aiming to obtain information about long non-coding RNA's I recommend performing ribosomal RNA depletion

Bacterial mRNAs are also not poly-adenylated

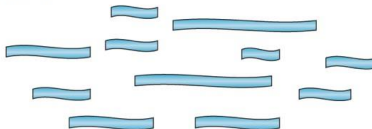# Illumina Library preparation


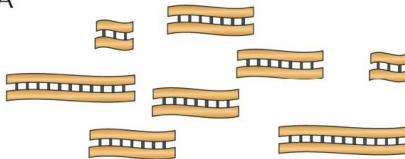
① mRNA or total RNA
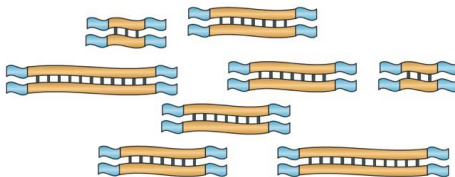
② Remove contaminant DNA

Remove rRNA?
Select mRNA?

③ Fragment RNA

④ Reverse transcribe into cDNA

⑤ Ligate sequence adaptors

# Strandedness

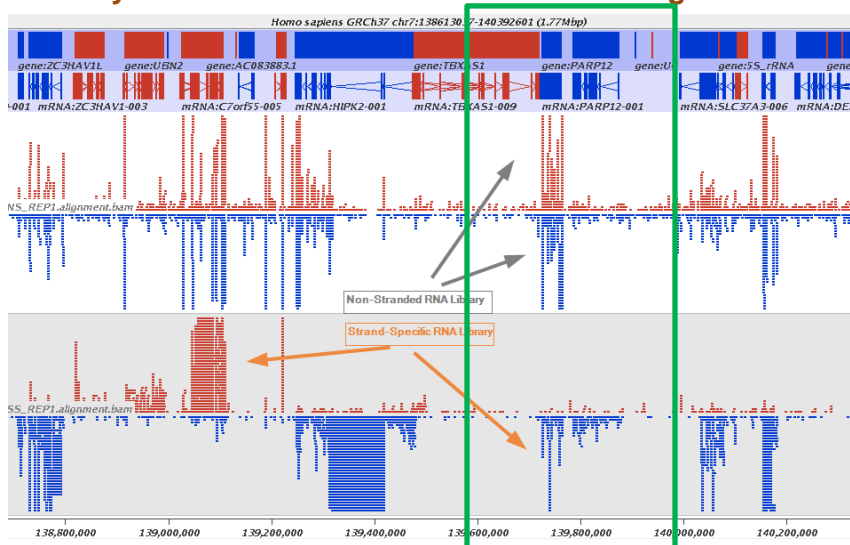Another consideration is whether to generate strand-preserving libraries

Libraries can be stranded or unstranded

The implication of **stranded** libraries is that you could distinguish whether the reads are derived from forward or reverse-encoded transcripts
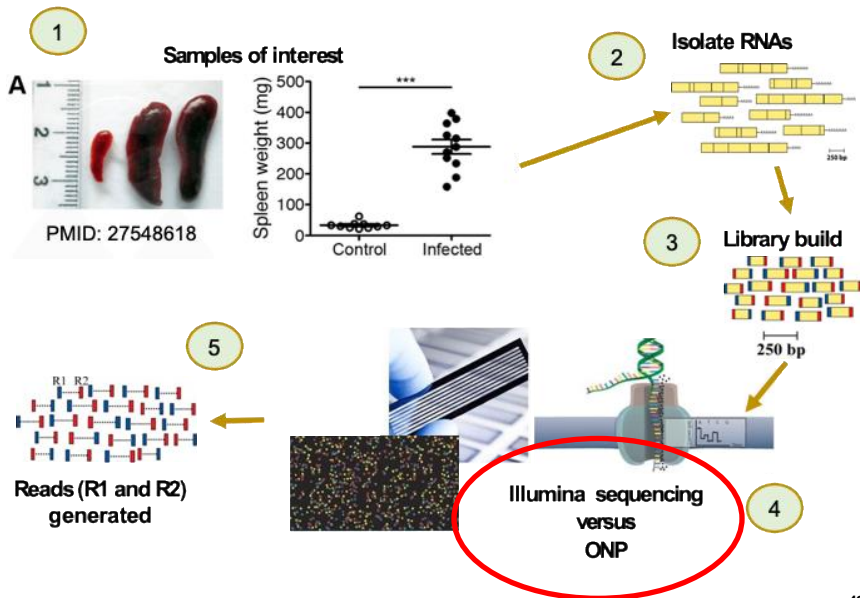
**Simple and accurate analysis of overlapping genes:**
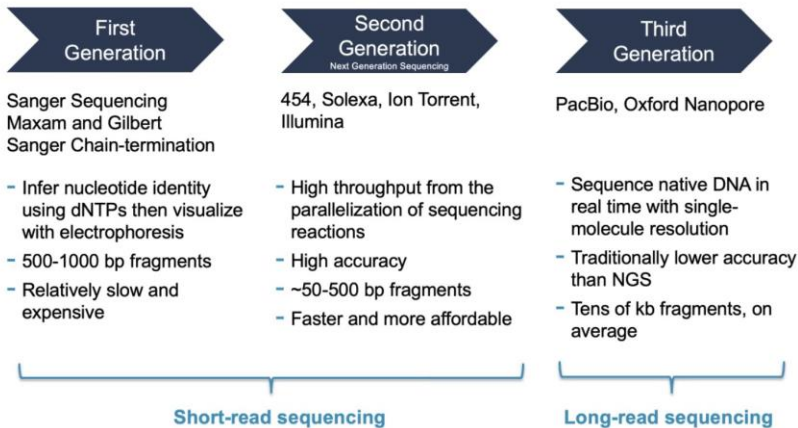**Clearly see that PARP12 is encoded on the negative strand**

Red = + strand    Blue = - strand

# Experimental workflow



1 — Samples of interest

PMID: 27548618

Spleen weight (mg) — Control, Infected (***)

2 — Isolate RNAs — 250 bp

3 — Library build — 250 bp

4 — Illumina sequencing versus ONP

5 — Reads (R1 and R2) generated — R1 R2

# Two main approaches in NGS: short-read vs long-read

**THE EVOLUTION OF SEQUENCING**

| First Generation | Second Generation Next Generation Sequencing | Third Generation |
|---|---|---|
| Sanger Sequencing Maxam and Gilbert Sanger Chain-termination | 454, Solexa, Ion Torrent, Illumina | PacBio, Oxford Nanopore |
| - Infer nucleotide identity using dNTPs then visualize with electrophoresis<br>- 500-1000 bp fragments<br>- Relatively slow and expensive | - High throughput from the parallelization of sequencing reactions<br>- High accuracy<br>- ~50-500 bp fragments<br>- Faster and more affordable | - Sequence native DNA in real time with single-molecule resolution<br>- Traditionally lower accuracy than NGS<br>- Tens of kb fragments, on average |

Short-read sequencing          Long-read sequencing

*The bioinformatic pipeline for these are different!*

## Single-end versus Paired-end

After preparation of the libraries, sequencing can be performed to generate the nucleotide sequences of the ends of the fragments, which are called **reads**. You will have the choice of sequencing a single end of the cDNA fragments (single-end reads) or both ends of the fragments (paired-end reads).



**Figure 10:** Paired End Reads

- SE => Only Read1 => one FASTQ file/sample
- PE => Read1 + Read2 => **two FASTQ files/sample**

# What is the Advantage of Longer and PE Reads?



> Reads mapping to junctions
> > With longer reads we will have more reads spanning exons
> > Isoforms or distinguishing paralogs

> Paired end reads

Knowing both ends of a fragment and an approximation of fragment size helps to determine the transcript from which it was derived.

# In Summary, to quantify Differential Gene Expression

- Technology: Illumina
- Read length: 50bp to 300 bp
- Paired vs single end: *doesn't matter but important to note*
- Number of reads: > 15 million per sample
- Replicates: 3 biological replicates *minimum*

*A well-planned experiment goes a long way!*

## Different sequencing platforms

There are a variety of Illumina platforms to choose from to sequence the cDNA libraries.



MiniSeq    MiSeq    NextSeq    HiSeq    HiSeq X

Benchtop    Production-Sca

**Final projects from the years have spanned the following topics:**


*Salmonella enterica*


Applications of organoids as research models


nf-core
https://nf-co.re

Deploy
- Stable pipelines
- Centralized configs
- List and update pipelines
- Download for offline use

Participate
- Documentation
- Slack workspace
- Twitter updates
- Hackathons

Develop
- Starter template
- Code guidelines
- CI code linting and tests
- Helper tools


Human microbiome
Archaea, bacteria, fungi and viruses

Dengue

Breast cancer

And more....!

*Original Published Work*

*Undergraduate student recreation*

**Green Trail**
UG credentials:
1-semester of intro bioinformatics

### Research Question

Due to the strong relationship between the kidney and the heart, which differentially expressed genes in bear kidneys are related to cardiac pathways?

**Black Trail**
UG credentials:
1-semester of intro bioinformatics

Research Question



Due to the strong relationship between the kidney and the heart, which differentially expressed genes in bear kidneys are related to cardiac pathways?

# Design



## How *Lactobacillus plantarum* shapes its transcriptome in response to contrasting habitats

"Aiming at elucidating how L. plantarum regulates and shapes its transcriptome in response to contrasting habitats."

Triplets from nine model media:
- A. mellifera L. worker bees
- D. melanogaster
- Human omnivore and vegan feces
- Table olives
- Tomato and pineapple juices
- Wheat flour hydrolysate
- Cheese broth.

Later cultivation on MRS broth with two reference strains: WCFS1 and LB16

**Green Trail**
UG credentials:
1-semester of intro
bioinformatics

## Citation

*This lesson has been developed by members of the teaching team at the [Harvard Chan Bioinformatics Core (HBC).](#) These are open access materials distributed under the terms of the [Creative Commons Attribution license](#) (CC BY 4.0), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.*
*Authors: Mary Piper, Meeta Mistry, Radhika Khetani*
*Other sources - [https://umich-brcf-bioinf.github.io/rnaseq_demystified_workshop/site/Module3a_Design_Prep_Seq#2_Experimental_Design_and_Practicalities](#)*