

# RNA-SEQ DATA ANALYSIS: TRANSCRIPTOME ASSEMBLY AND DIFFERENTIAL EXPRESSION ANALYSIS

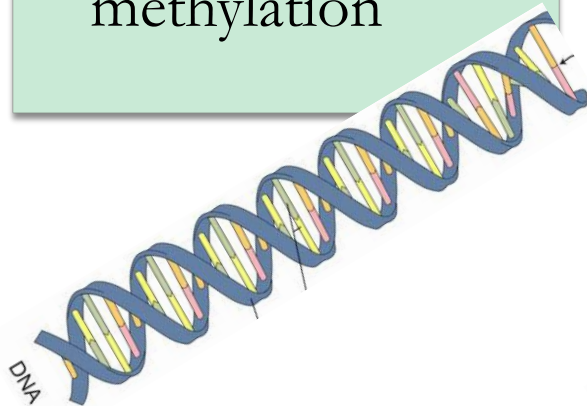
Ashley Sawle (Ashley.Sawle@cruk.cam.ac.uk)  
Guillermo Parada (guillermo.parada@sanger.ac.uk)

# HTS Applications - Overview

2

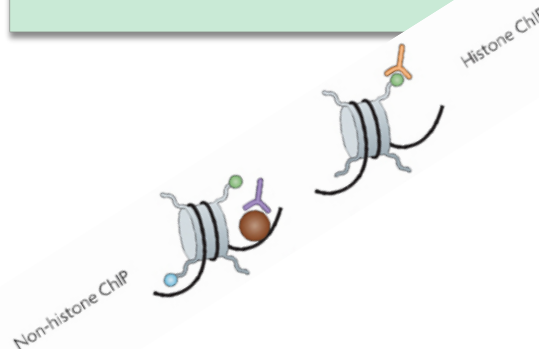
## DNA Sequencing

- Genome Assembly
- SNPs
- DNA methylation



## ChIP-sequencing

- Transcription Factor Binding Sites
- Chromatin Modification Regions



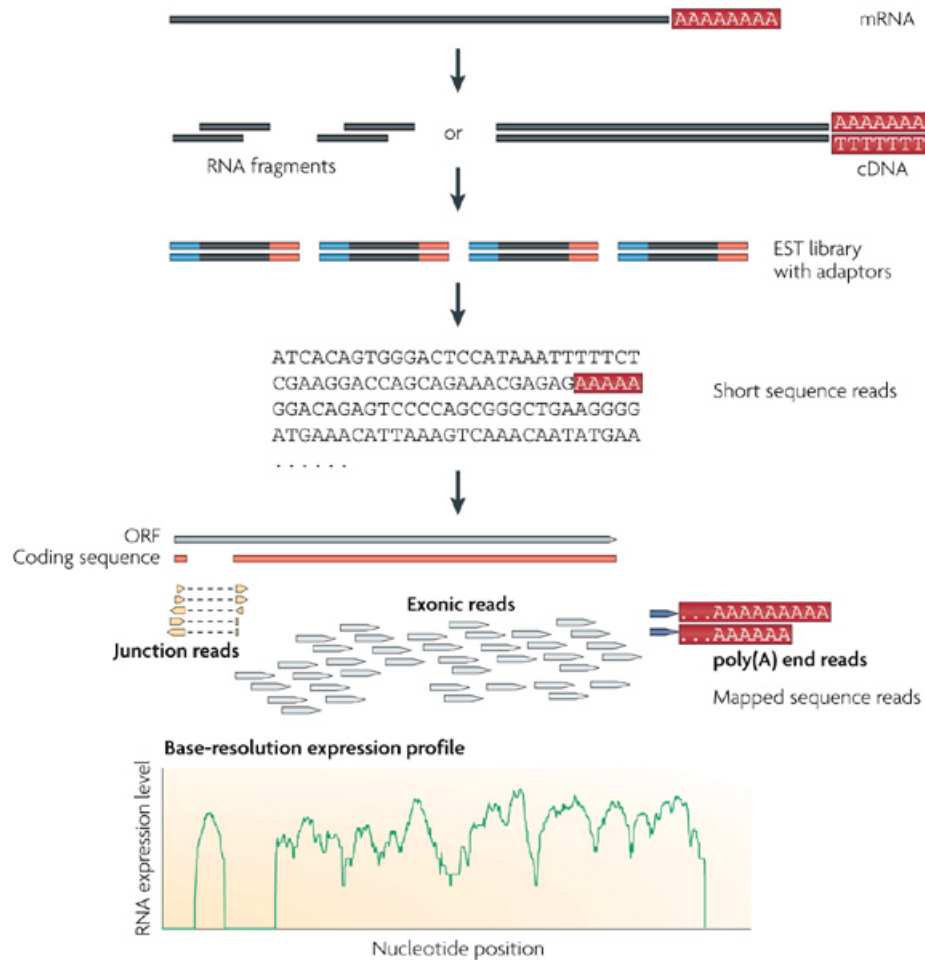
## RNA-sequencing

- Transcriptome Assembly
- Gene Expression
- Differential Expression



# RNA-seq workflow

3



Step1

Library Preparation

Step2

Sequencing

Step3

Bioinformatics Analysis

# Designing the right experiment

4

## 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.

Comic by Christine Ambrosino <http://www.hawaii.edu/fishlab/Nearside.htm>

# Designing the right experiment

5

- The design of the experiment is the first step and it is obviously determinant for all downstream analyses
- You have to evaluate all the eventualities and limitations of available technologies, designing the experiment according to your goals

# Designing the right experiment

6

## **COVERAGE: How many reads do we need?**

The coverage is defined as  $C = (R_{\text{length}} \times R_{\text{num}}) / A_{\text{length}}$

$R_{\text{length}}$  = length in nucleotides of the reads

$R_{\text{num}}$  = number of sequenced reads

$A_{\text{length}}$  = number of nucleotides of sequenced subject (genome, transcriptome, exome)

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

# Designing the right experiment

7

## READ LENGTH: long or short reads?

The answer depends again on the experiment:

GENOME RESEQUENCING

De novo TRANSCRIPTOME

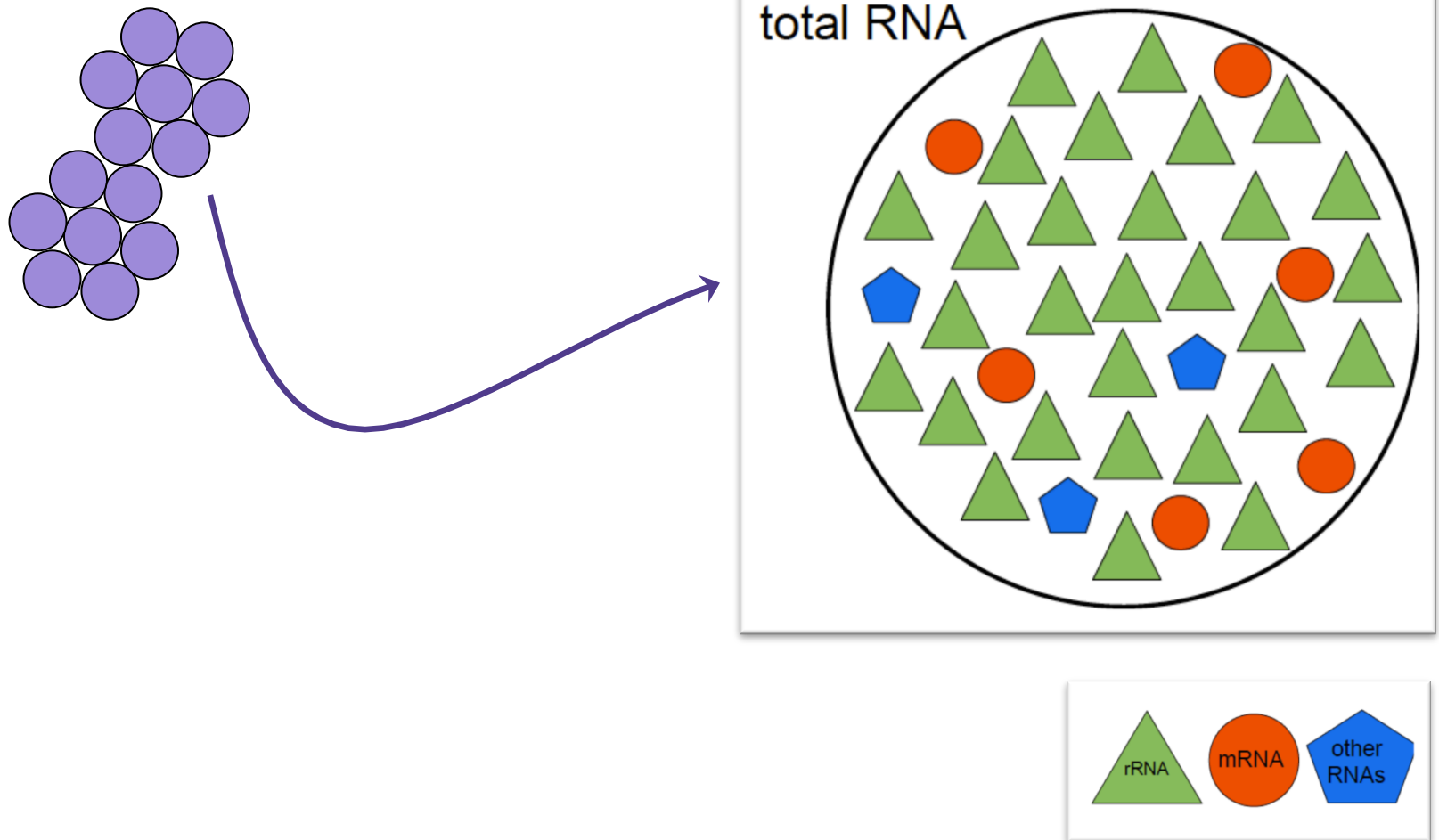
TRANSCRIPTOME seq

ChIP seq

Read length is inversely proportional to the multi-mappability of a read, in a sample of 50 nt reads there is a small fraction ( $<0.01\%$ ) that can be mapped to multiple positions of the human genome.

# Step1 – Library Preparation

8

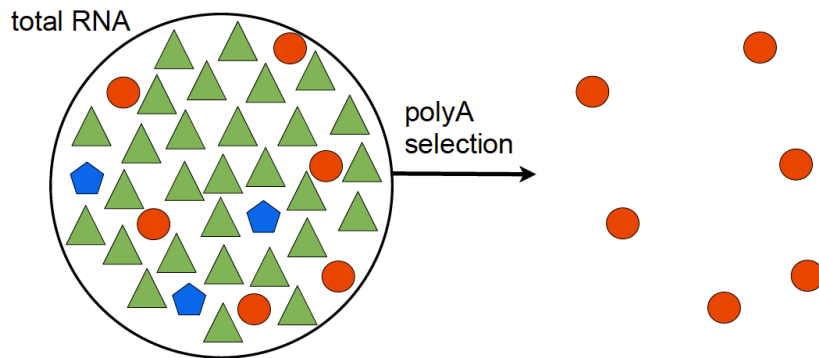




# Step 1 – Library Preparation

9

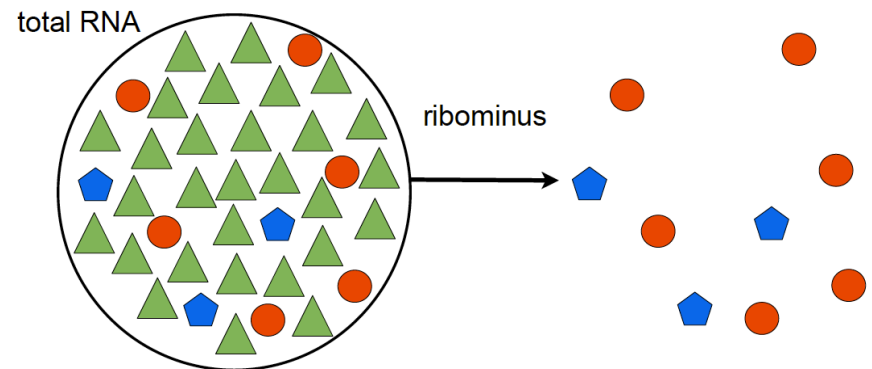
## polyA selection



poly(A<sup>+</sup>)-transcripts:

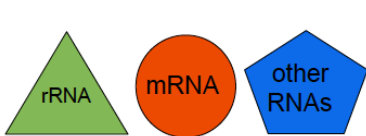
- mRNAs
- immature microRNAs
- snoRNAs

## ribominus selection



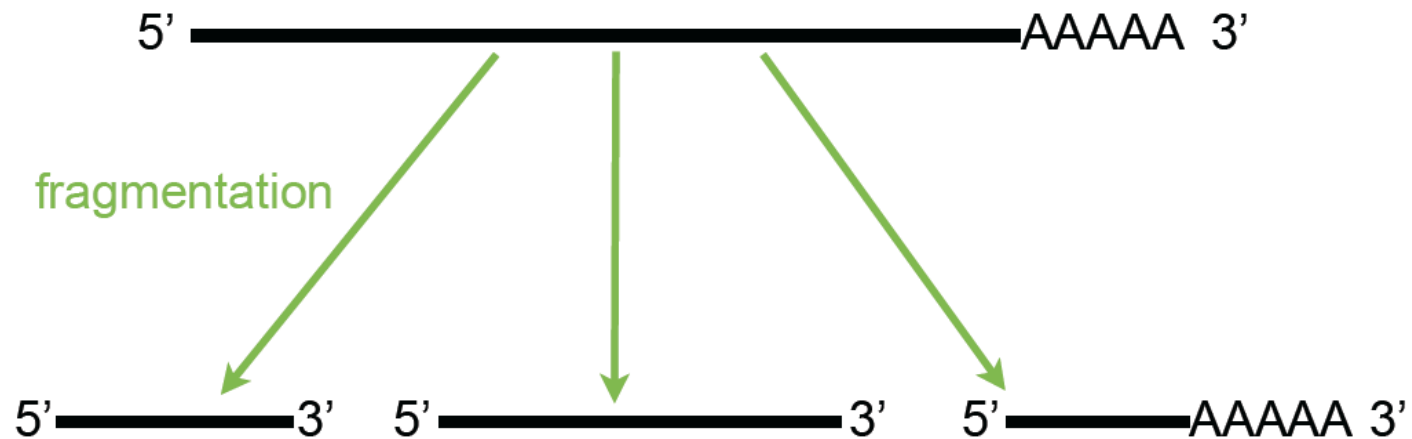
non poly(A<sup>+</sup>)-transcripts:

- mRNAs
- histone mRNAs
- tRNAs
- other small RNAs



# Step 1 – Library Preparation

10



# Step 1 – Library Preparation

11

## 1 strand cDNA synthesis



## remove RNA strand



# Step 1 – Library Preparation

12

## 2nd strand cDNA synthesis



Unstranded protocol

# Step 1 – Library Preparation

13

## 2nd strand cDNA synthesis



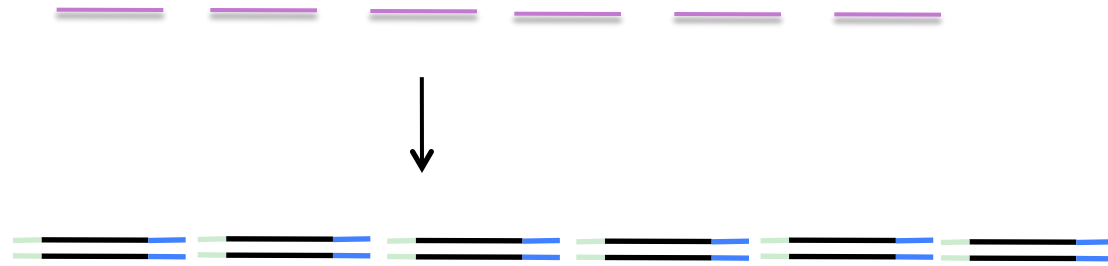
stranded protocol

# Step 1 – Library Preparation

14

Fragmented  
cDNA

cDNA with  
adaptors

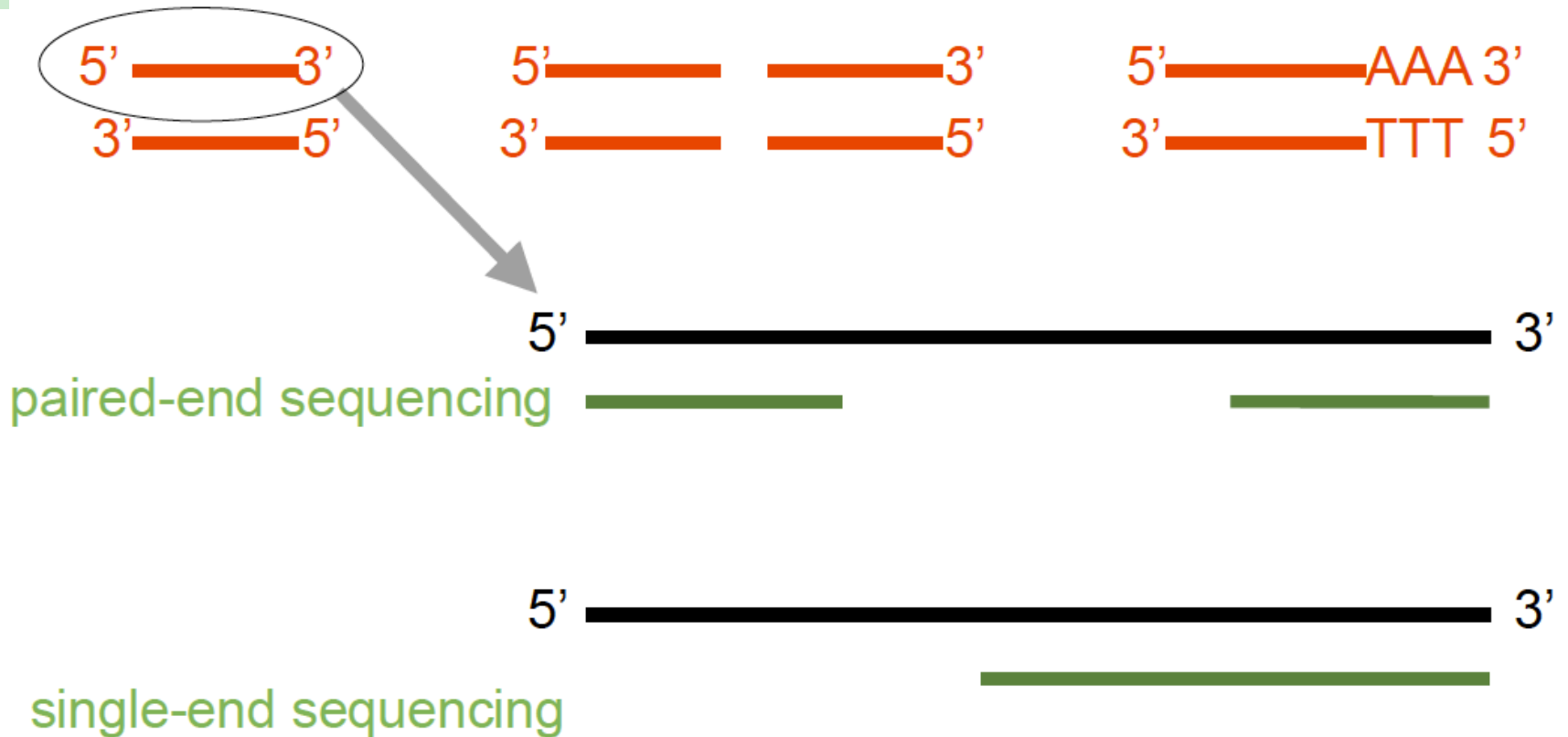


Size selection

PCR amplification

# Step 2 – Sequencing

15



# Single- vs paired-end sequencing

16

my\_sequence.fastq

```
@HWI-BRUNOP16X_0001:1:1:1466:1018#0/1
AAGGAAGTGCTTGTCTGGCTAACACAGCNAGNCACGTGAC
+
aVfbe`^^^_TTTSSdffffdfffabbZbbfebafbbbbbb
```

SE

my\_sequence\_1.fastq

```
@HWI-BRUNOP16X_0001:1:1:1278:989#0/1
NAAATTTTGAATTTCTGTGAAGTAAGCATCTTCTTTGTCAT
+
BJJGGKIINN^^^^QQNTUQ00TTTTRTOTY^^Y^^\^^^\
```

my\_sequence\_2.fastq

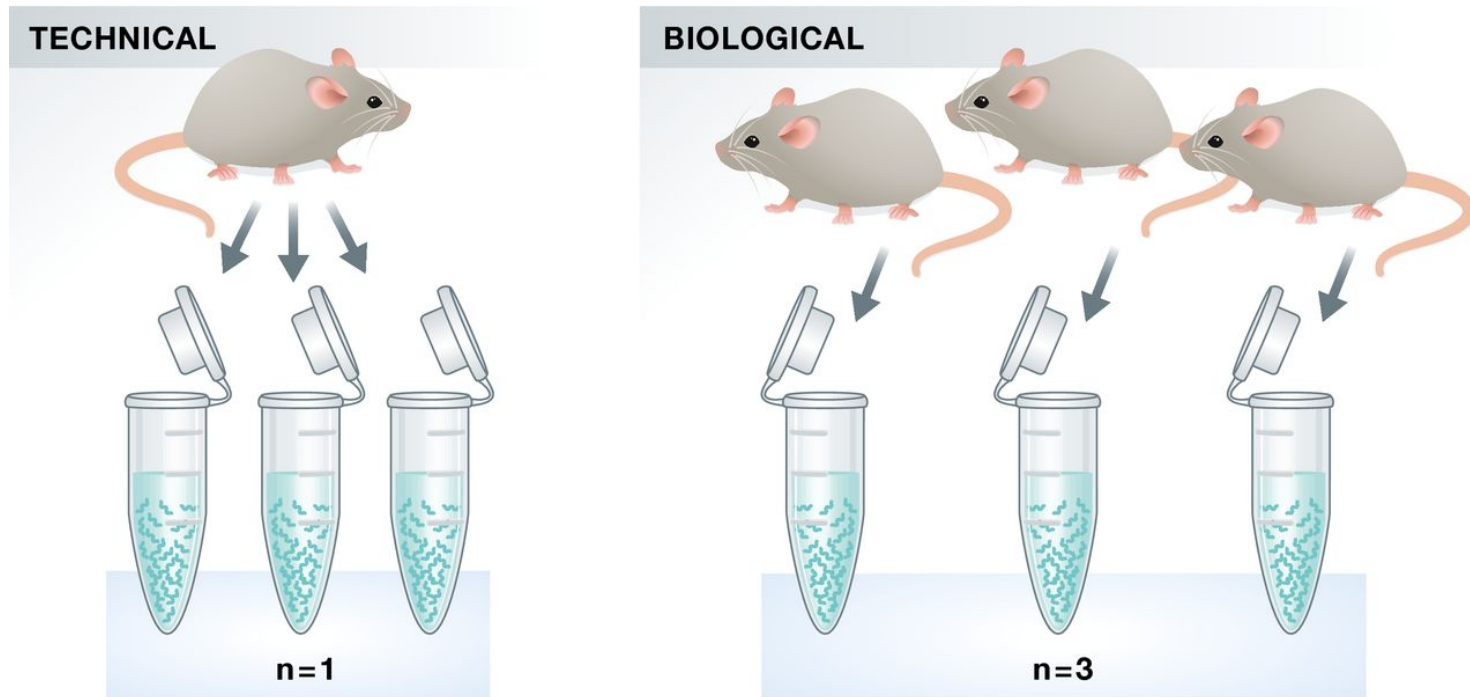
```
@HWI-BRUNOP16X_0001:1:1:1278:989#0/2
AACCCACACAGGAGAGCAGCCTTACAGATGCAAATACTGTG
+
]K____fffffggghgegghggggggdgggggfggggggegghh
```

PE



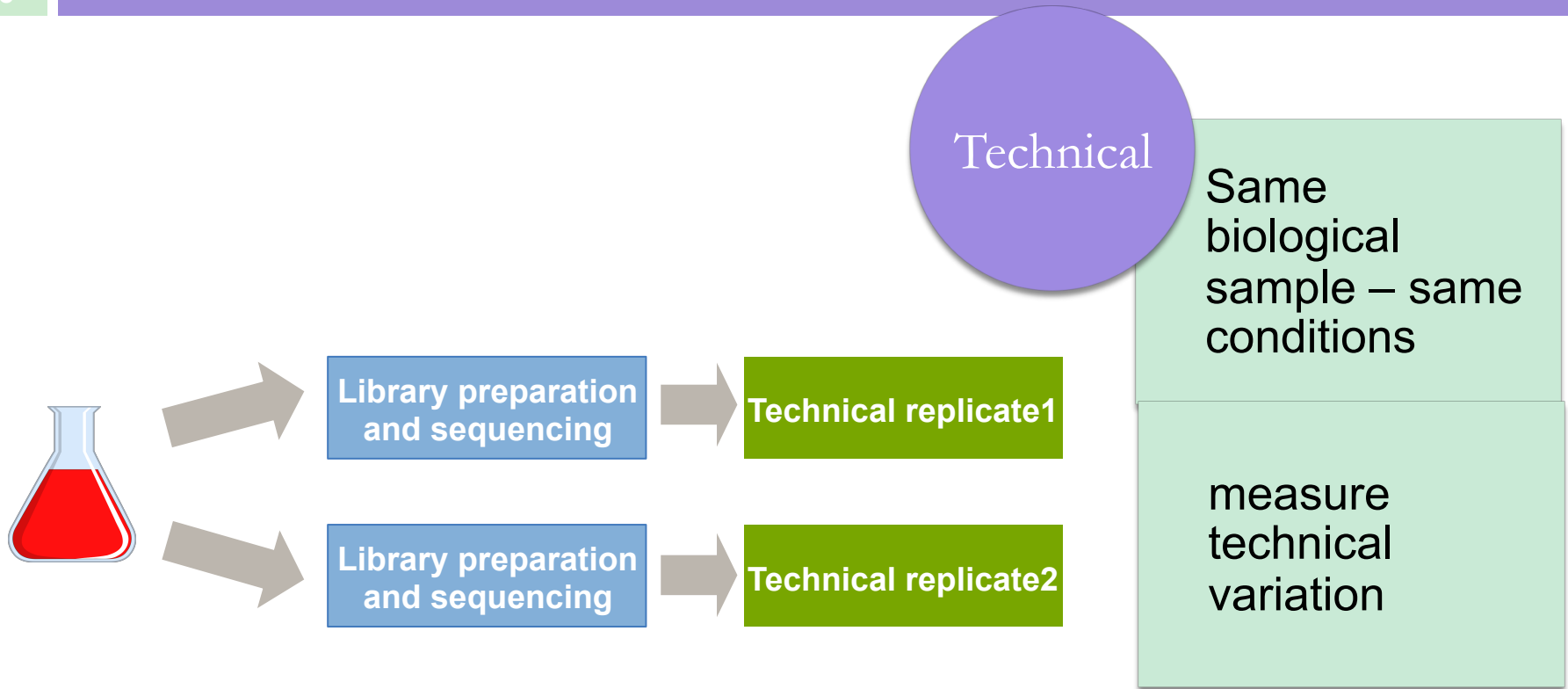
# Replicates – do I need them?

17



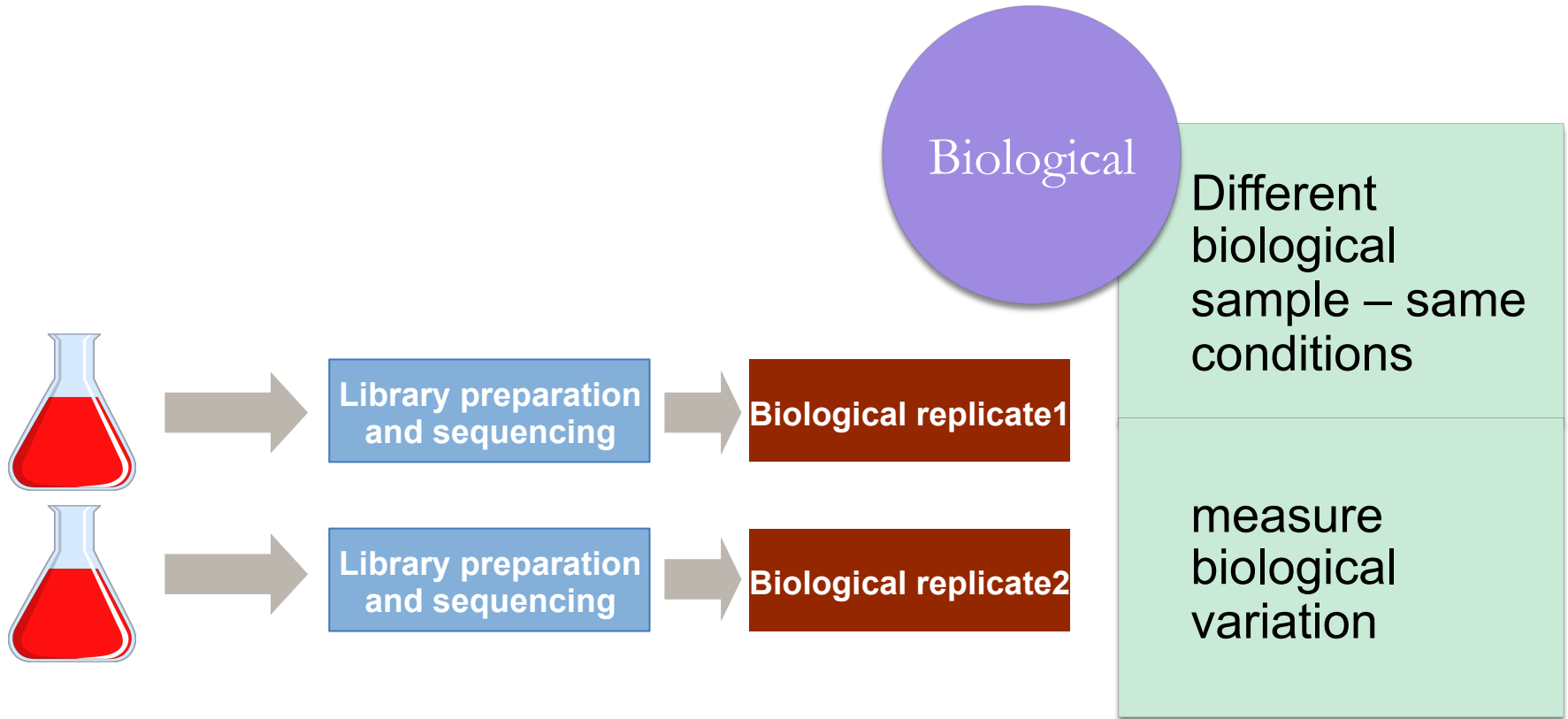
# Replicates – do I need them?

18



# Replicates – do I need them?

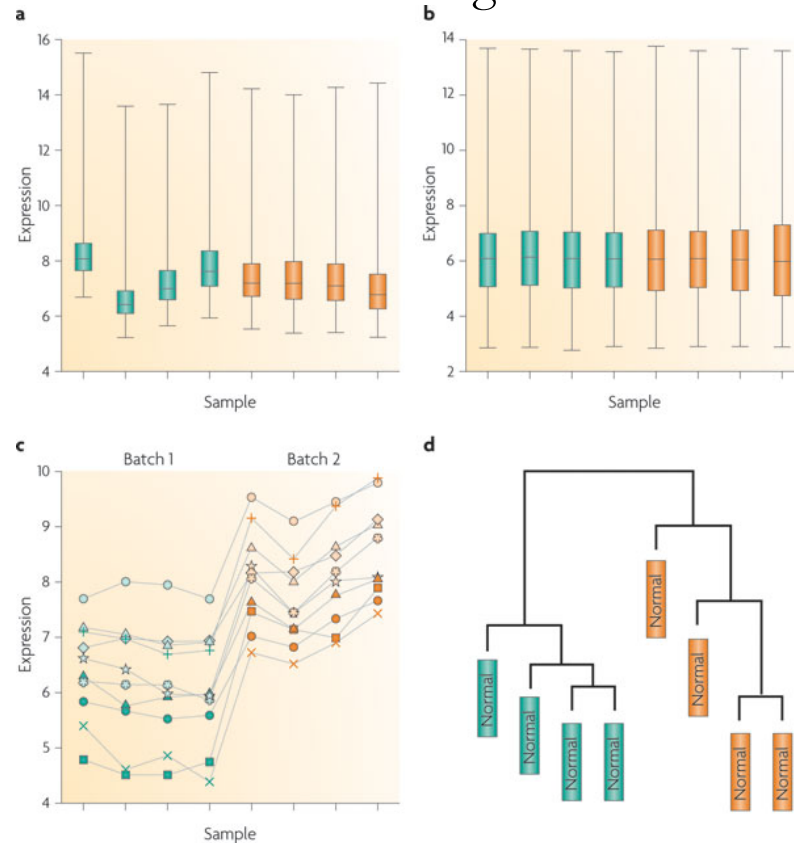
19



# Controlling batch effects

20

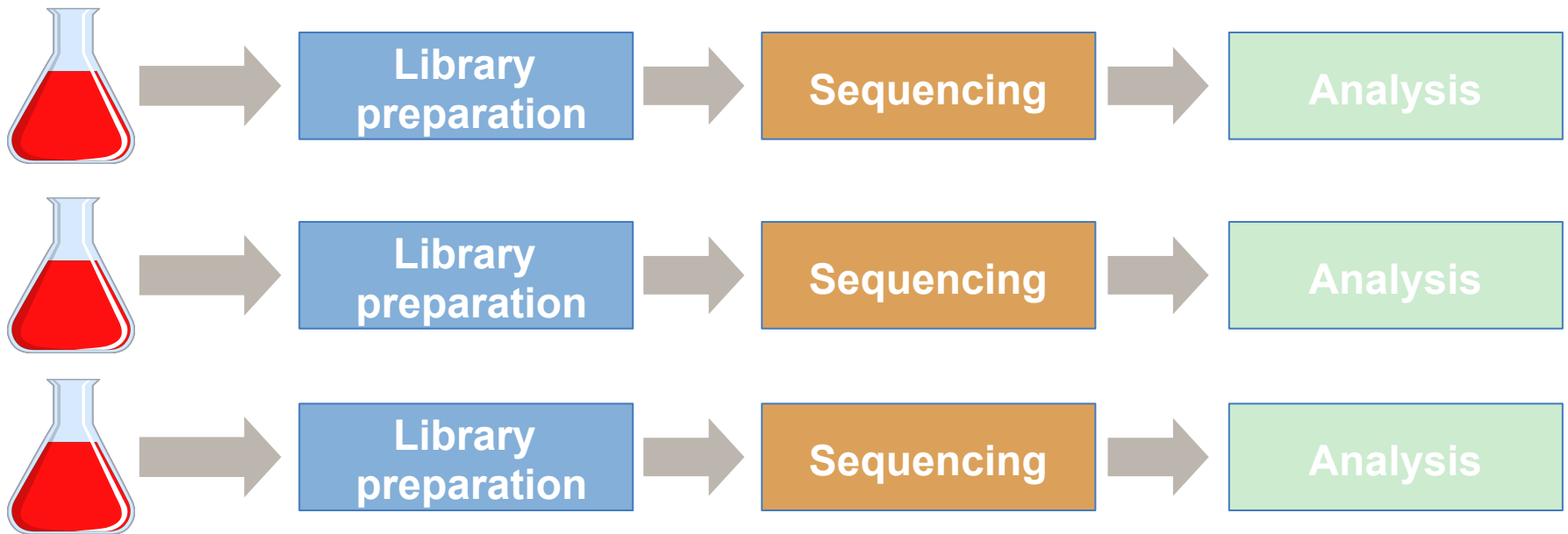
Batch effects are sub-groups of measurements that have qualitatively different behavior across conditions and are unrelated to the biological or scientific variables in a study



Nature Reviews | Genetics

# Controlling batch effects

21



# Example of experimental design

22



Group A



Group B

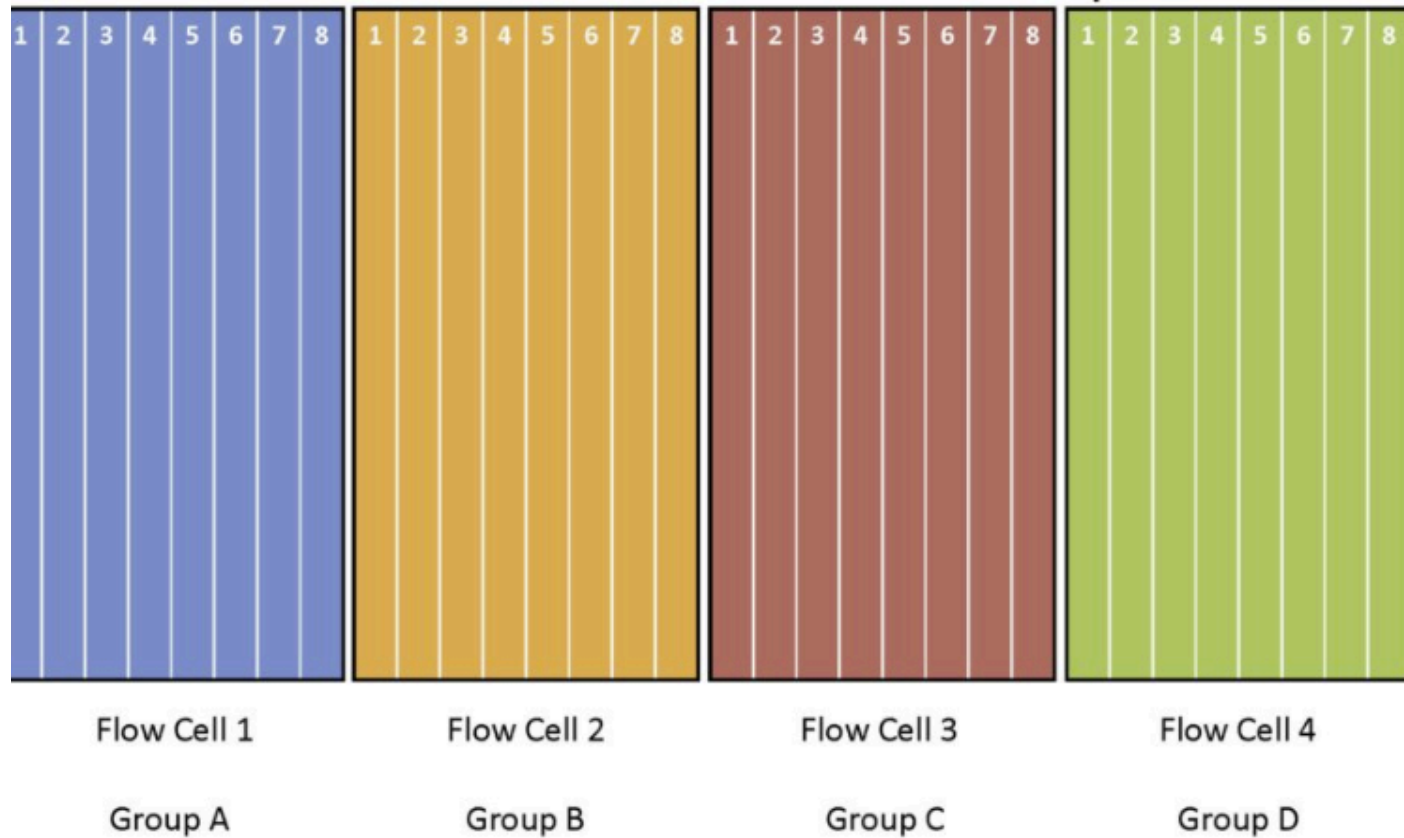


Group C



Group D

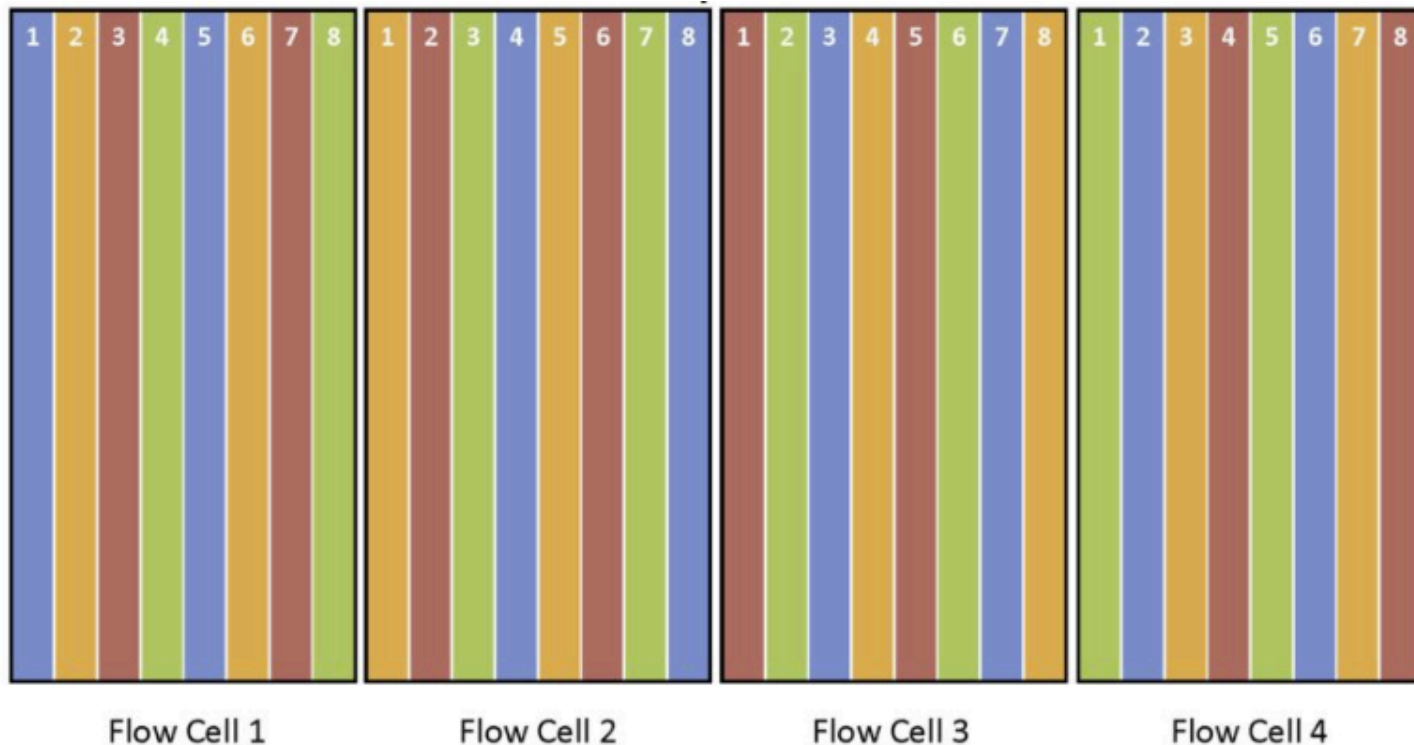
# Example of experimental design



# ...better experimental design

24

- Randomize samples with respect to the flow cell

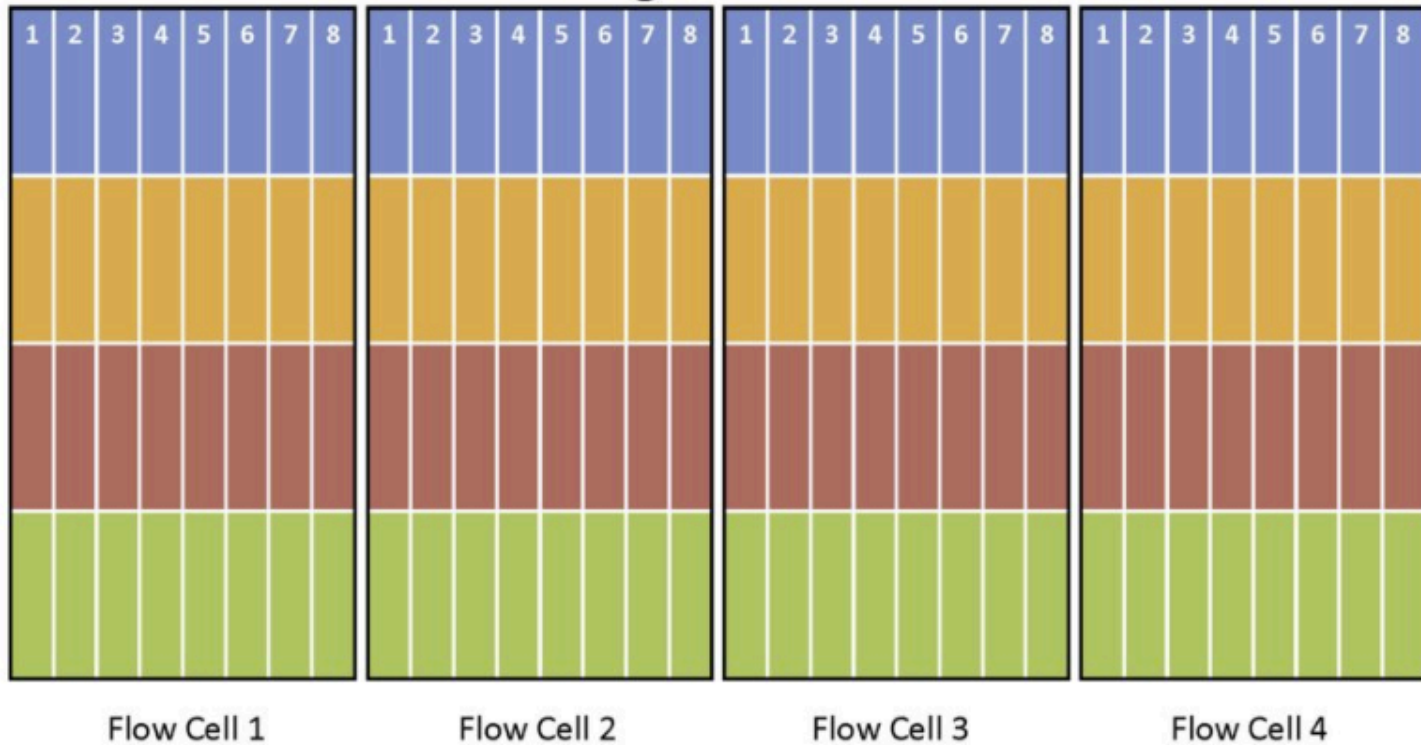




# ...even better experimental design

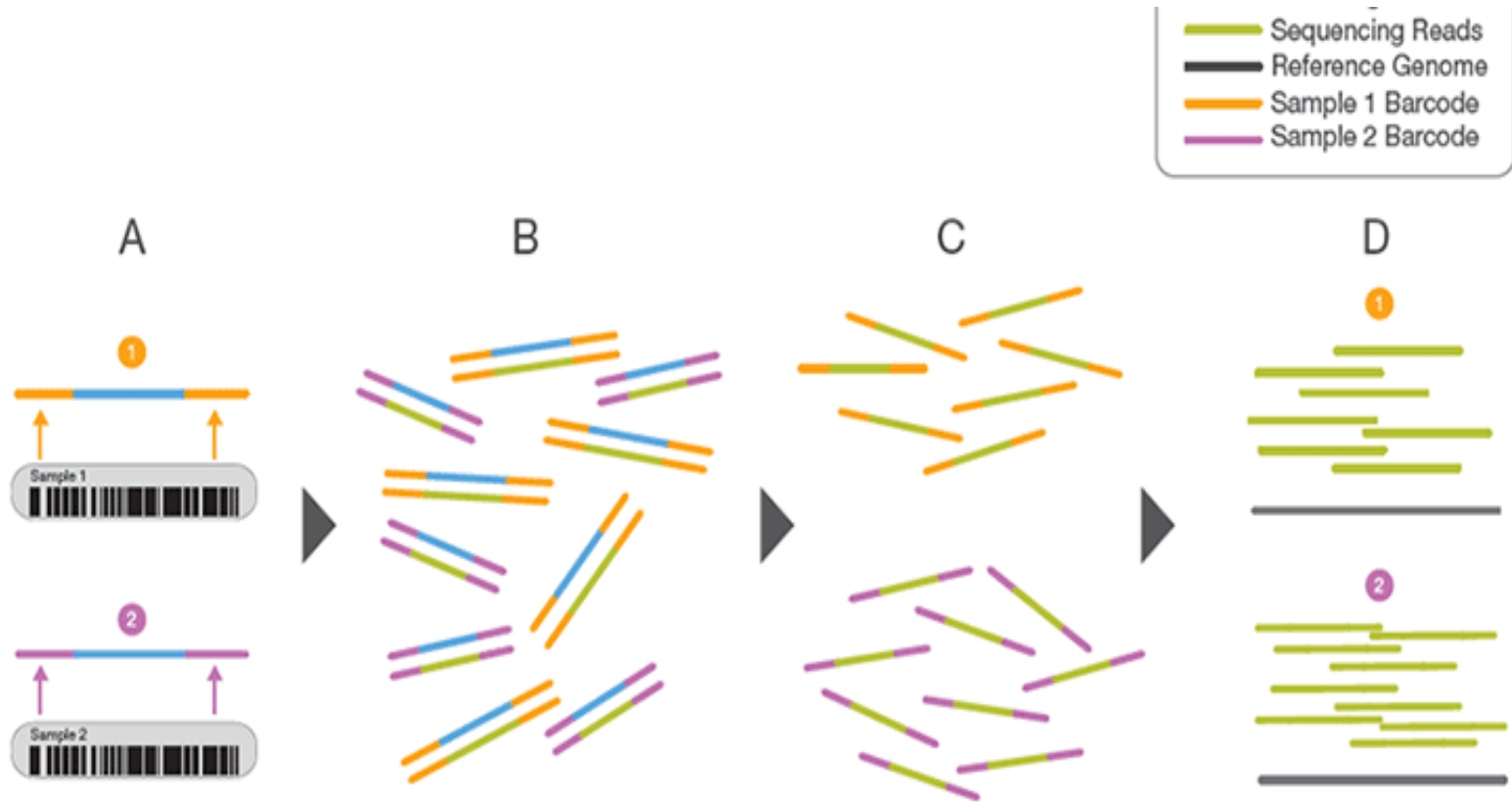
25

Barcoding vs. Lane Effect



# Multiplexing to prevent batch effects

26



- A. Two representative DNA fragments from two unique samples, each attached to a specific barcode sequence that identifies the sample from which it originated.
- B. Libraries for each sample are pooled and sequenced in parallel. Each new read contains both the fragment sequence and its sample-identifying barcode.