# RNA-Seq

# UCSC genome browser

- https://genome.ucsc.edu/

- Authored by Jim Kent

- Came out of human genome project

- Major browser for vertebrate data

- Specific instances exist

# UCSC browser

- Select Genome
- Select Version
- Hide All
- Tracks
- Display options

# Problem

Your lab work on the SCN5A gene. The SCN5A gene belongs to a family of genes that provide instructions for making sodium channels. The sodium channels produced from the SCN5A gene are abundant in heart (cardiac) muscle and play a major role in signaling the start of each heartbeat, and maintaining a normal heart rhythm. Your lab has identifed 5 SNPs in SCN5A in Human which are related to heart failure Your advisor as asked you to choose 2 SNPs for mouse testing. Choose 2 of the 5. Hint: The bed file for these snps is available in the folder Lect_04/SCN5A_hg19.bed.

Also find:

How many splice variants does SCN5A have ?

Do all of them code for a protein ?

Which Strand is SCN5A in ?

Disclaimer: This problem is completely imaginary! Any resemblance to real data is purely coincidental!

# Galaxy

# Exercise 1

- Create a Galaxy account

- Try **not** using TUM-ids instead use real name based email ( firstname.lastname@tum.de) or any other email

# Study Structure

S → E → D          S → E → D

S → E → D          S → E → D

S → E → D

Sample → Experiment → Data

BioProject
PRJNA*

**BioSample**
- SRS*/ERS*
- SAM
- Organism
- Tissue

- ...

**Experiment**
- SRX*/ERX*
- Machine
- Protocol

- ...

**Study**
- SRP*/ERP*
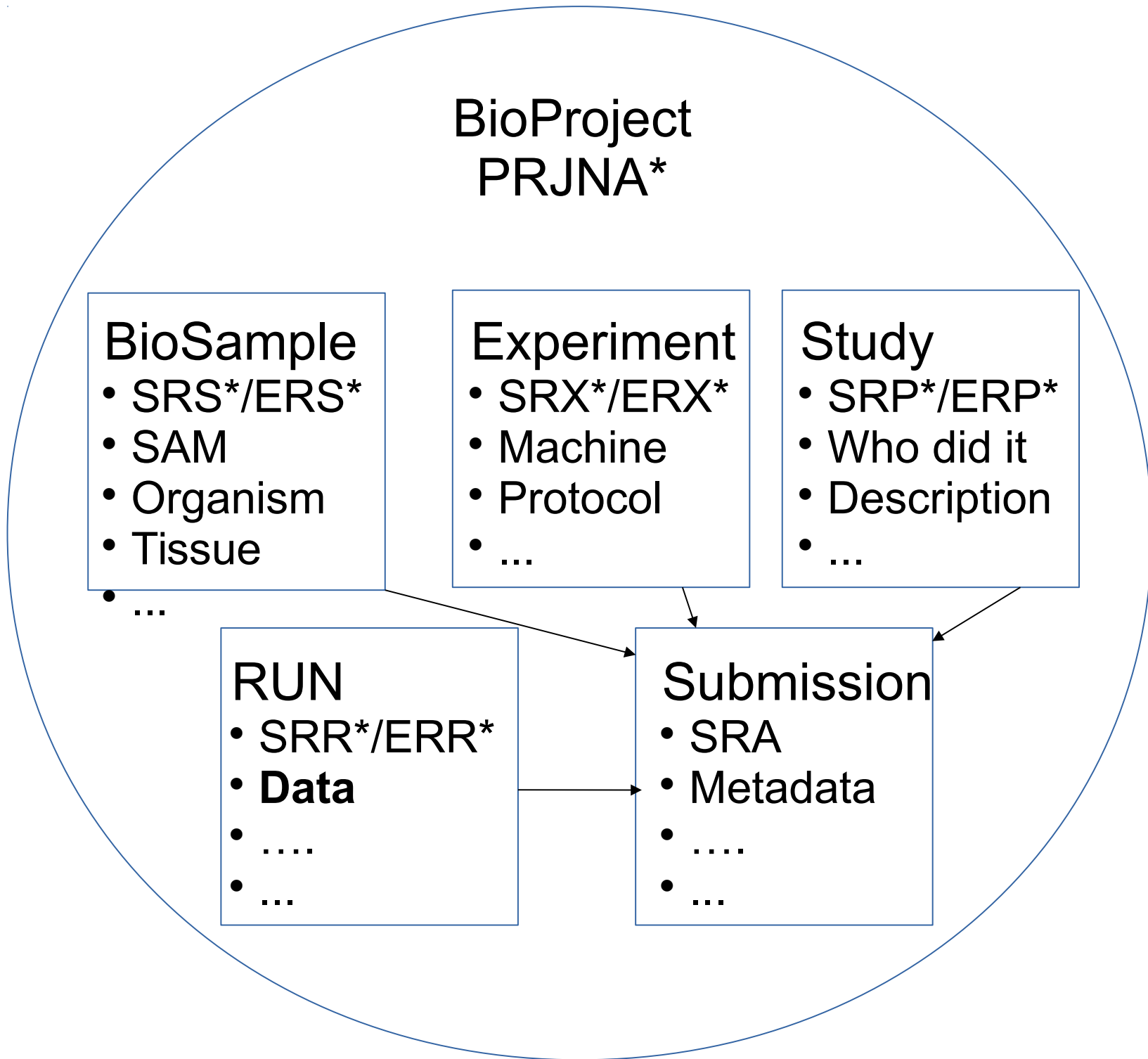- Who did it
- Description

- ...

**RUN**
- SRR*/ERR*
- **Data**
- ....
- ...

**Submission**
- SRA
- Metadata
- ....
- ...

# Splicing misregulation of *SCN5A* contributes to cardiac-conduction delay and heart arrhythmia in myotonic dystrophy

# ARTICLE

# Splicing misregulation of *SCN5A* contributes to cardiac-conduction delay and heart arrhythmia in myotonic dystrophy

Myotonic dystrophy (DM) is caused by the expression of mutant RNAs containing expanded CUG repeats that sequester muscleblind-like (MBNL) proteins, leading to alternative splicing changes. Cardiac alterations, characterized by conduction delays and arrhythmia, are the second most common cause of death in DM. Using RNA sequencing, here we identify novel splicing alterations in DM heart samples, including a switch from adult exon 6B towards fetal exon 6A in the cardiac sodium channel, *SCN5A*. We find that MBNL1 regulates alternative splicing of *SCN5A* mRNA and that the splicing variant of *SCN5A* produced in DM presents a reduced excitability compared with the control adult isoform. Importantly, reproducing splicing alteration of *Scn5a* in mice is sufficient to promote heart arrhythmia and cardiac-conduction delay, two predominant features of myotonic dystrophy. In conclusion, misregulation of the alternative splicing of *SCN5A* may contribute to a subset of the cardiac dysfunctions observed in myotonic dystrophy.

# ARTICLE

# Splicing misregulation of *SCN5A* contributes to cardiac-conduction delay and heart arrhythmia in myotonic dystrophy
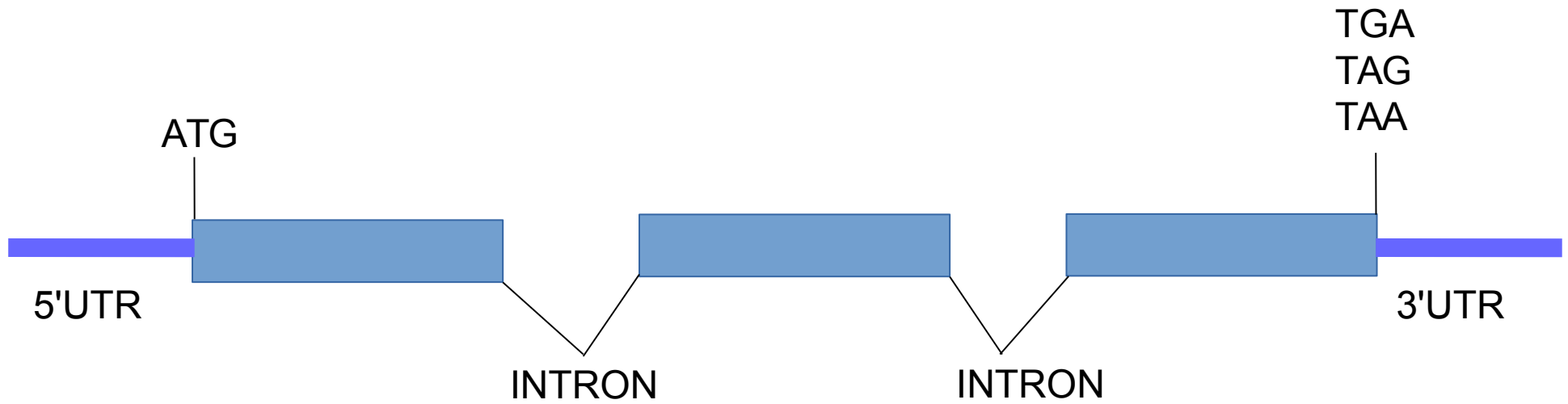
**PRJNA280990**

- Locate the study in SRA run selector
- How many case and control samples
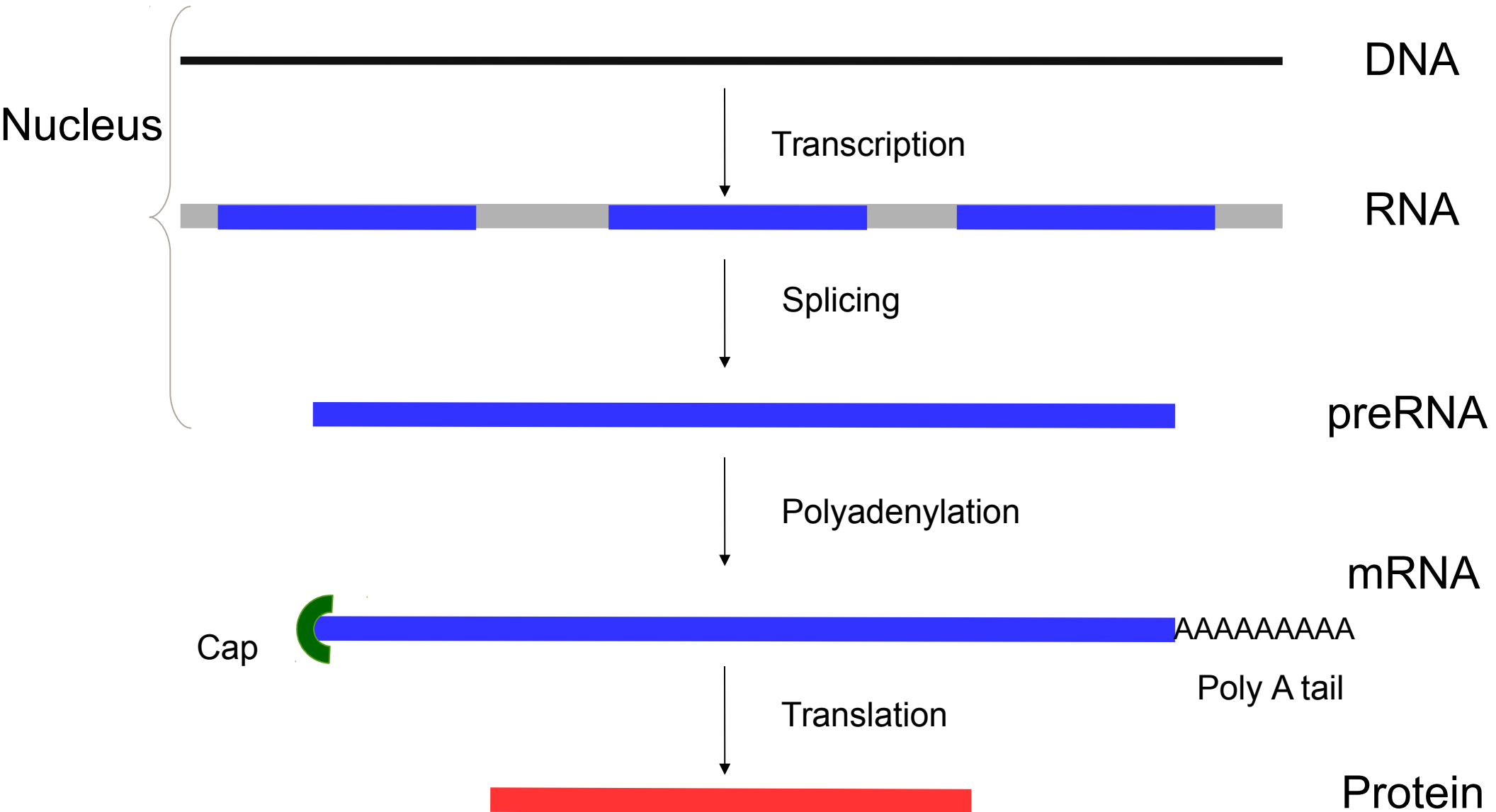- How many technical replicates ?

# RNA-Seq

# What can be done

- Unbiased Expression profiling
- Differential expression analysis
- Gene discovery
- Splice variant Profiling/Discovery
- Regulatory network
- Transcription start site discovery (CAGE)
- Single cell Genomics
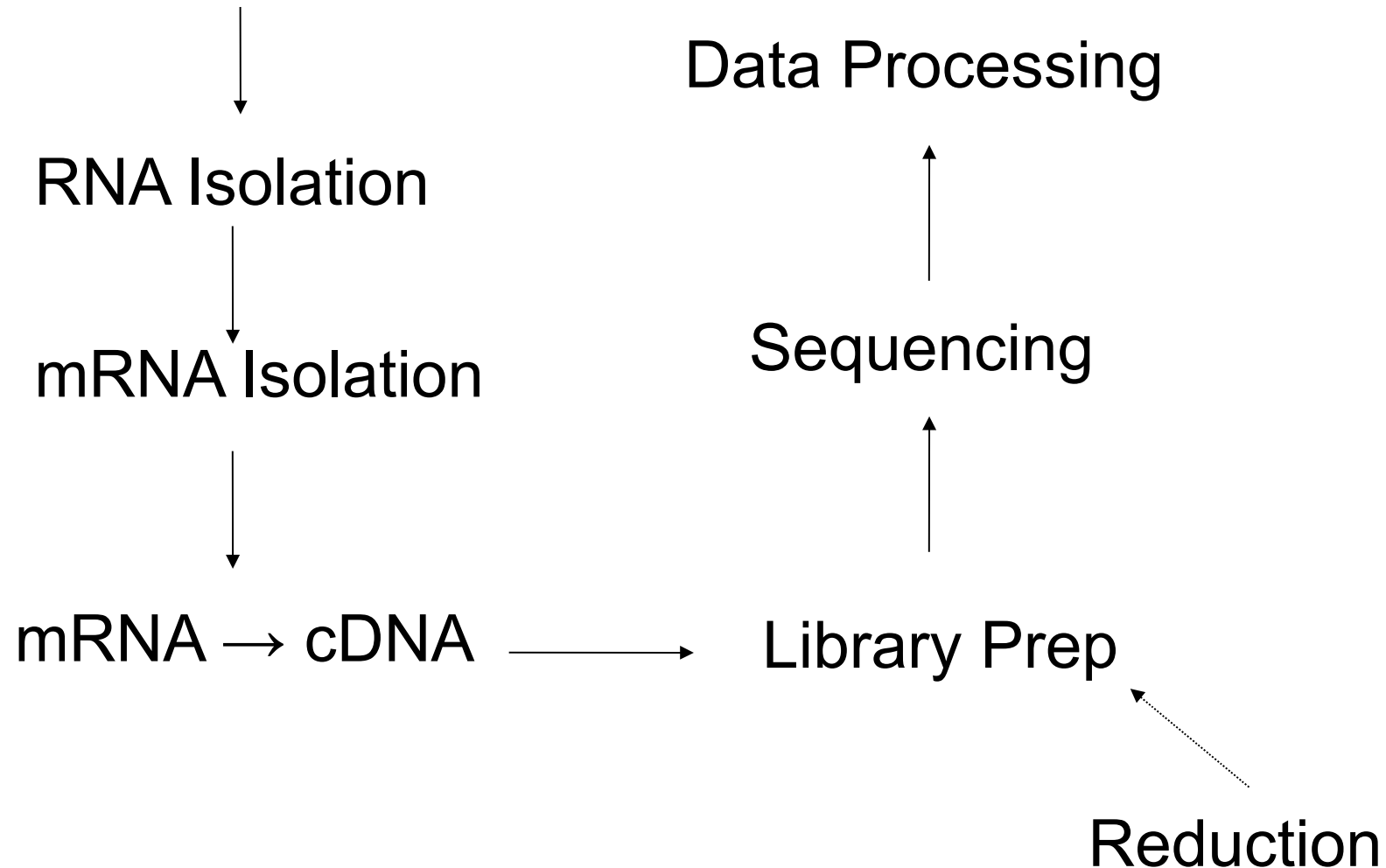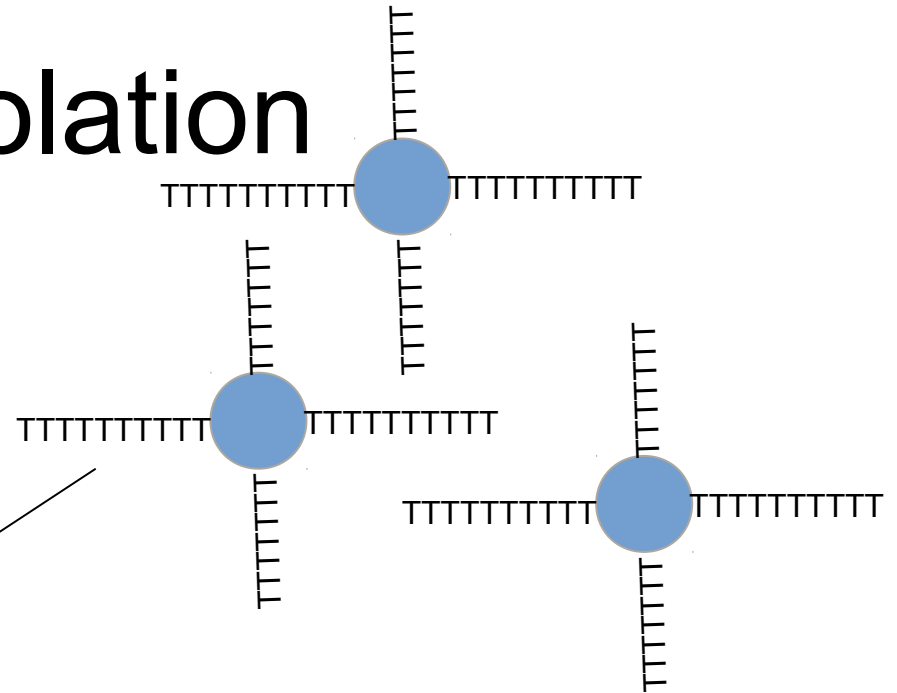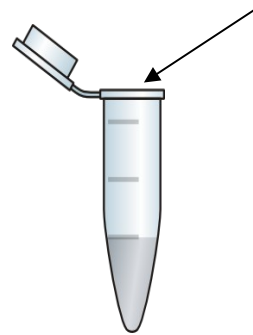- Strand specific RNA-Seq
- …..

# Canonical Gene

# Central Dogma

# RNA-Seq Protocol

Sample Homogenization/Prep

Data Processing

RNA Isolation

mRNA Isolation

Sequencing

mRNA → cDNA ⟶ Library Prep

Reduction

MRNA Isolation

# mRNA → cDNA

mRNA $\xrightarrow[\text{Random Hexamers}]{\text{Reverse Transcriptase}}$ cDNA

cDNA $\downarrow$ Library Prep

Sequencing

RNA-Seq is actually cDNA-Seq

# Splicing

# Problem

Which case has an intron retention event ?


Gene


Gene

# Gold Standards

# Deciding Factors

- Study design
  - Differing read lengths for case and control
  - Differing insert sizes
  - Cases and Controls in different flowcells
  - Differing Library prep
- Technical
  - Differing FastQC results for cases and controls
  - Low mapping in one/more samples

# Deciding Factors

- Biological
  - Ts/Tv
  - Previous Knowledge
  - Hardy-Weinberg Eq
  - Too many non-sense mutations
  - Heterozygous SNPs
  - Unexpected Read Mappings
  - Housekeeping genes

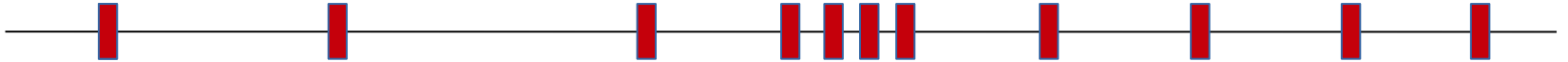# Gold Standards

- Sanger Sequencing

- RT-PCR

- Calculate False Positives and False Negatives

# Problem



Sanger

NGS Param 1

NGS Param 2

Identify the number of false positives and false
negatives in both cases and try to guess the
possible problems and their solutions

# Problem

Blue  Green  Red

100  40  10

You sample 20 reads from the cell. Assuming that all 3 genes R G and B have same length. How many reads do you expect for each gene ?

|   | 20 Reads | 50 Reads |
|---|----------|----------|
| A |          |          |
| B |          |          |
| C |          |          |

# To Do

There is a file shuffle.sh in you Home directory. It randomly takes reads from each gene and prints the counts. The command is ./shuffle.sh n . Where N is the number of reads to sample. Run it for sample size 10, 50 ,100.

# Coverage vs Discovery and accuracy

Discovery/
Accuracy/
Improvement

Low Advantage

Hi Advantage

Coverage/Cost

# Problem

Two genes (A and B) were sequenced in the same library. Gene B is twice as long as gene A. Gene A has 100 reads mapped to it. How many reads would you expect for gene B to map ?

Gene A

Gene B

# Replication

# Technical Replication

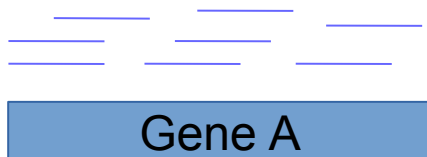# Biological Replication

# Replication

- Microarrays need technical replication
- Technical replicates for RNA-Seq are highly correlated
- Technical replication is generally not required
- Biological replication is considered good for differential expression
- Pooling by barcoding is highly recommended

# Coverage vs Replicates

# Coverage vs Replicates



Low Coverage

Rep1

Low Coverage

Rep2

Low Coverage

Rep1

Low Coverage

Rep2

# Which One to Choose ??

Think Over it!

# Coverage vs Replicates

- How many results you want (Sensitivity) ?

- Is the project exploratory ?

- Monitory Considerations (Barcoding costs money)

- Type of sample

- Noise Level

- Quality of annotation

- More replicates is better

- Typically 3 replicates (more the better)

# Randomizing

- Gene expression is very noisy

- Very plastic to conditions

# Problem

You have 6 (C1 …. C6) case and 6 control samples (O1.... O6). You also have access to 3 lanes and 4 barcodes. How would you arrange the samples ?

|    | Lane1 | Lane2 | Lane3 |
|----|-------|-------|-------|
| B1 |       |       |       |
| B2 |       |       |       |
| B3 |       |       |       |
| B4 |       |       |       |

# Checking for Splicing

- Discovery of splice variants

- Differential splicing between case and control

# Coverage

- Splice variants typically have varying expression

- Coverage can be used to predict a splicing event

- Tissue/Sample wise loss of coverage in an exon

- Differential splicing has been linked to variety of biological effects

# Problem

Can You guess the splice variants and expression of each variant (Qualitatively)  ?

18X     5X     10X

Highly expressed ~10X

Low expressed ~5X

# Second Approach



Intron    Intron

Junction Read

# Junction Read



A well mapped Junction Read is a strong indicator of a splicing event

# TOPHAT Aligner



**TopHat**
A spliced read mapper for RNA-Seq

JOHNS HOPKINS UNIVERSITY
CENTER FOR COMPUTATIONAL BIOLOG
C C B

**TopHat** is a fast splice junction mapper for RNA-Seq reads. It aligns RNA-Seq reads to mammalian-sized genomes using the ultra high-throughput short read aligner Bowtie, and then analyzes the mapping results to identify splice junctions between exons.
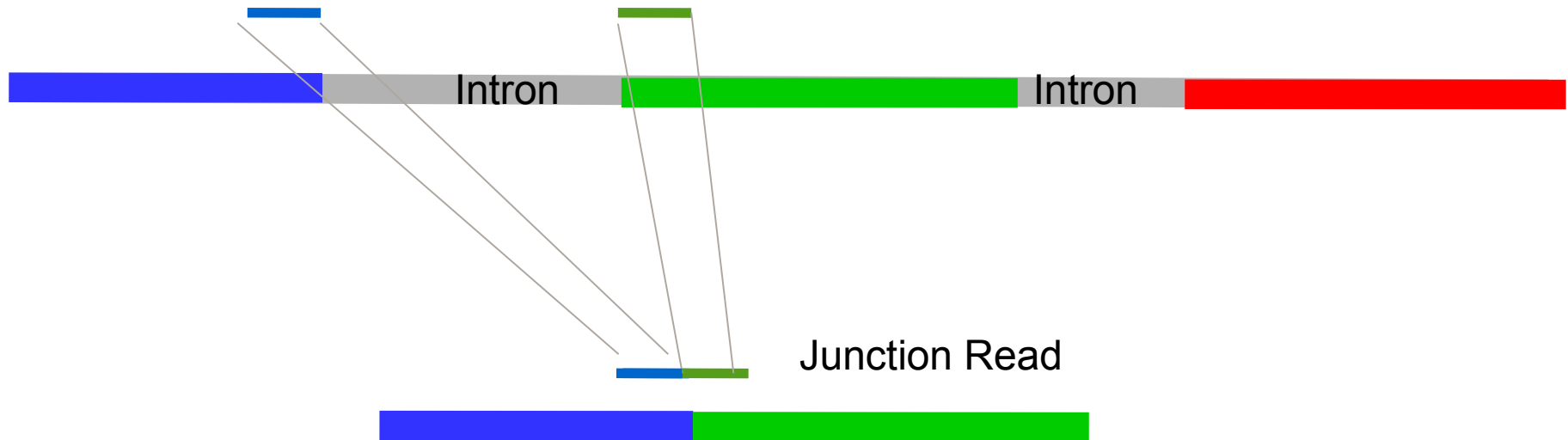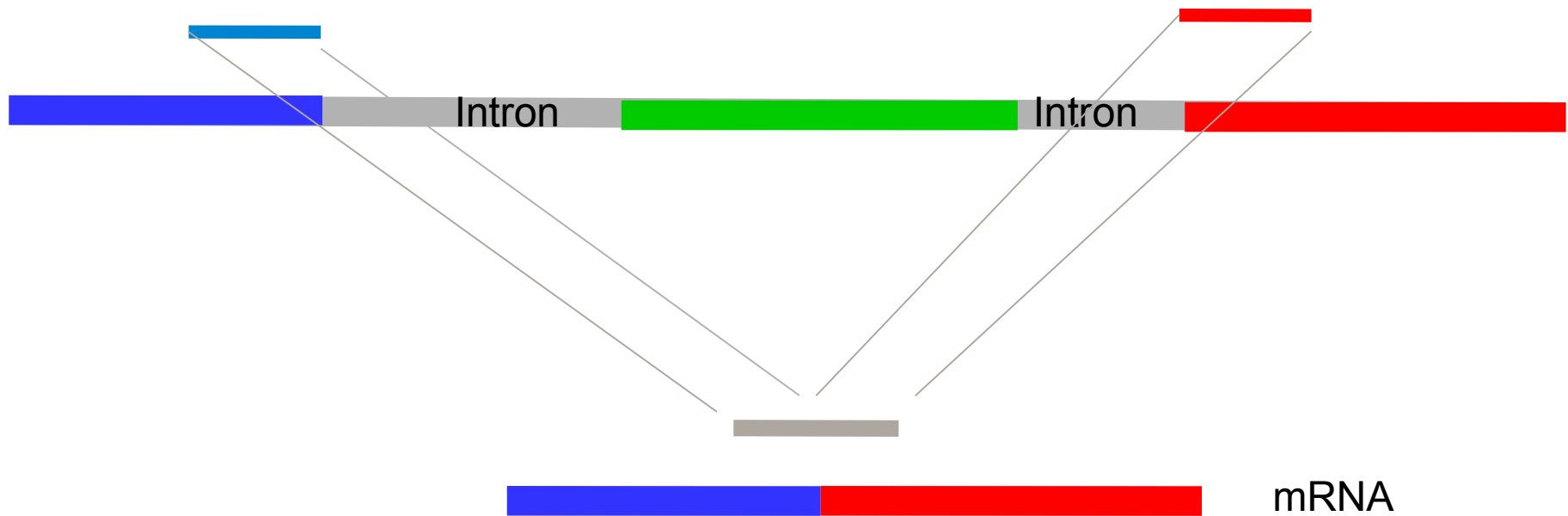
TopHat is a collaborative effort among Daehwan Kim and Steven Salzberg in the Center for Computational Biology at Johns Hopkins University, and Cole Trapnell in the Genome Sciences Department at the University of Washington. TopHat was originally developed by Cole Trapnell at the Center for Bioinformatics and Computational Biology at the University of Maryland, College Park.

OSI certified

- Splice Junction Mapper
- Will use bowtie to map reads normally
- The unmapped reads will then be realigned with large gaps

# Problem

How many splicing events. Which splice pattern does not occur ?
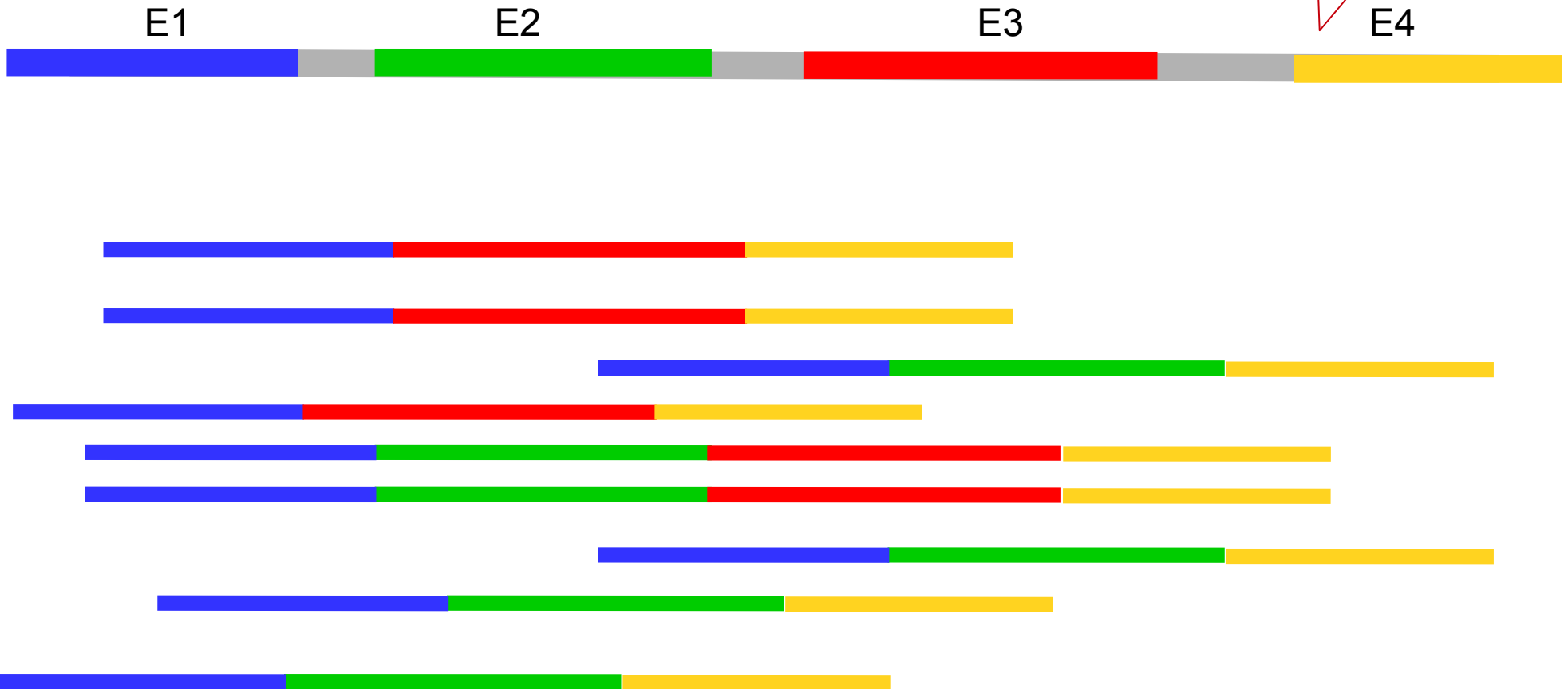
# Expression Quantification

- One value which estimates the abundance of a gene/transcript
- Easy to compute
- Unbiased
  - Gene Length
  - Depth
- Comparable between different samples
- Easily understood

# No of Reads

- Number of reads mapping to a given transcript is a good measure of abundance

- But larger the gene more the reads mapping to it

- More the sequencing depth more reads would be aligned
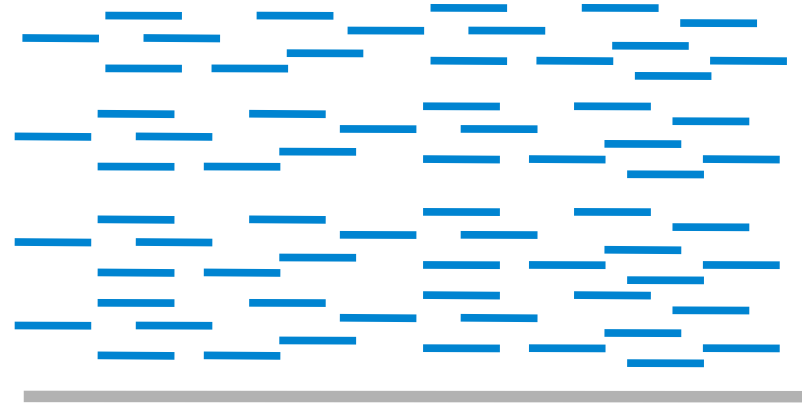
# No of Reads



1Kb

2KB

# Total Sequencing Output

Sample1

Sample2

1 Million Reads

10 Million Reads

# RPKM and FPKM

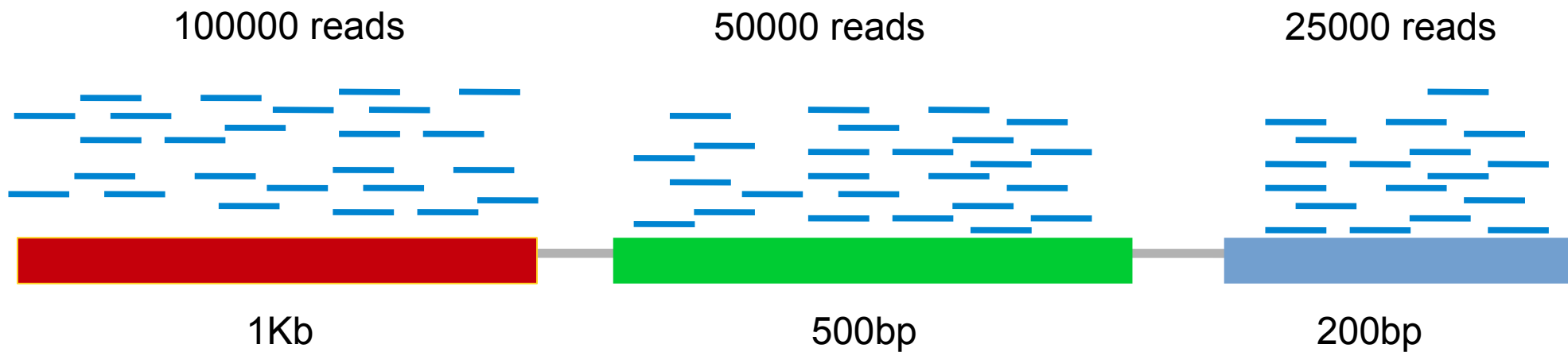Reads Per Kilobase (Transcript) Per Million Mapped Reads.

Fragments Per Kilobase (Transcript)  Per Million Mapped Reads.

$$RPKM = \frac{\text{Total Reads mapped to a transcript}}{(\text{Length of region (Kb)}*(\text{Total Mapped reads in Mb}))}$$

# Problem

The genome consists of 3 genes. Calculate RPKM for the green gene.

100000 reads                    50000 reads                    25000 reads

1Kb                             500bp                          200bp

# Problem

You are analyzing a case control study. The case bam file has 2 million mapped reads whereas the control bam file has 1.5 million mapped reads. Is gene A differently expressed and what is the ratio of case/control expression ?

50000 reads

30000 reads

Gene A

Gene A

Control

Case