

RNA-seq: From (good) experimental design to (accurate) gene expression abundance.

Steve Munger
Narayanan Raghupathy

The Jackson Laboratory
21st Century Mouse Genetics
11 August 2016

Outline

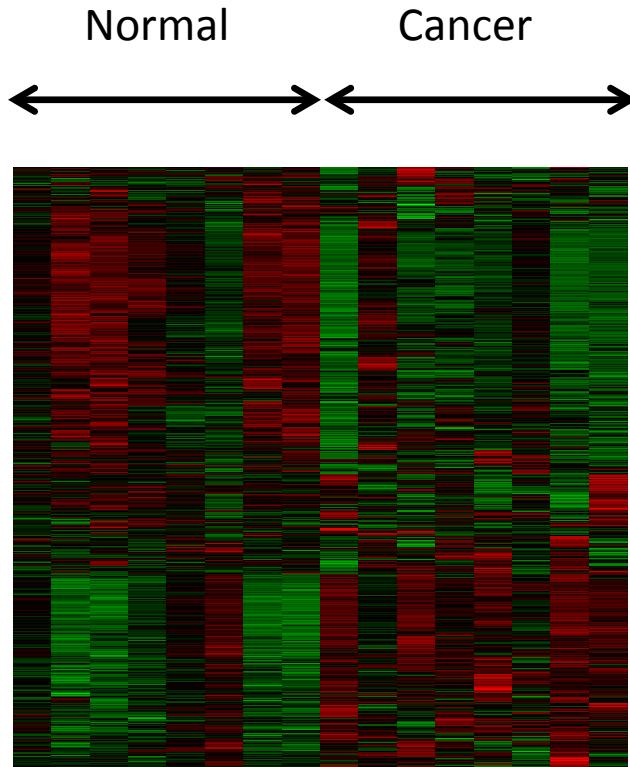
General overview of RNA-seq analysis.

- Introduction to RNA-seq
- The importance of a good experimental design
- Quality Control
- Read alignment
- Quantifying isoform and gene expression
- Normalization of expression estimates

RNA-seq: Sequencing Transcriptomes

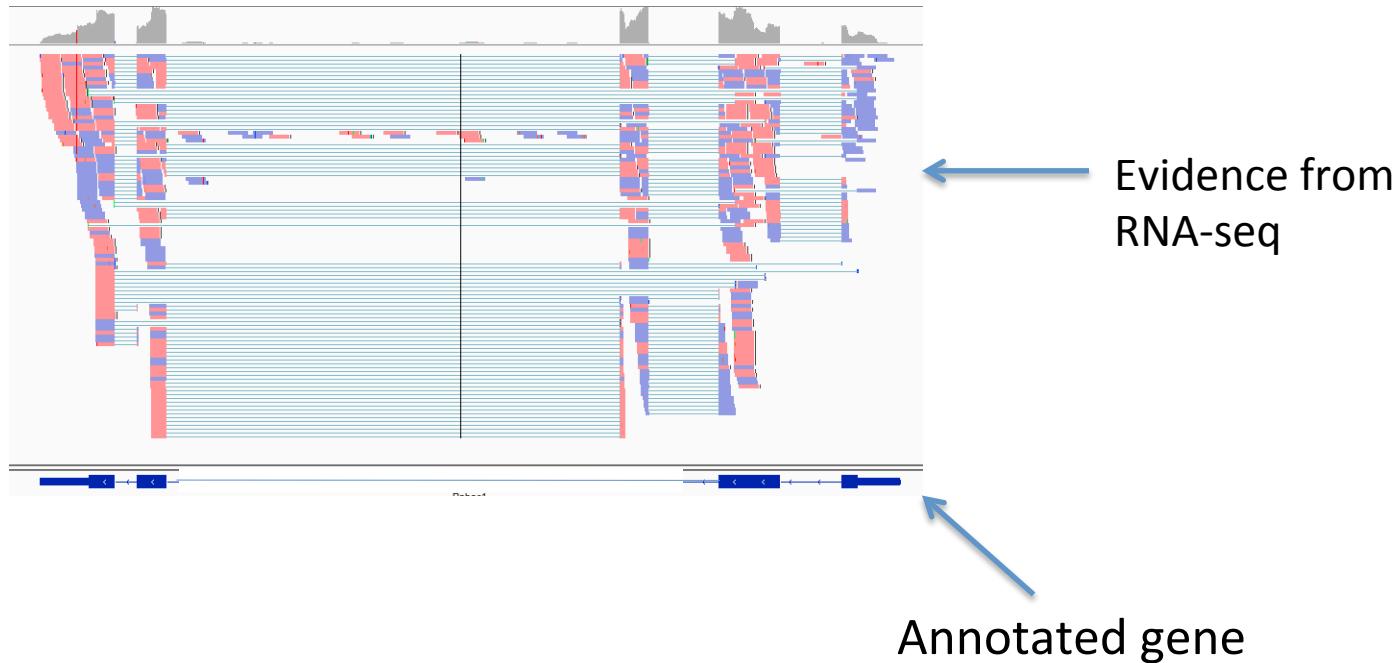


Applications of RNA-seq Technology



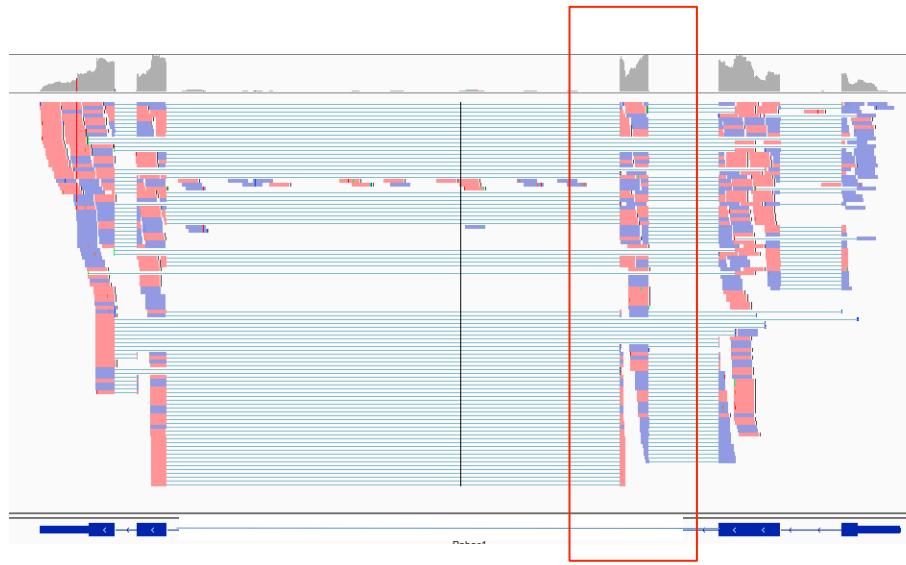
Differential Gene expression analysis

Applications of RNA-seq Technology



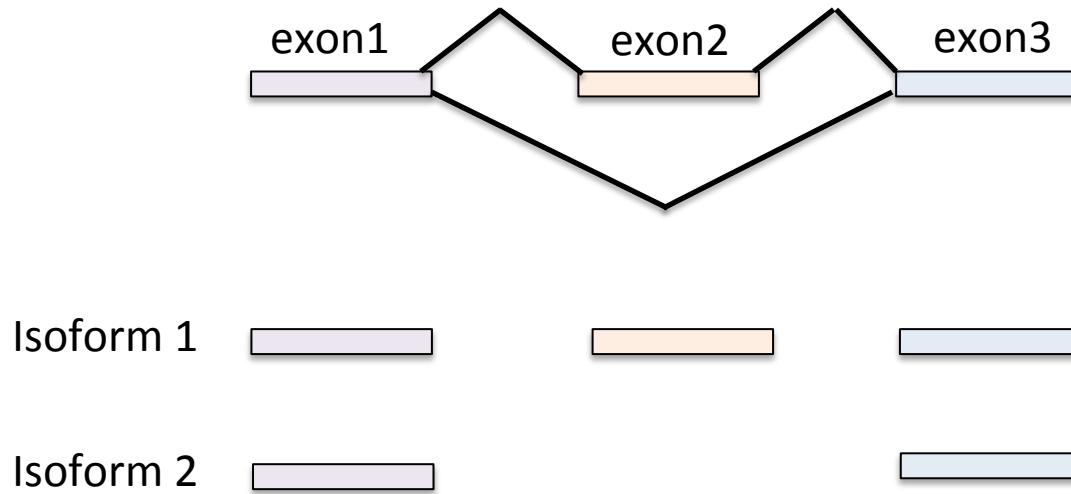
Novel exon discovery

Applications of RNA-seq Technology



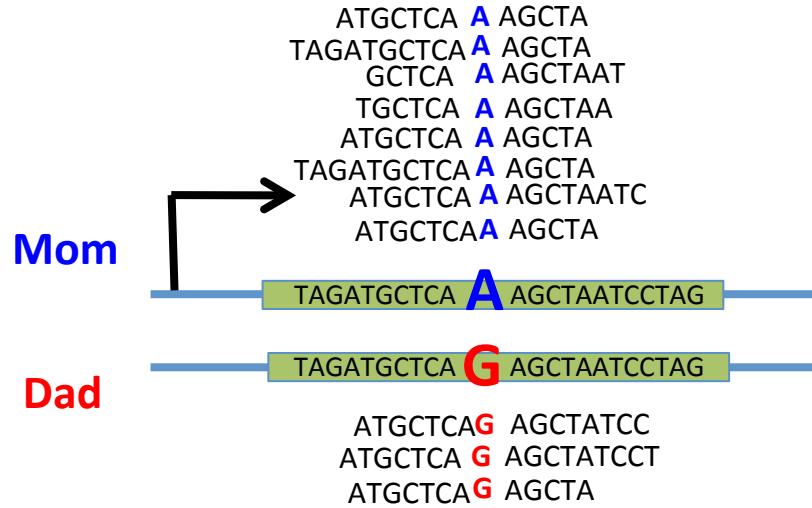
Novel exon discovery

Applications of RNA-seq Technology



Alternative splicing

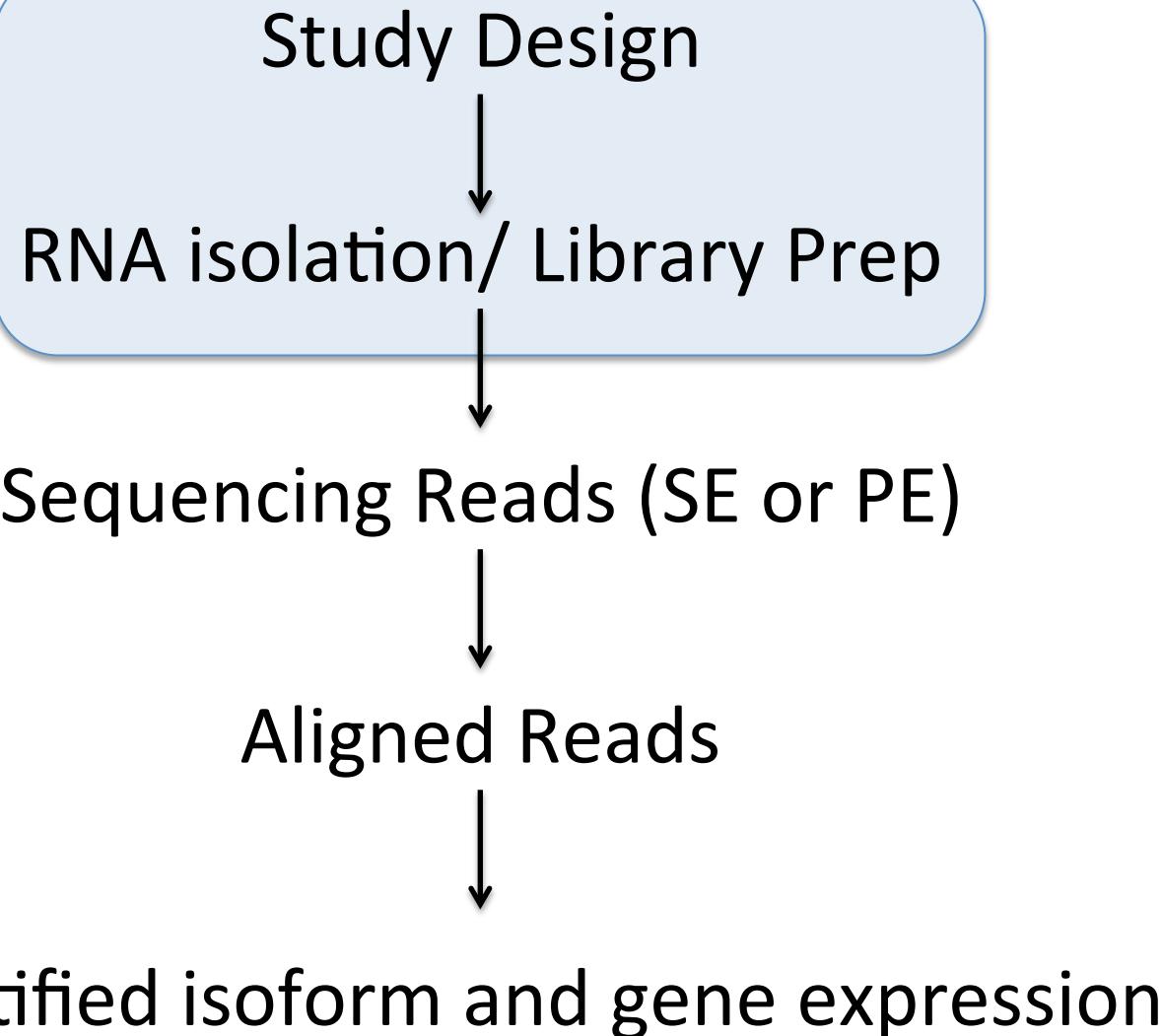
Applications of RNA-seq Technology



Allele-Specific gene Expression (ASE)

Preferential expression of one allele over the other.

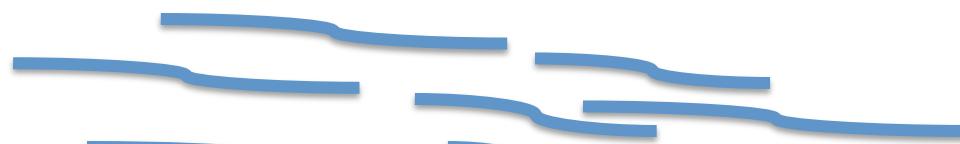
RNA-seq Work Flow



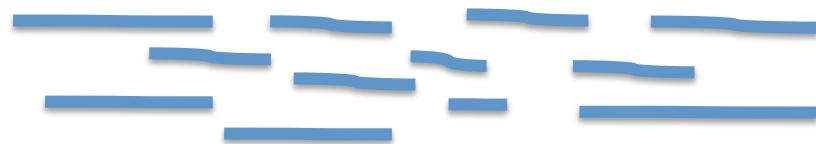
Total RNA

RNA-Seq

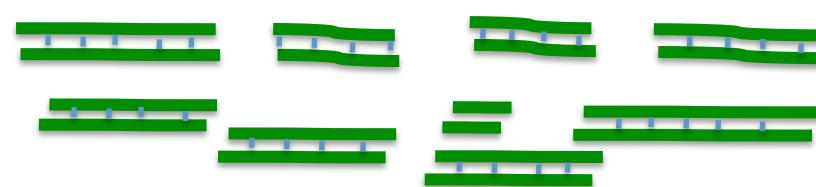
mRNA



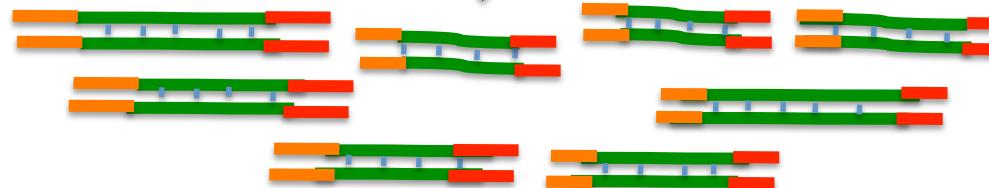
mRNA after
fragmentation



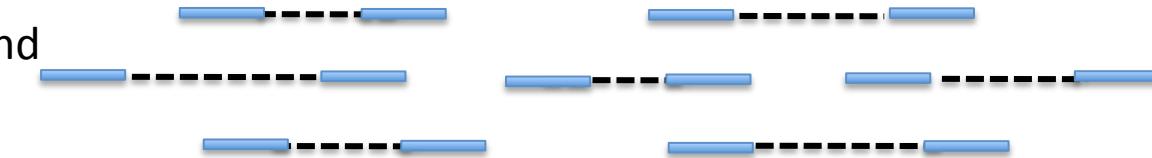
cDNA



Adaptors ligated
to cDNA



Single/ Paired End
Sequencing



Know your application – Design your experiment accordingly

- How many reads? Read depth
- Single-end or Paired-end sequencing?
- Read length?
- How many samples?

RNA-seq Experimental design

- Differential expression of highly expressed and well annotated genes?
 - Smaller sample depth; more biological replicates
 - No need for paired end reads; shorter reads (50bp) may be sufficient.
 - Better to have 20 million 50bp reads than 10 million 100bp reads.
- Looking for novel genes/splicing/isoforms?
 - More read depth, paired-end reads from longer fragments.

Good Experimental Design

Multiplexing
Replication
Randomization

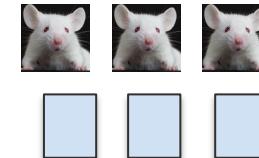


Illumina flowcell

RNA-Seq Experimental Design: Randomization

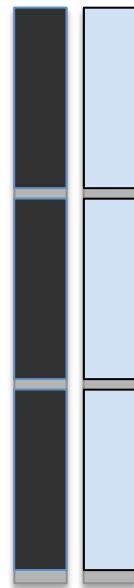


Experimental Group 1



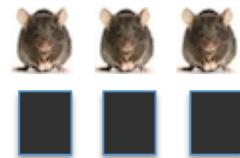
Experimental Group 2

Two Illumina Lanes



Bad Design

RNA-Seq Experimental Design: Randomization

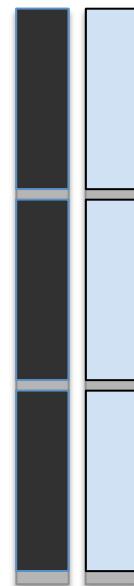


Experimental Group 1

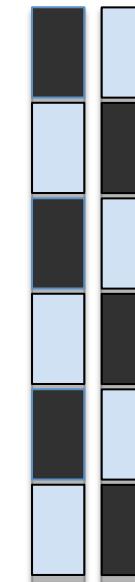


Experimental Group 2

Two Illumina Lanes

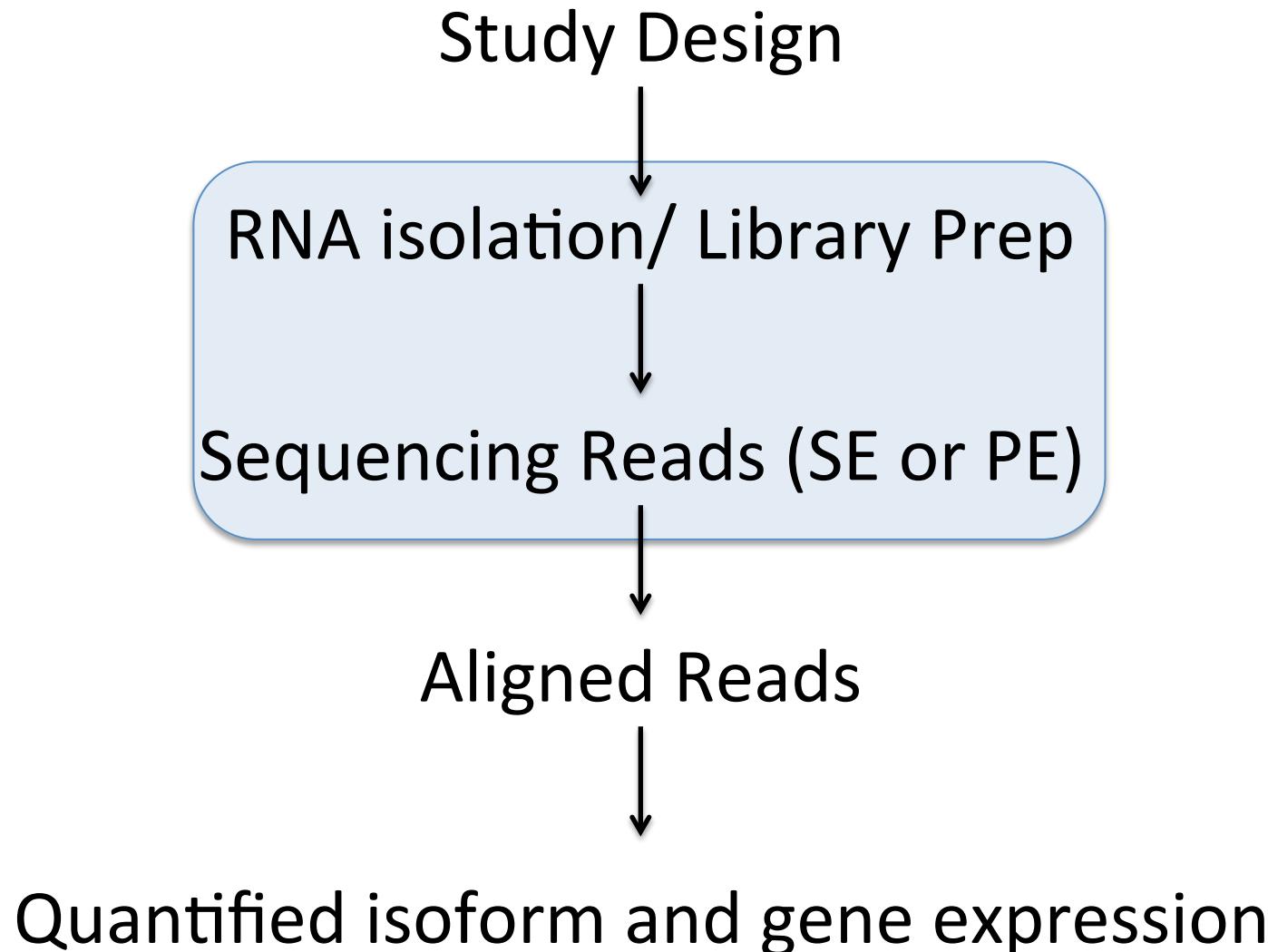


Bad Design



Better Design

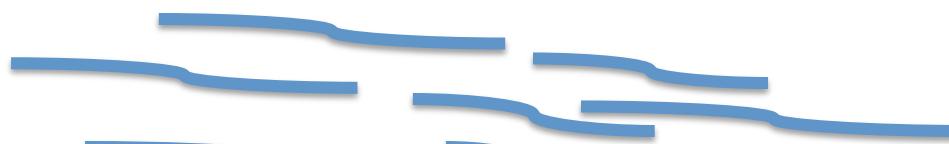
RNA-seq Work Flow



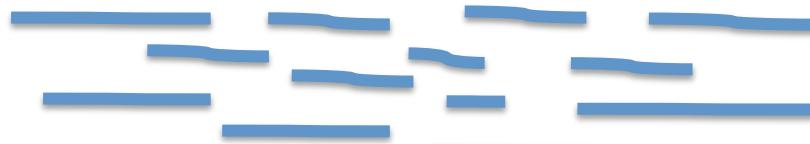
Total RNA

RNA-Seq

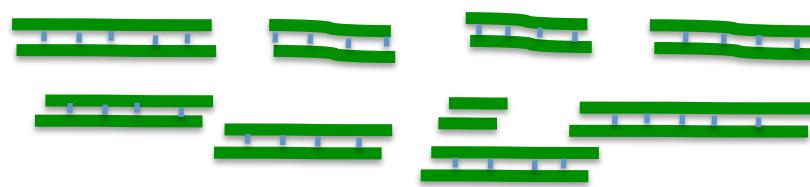
mRNA



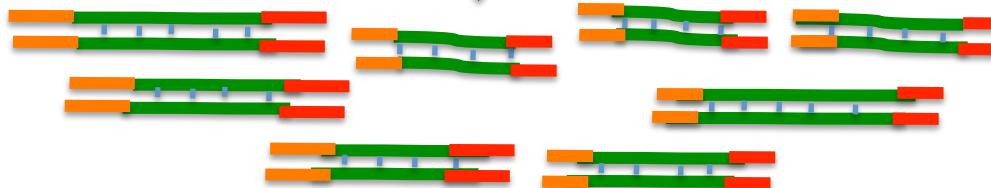
mRNA after
fragmentation



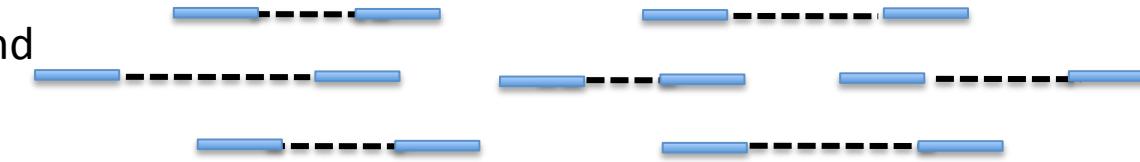
cDNA



Adaptors ligated
to cDNA



Single/ Paired End
Sequencing



Millions and millions of reads...

```
@HISEQ2000_0074:8:1101:7544:2225#TAGCTT/1  
TCACCCGTAAGGTAAACAAACCGAAAGTATCCAAAGCTAAAAGAAGTGGACGACGTGCTGGTGGAGCAGCTGCATG  
+  
CCCCFFFFFHHHHDDHHJJJJJJJJ?FGIIIIIIIIIIIIJJFHIIJJIJHHHFFFFD>AC?B??C?ACCAC>BB<<>C@CCCACCCDCCIJ
```

@HISEQ2000_0074:8:1101:7544:2225#TAGCTT/1

Instrument: run/flowcell id

Flowcell lane and tile number

X-Y Coordinate in flowcell

The member of a pair

Index Sequence

$$Q = -10 \log_{10} P$$

10 indicates 1 in 10 chance of error

20 indicates 1 in 100,

30 indicates 1 in 1000,

Phred Score:

SN

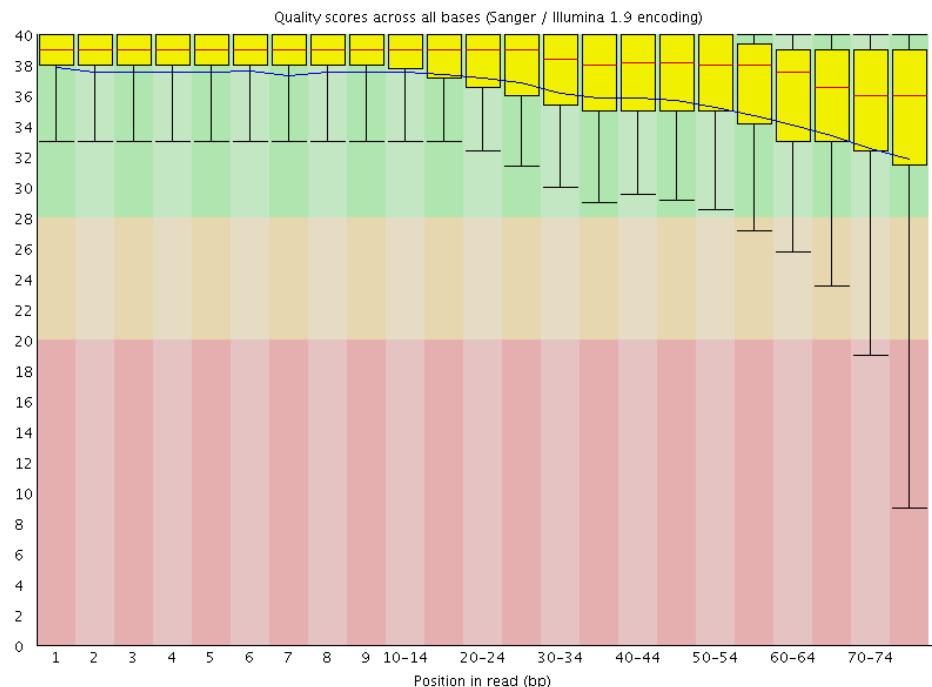
Quality Control: How to tell if your data is clean

- FASTX-Toolkit
 - http://hannonlab.cshl.edu/fastx_toolkit/
- FastQC
 - <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

RNA-seq Data: <ftp://ftp.jax.org/dgatti/MouseGen2016/>

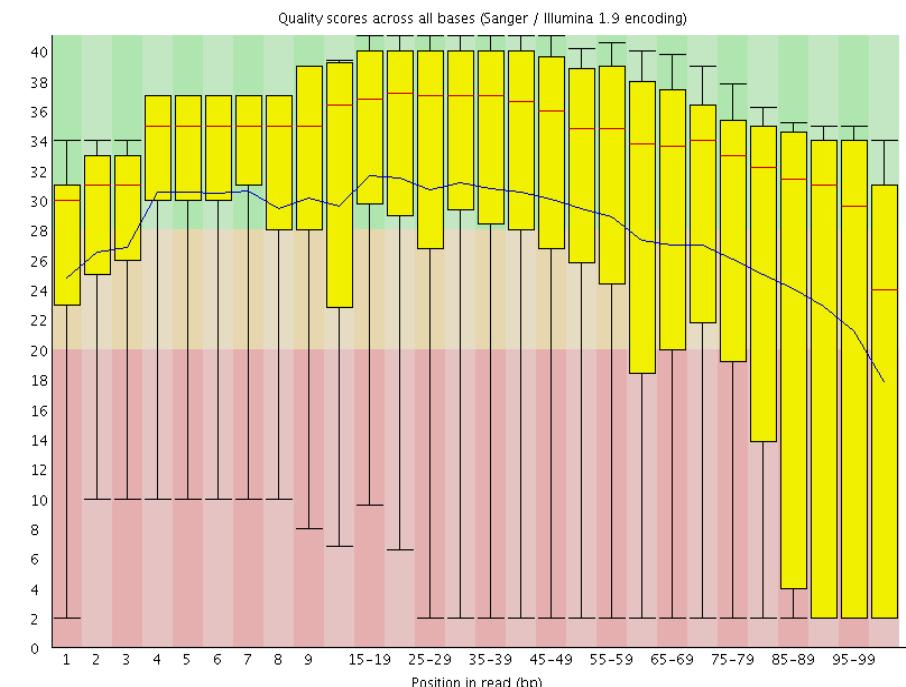
- B6-100K.fastq and Cast-100K.fastq

Quality Control: How to tell if your data is clean



Good data

- Consistent
- High Quality Along the reads



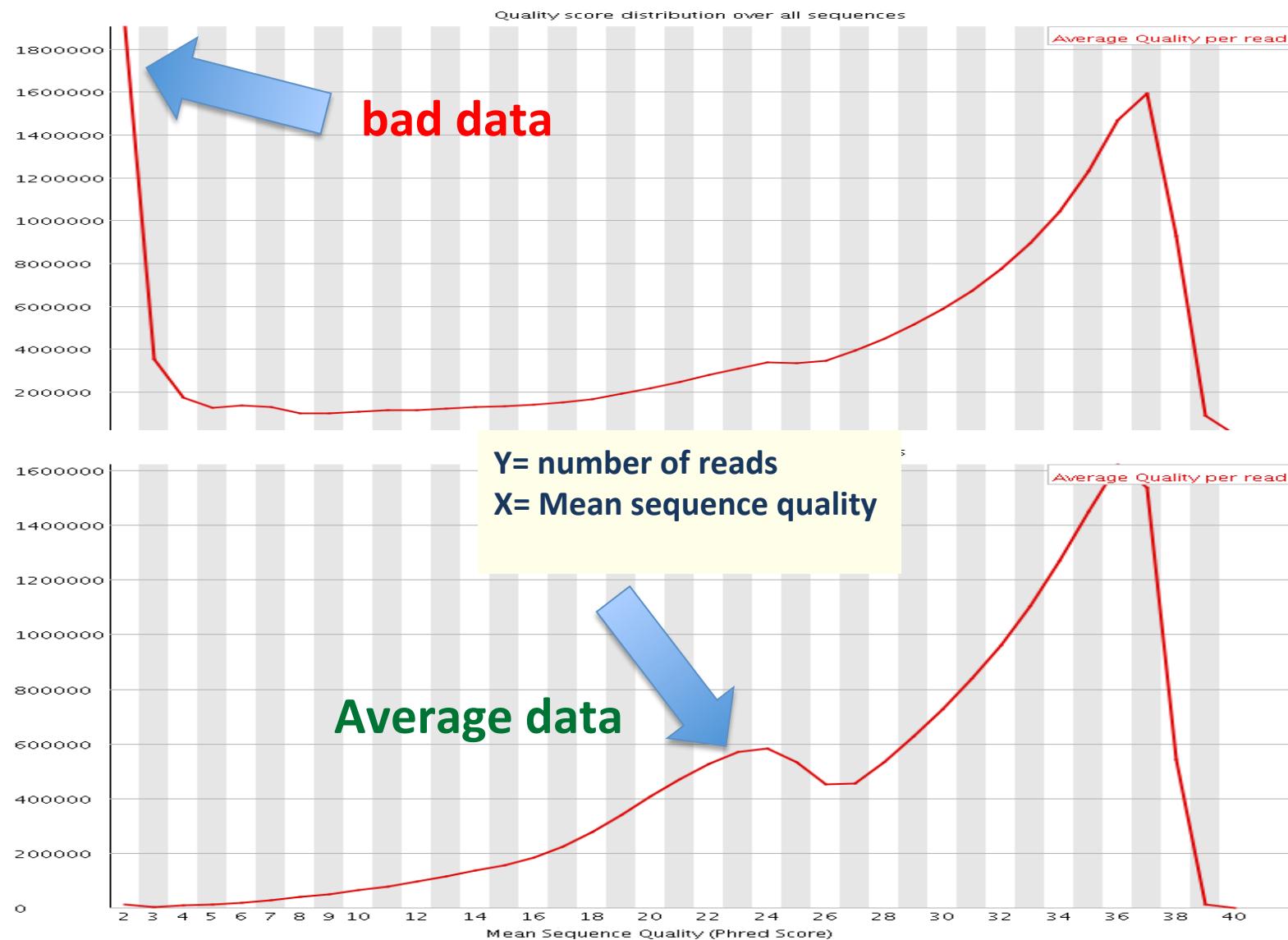
Bad data

- High Variance
- Quality Decrease with Length

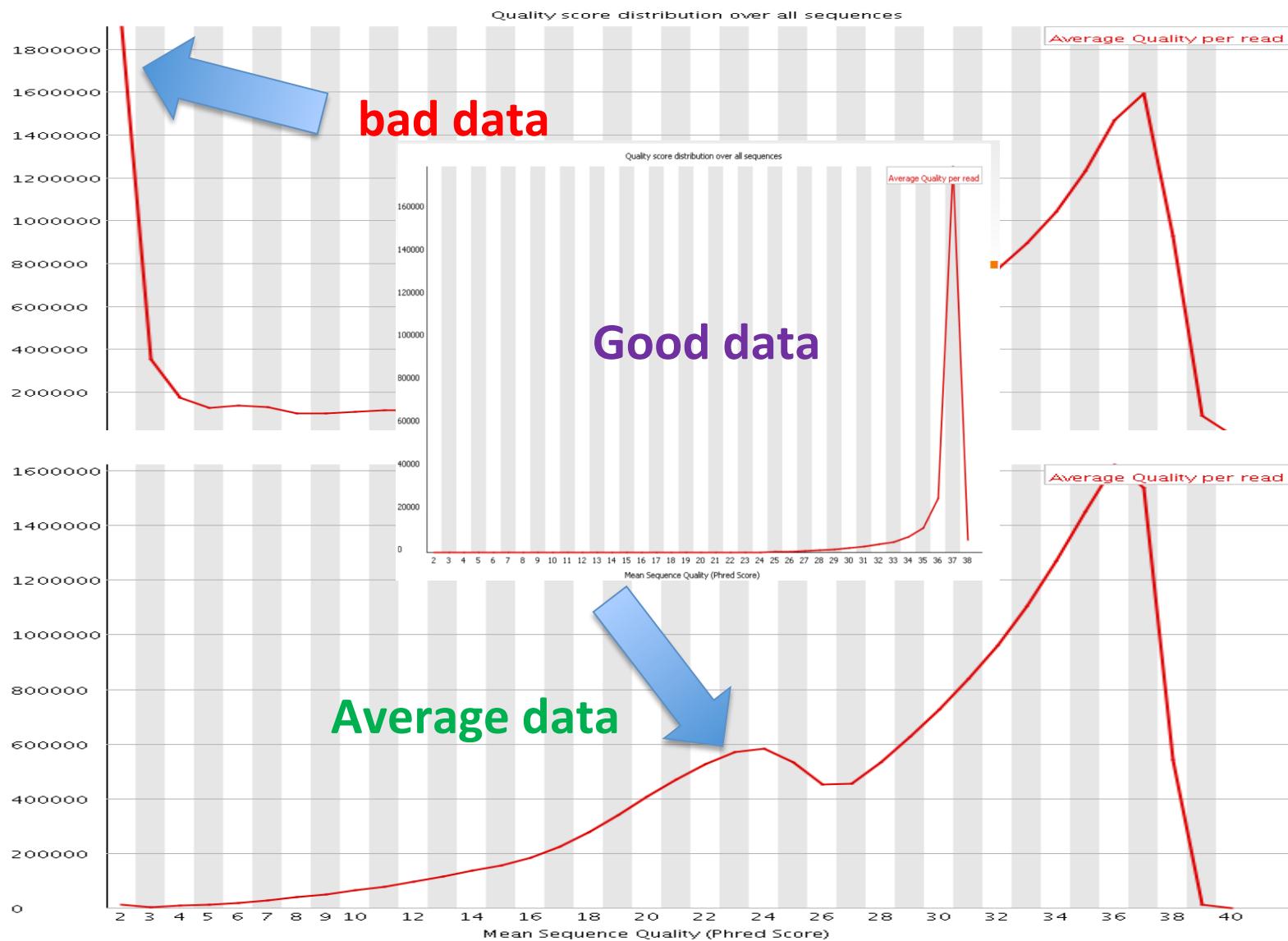
RNA-seq Data: <ftp://ftp.jax.org/dgatti/MouseGen2016/>

- B6-100K.fastq and Cast-100K.fastq

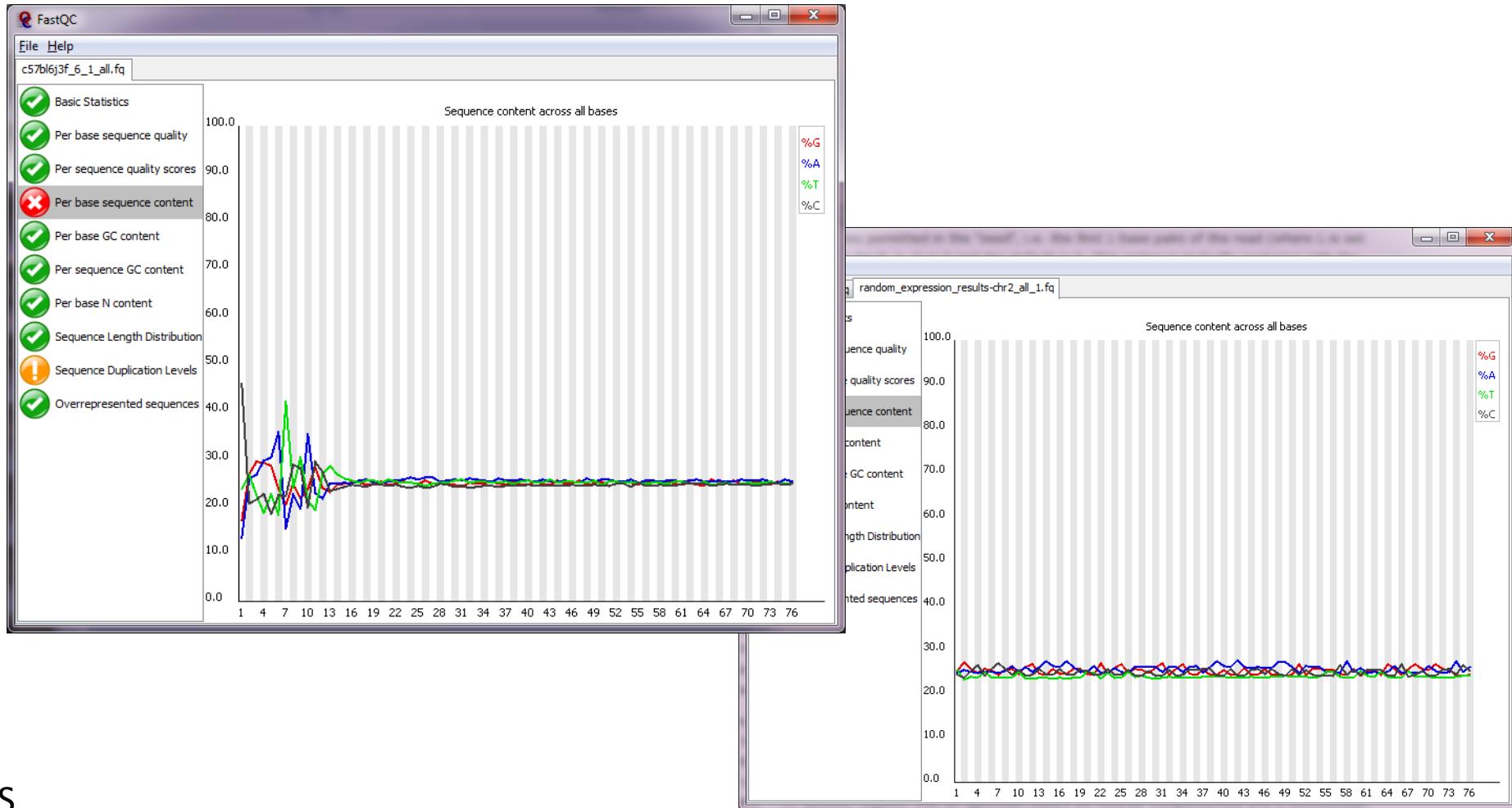
Per sequence quality distribution



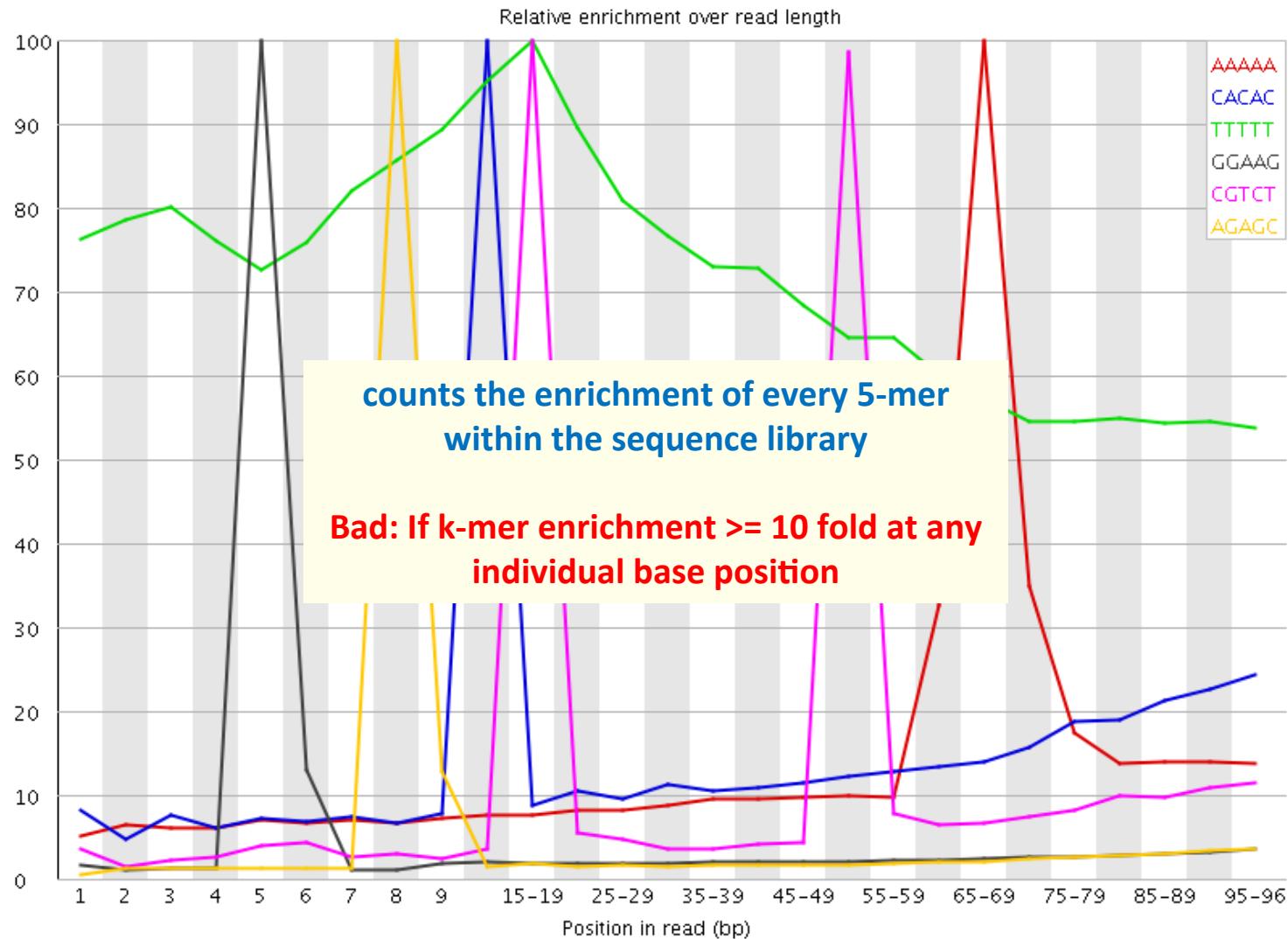
Per sequence quality distribution



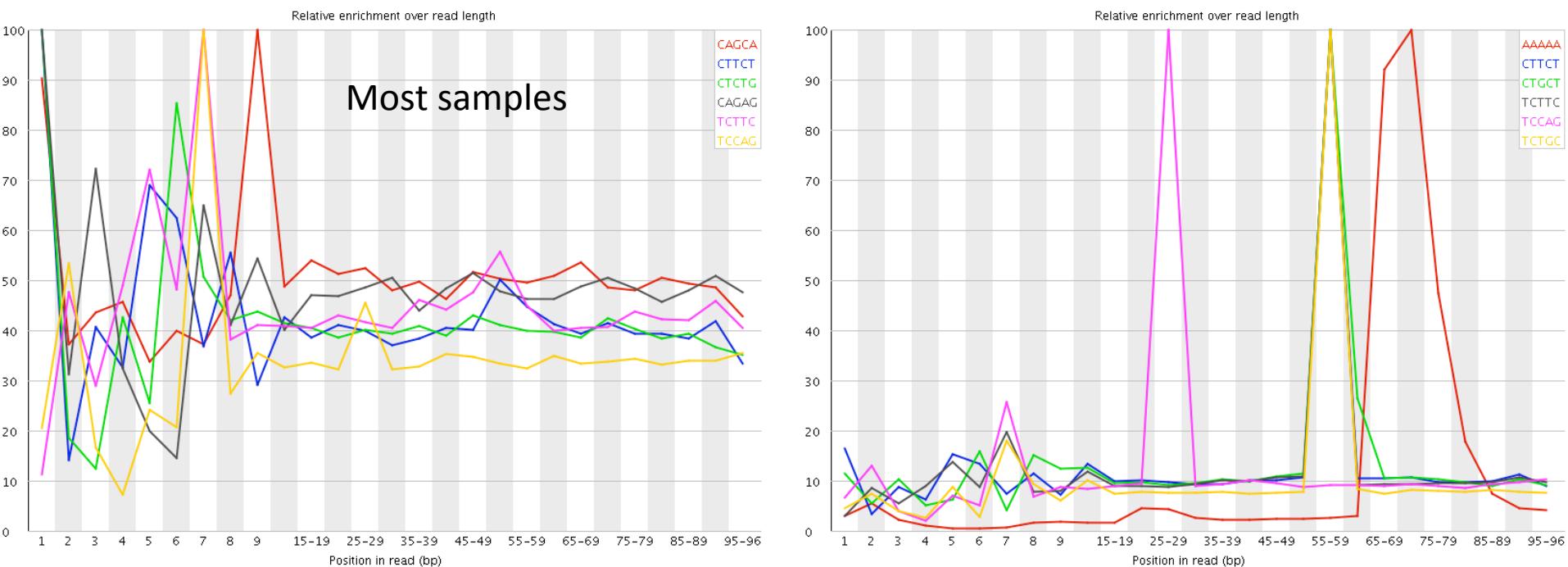
Quality Control: Sequence Content Across Bases



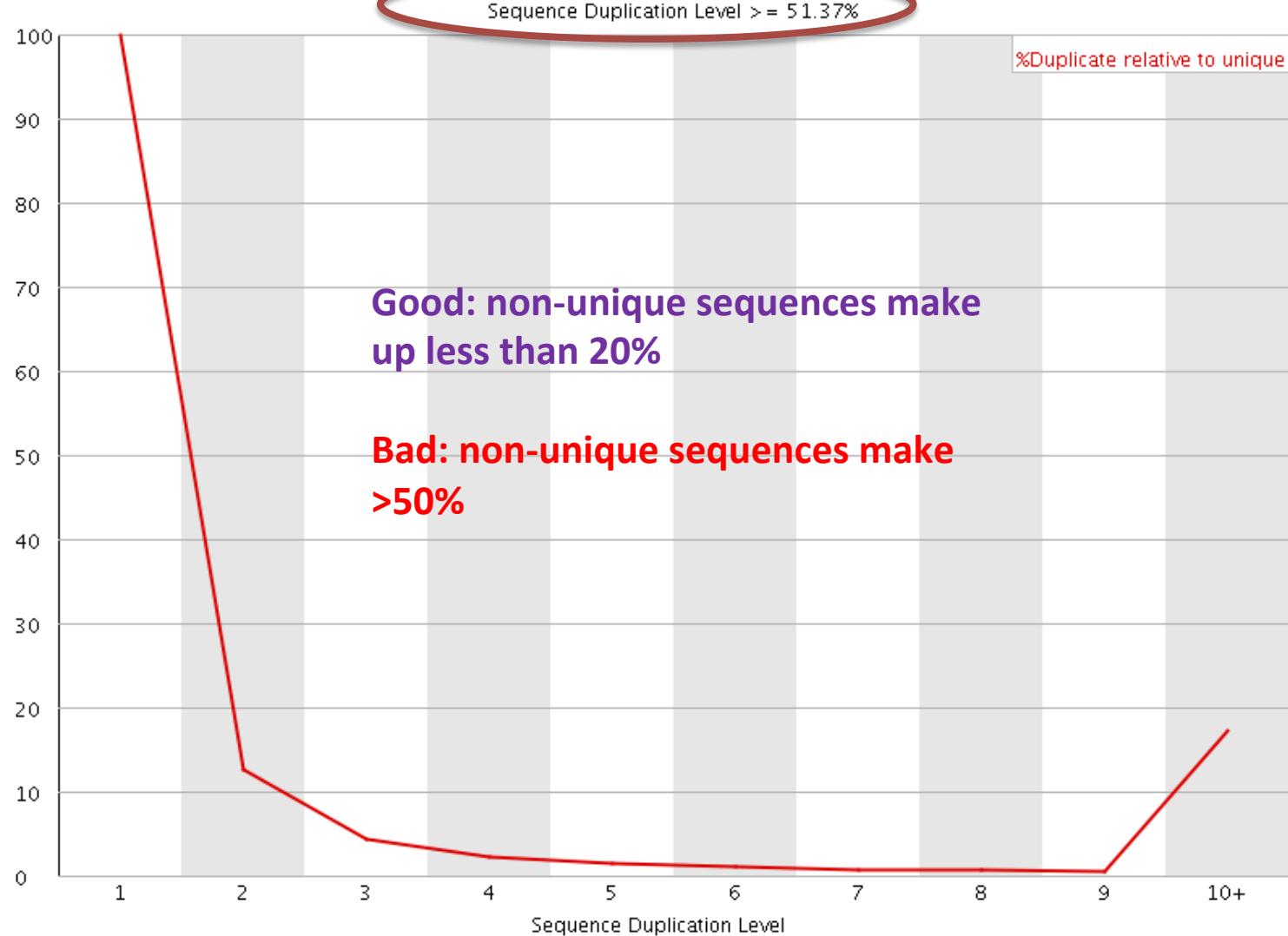
K-mer content



K-mer content



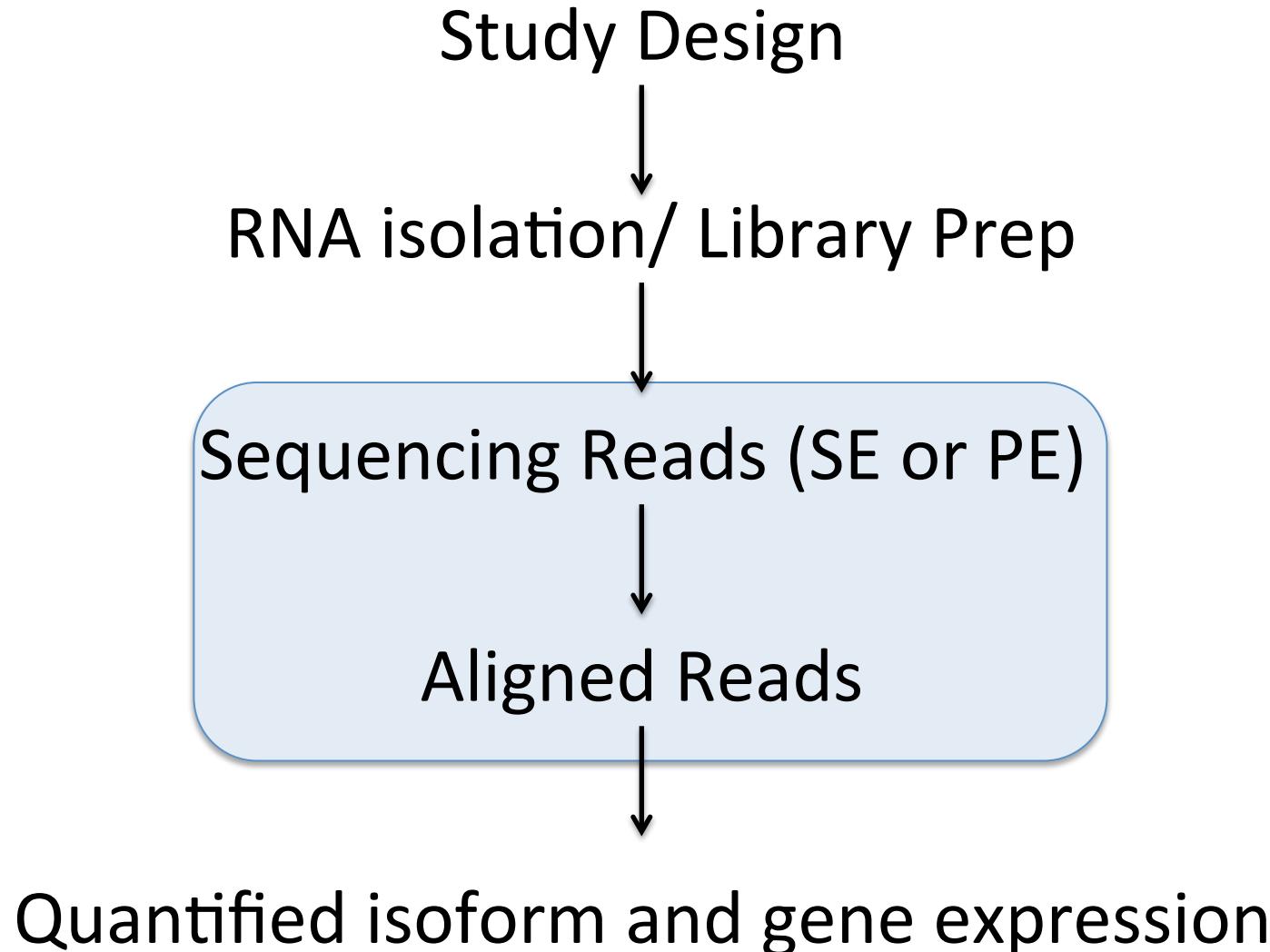
Duplicated sequences



Tradeoffs to preprocessing data

- Signal/noise -> Preprocessing can remove low-quality “noise”, but the cost is information loss.
 - Some uniformly low-quality reads map uniquely to the genome.
 - Trimming reads to remove lower quality ends can adversely affect alignment, especially if aligning to the genome and the read spans a splice site.
 - Duplicated reads or just highly expressed genes?
 - Most aligners can take quality scores into consideration.
 - Currently, we do not recommend preprocessing reads aside from removing uniformly low quality samples.

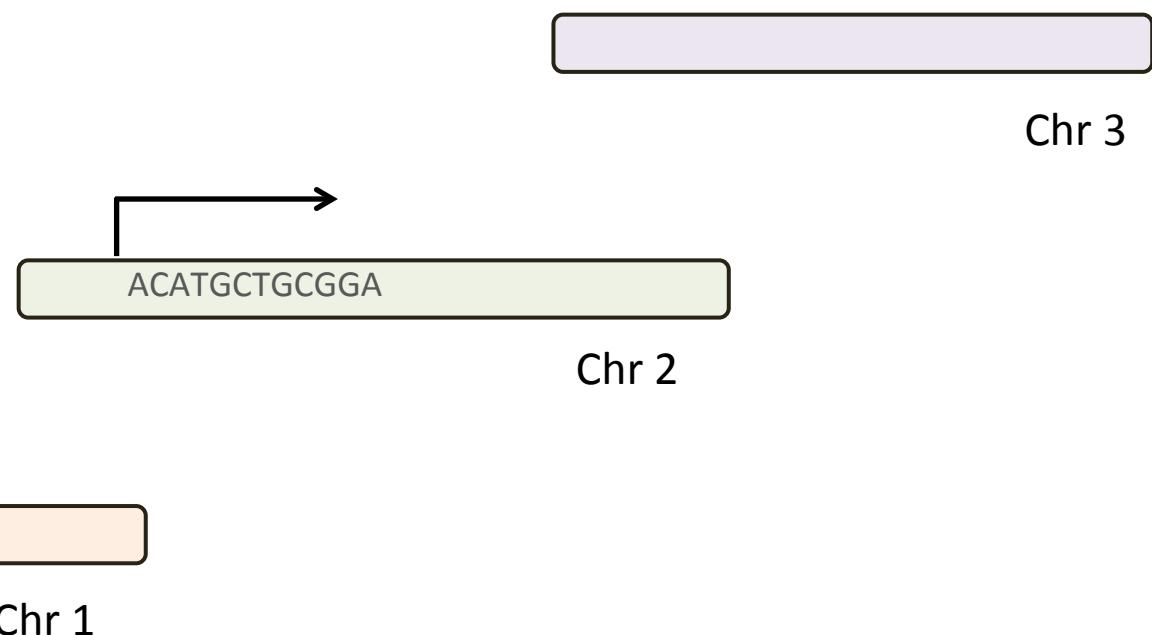
RNA-seq Work Flow



Alignment 101

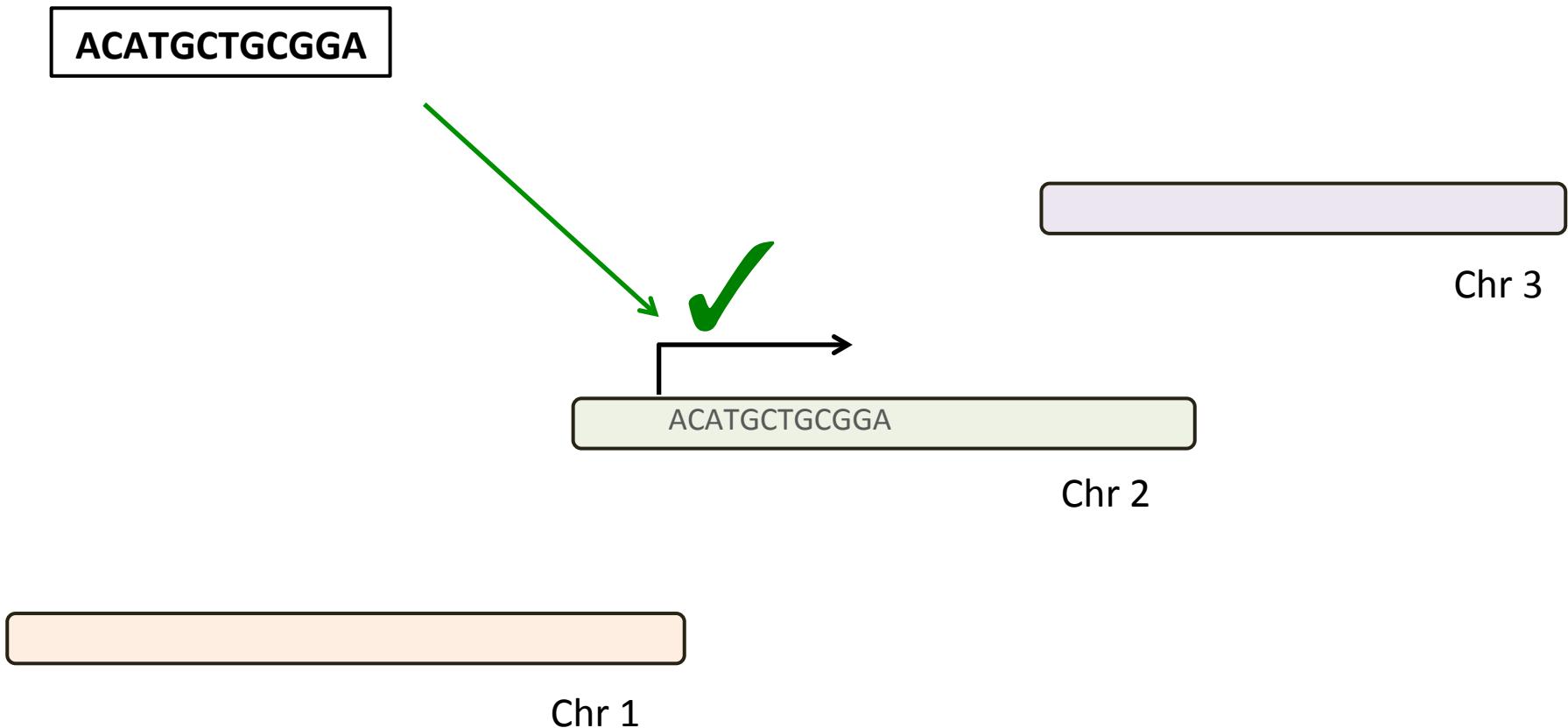
100bp Read

ACATGCTGCGGA



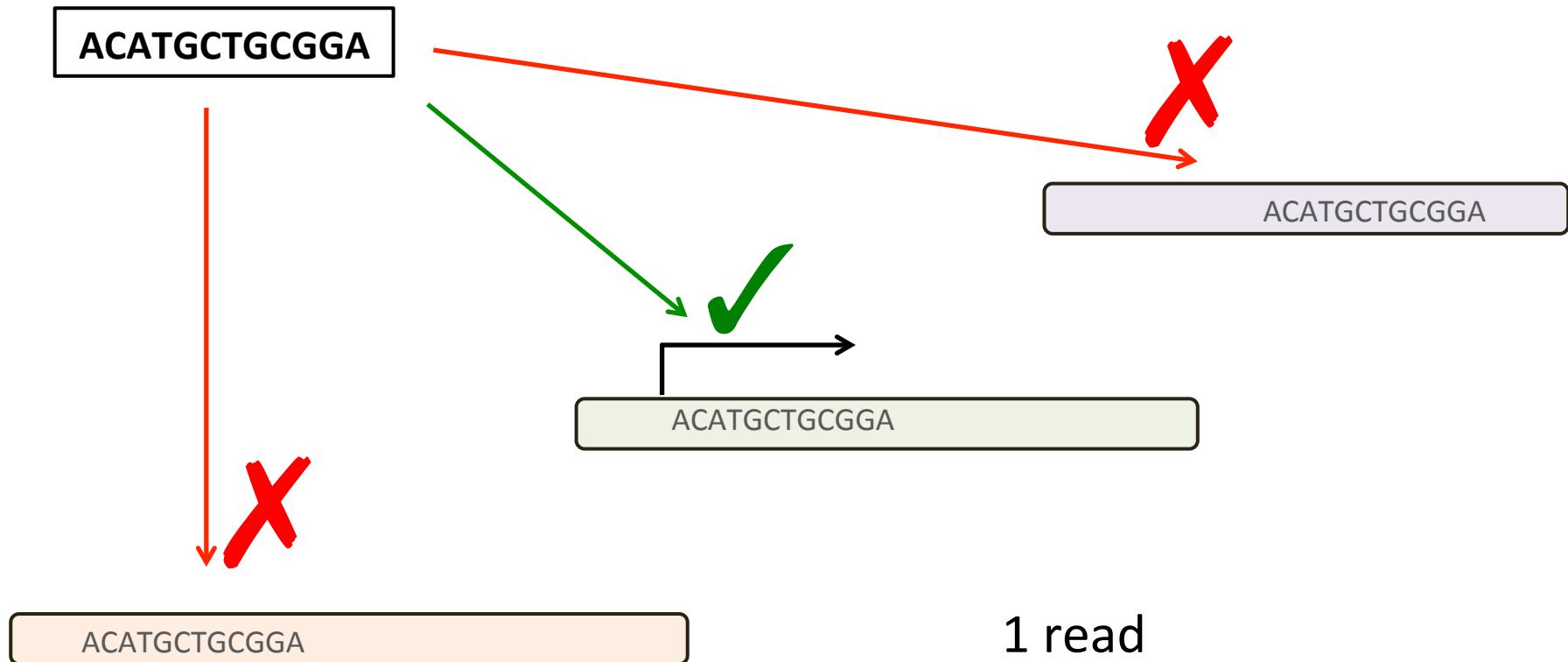
The perfect read: 1 read = 1 unique alignment.

100bp Read



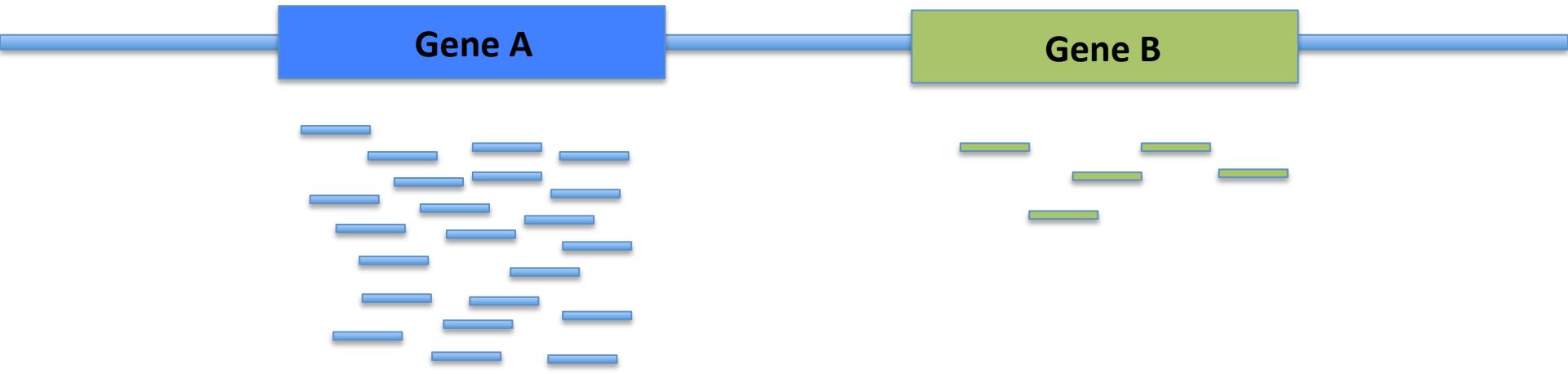
Some reads will align equally well to multiple locations. “Multireads”

100bp Read



1 read
3 valid alignments
Only 1 alignment is correct

Aligning Millions of Short Sequence Reads



Aligners: Bowtie, GSNAP, STAR, BWA, BLAT,
HISAT2, Bowtie2, Kallisto, Salmon

Align to Genome or Transcriptome?

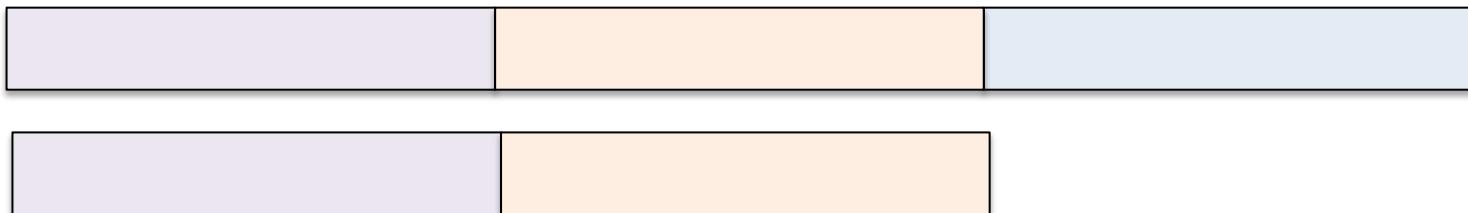
Genome



Advantages: Can align novel isoforms.

Disadvantages: Difficult, Spurious alignments, spliced alignment, gene families, pseudo genes

Transcriptome



Align to Genome or Transcriptome?

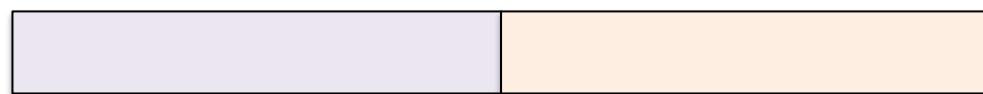
Genome



Advantages: Can align novel isoforms.

Disadvantages: Difficult, Spurious alignments, spliced alignment, gene families, pseudo genes

Transcriptome



Advantages: Easy, Focused to the part of the genome that is known to be transcribed.

Disadvantages: Reads that come from novel isoforms may not align at all or may be misattributed to a known isoform.

Visualization of alignment data (BAM/SAM)

Genome browsers – IGV and UCSC



Integrative Genome Viewer (IGV)

<http://software.broadinstitute.org/software/igv/download>

RNA-seq Data: <ftp://ftp.jax.org/dgatti/MouseGen2016/>

- DO.chr1XY.sorted.bam and DO.chr1XY.sorted.bam.bai

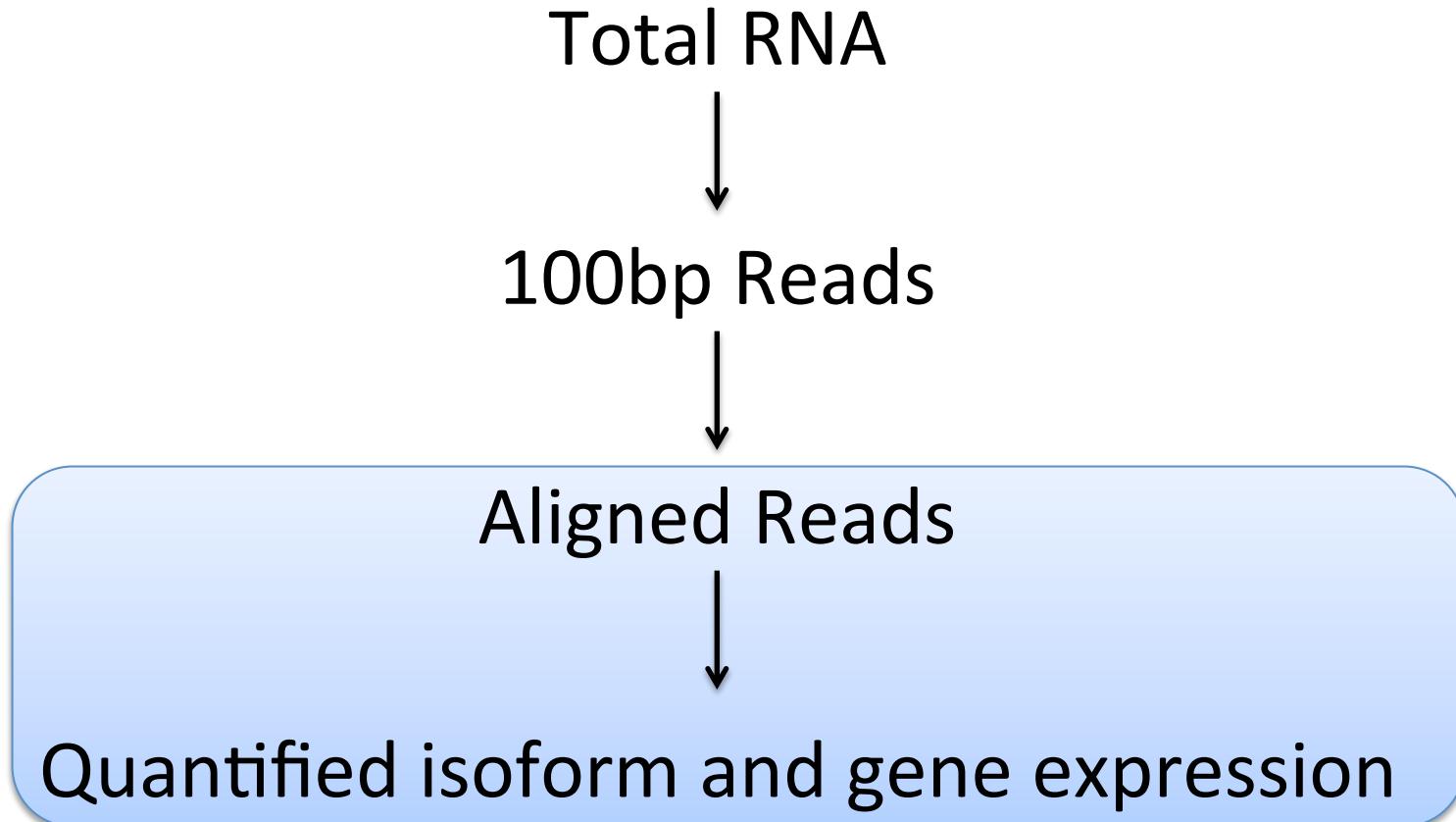
IGV is your friend.



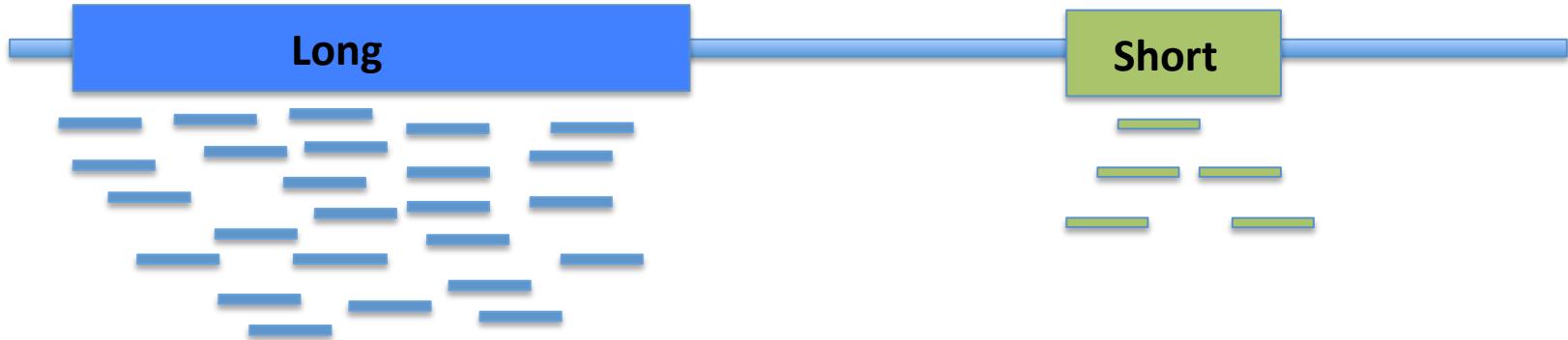
Example genes to look at in IGV

1. Tsn
2. Gorab
3. Fmo1, Fmo2, Fmo3, Fmo4, Fmo6
4. Ids
5. Zfx
6. Ssty1, Ssty2

Aligned Reads to Gene Abundance

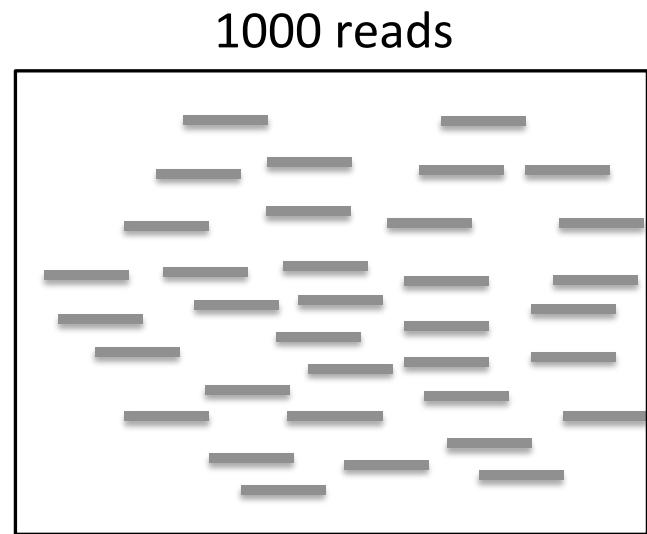
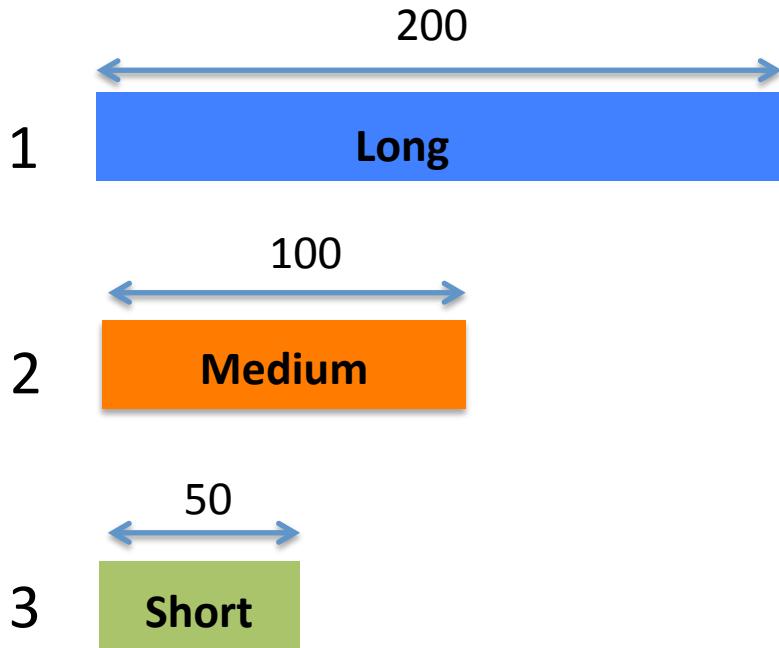


Aligned Reads to Gene Abundance: Challenges



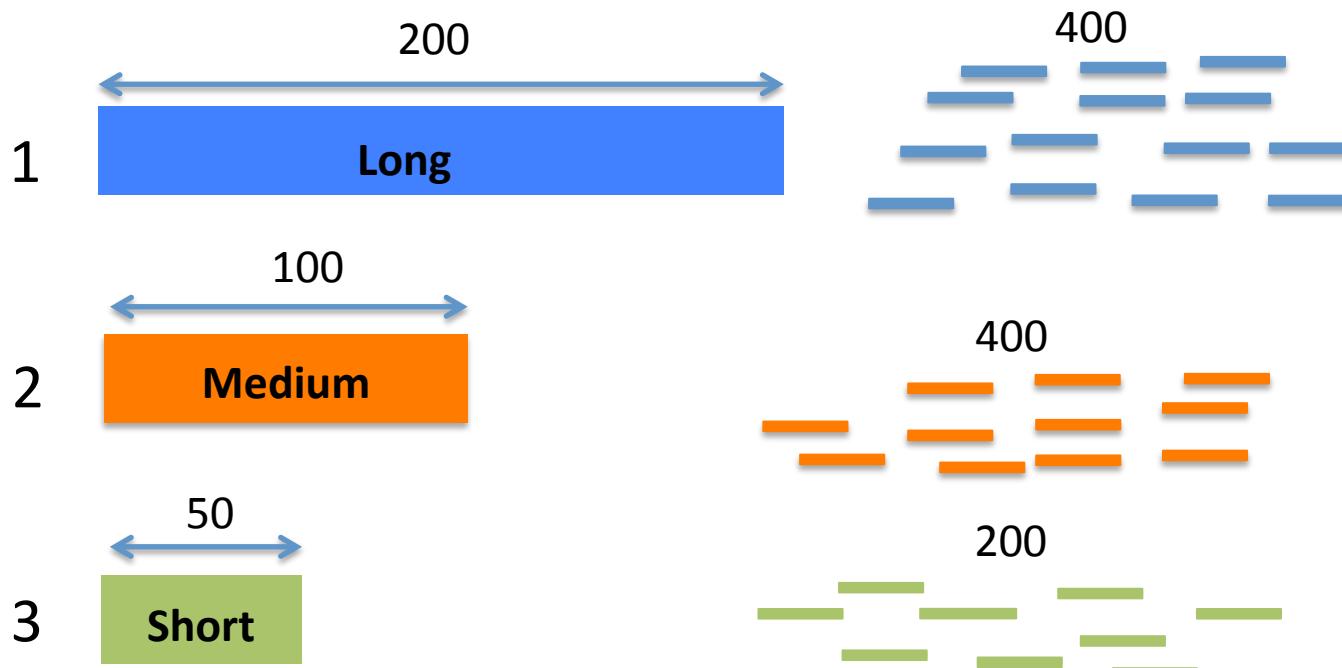
Many approaches to quantify expression abundance

Aligned Reads to Gene Abundance: Challenges



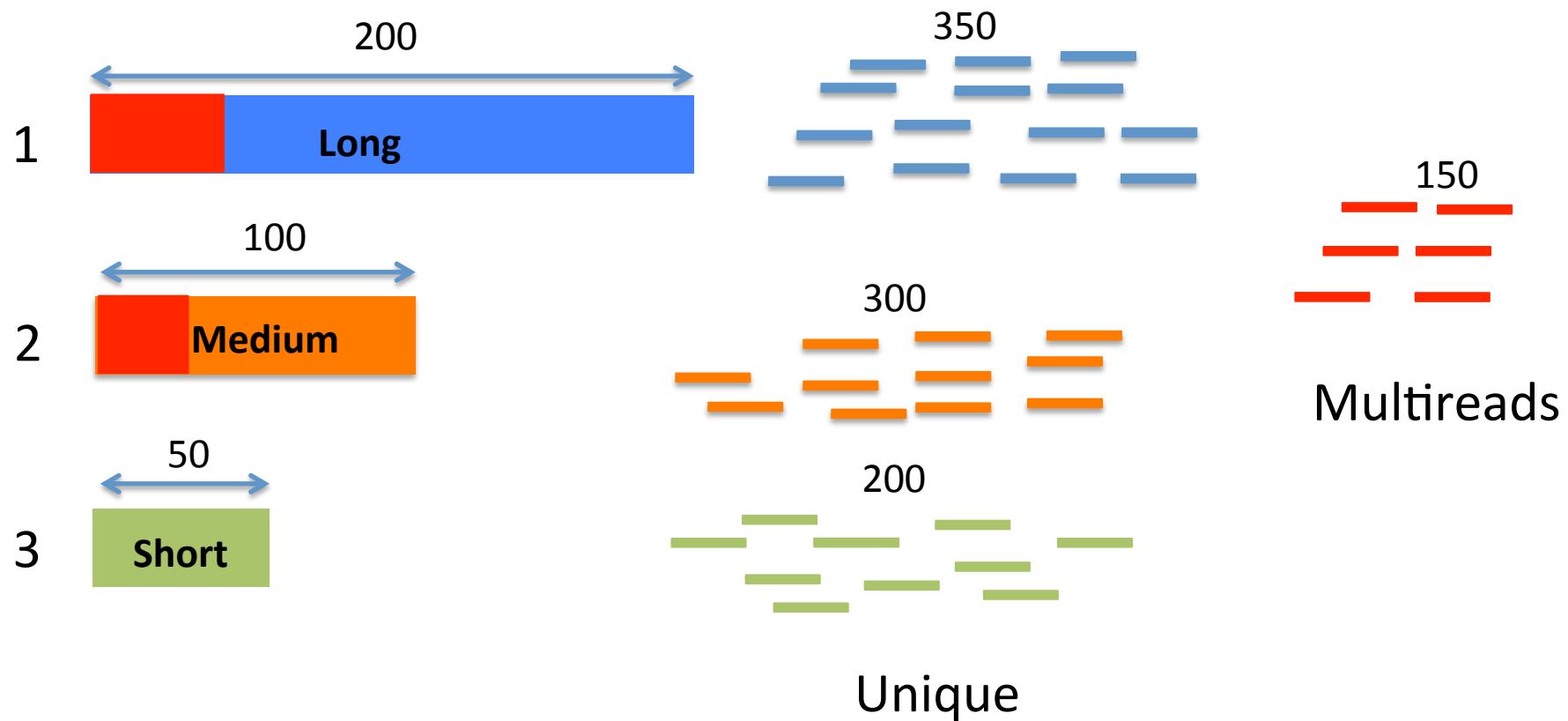
Relative abundance for these genes, f_1 , f_2 , f_3

Aligned Reads to Gene Abundance: Challenges



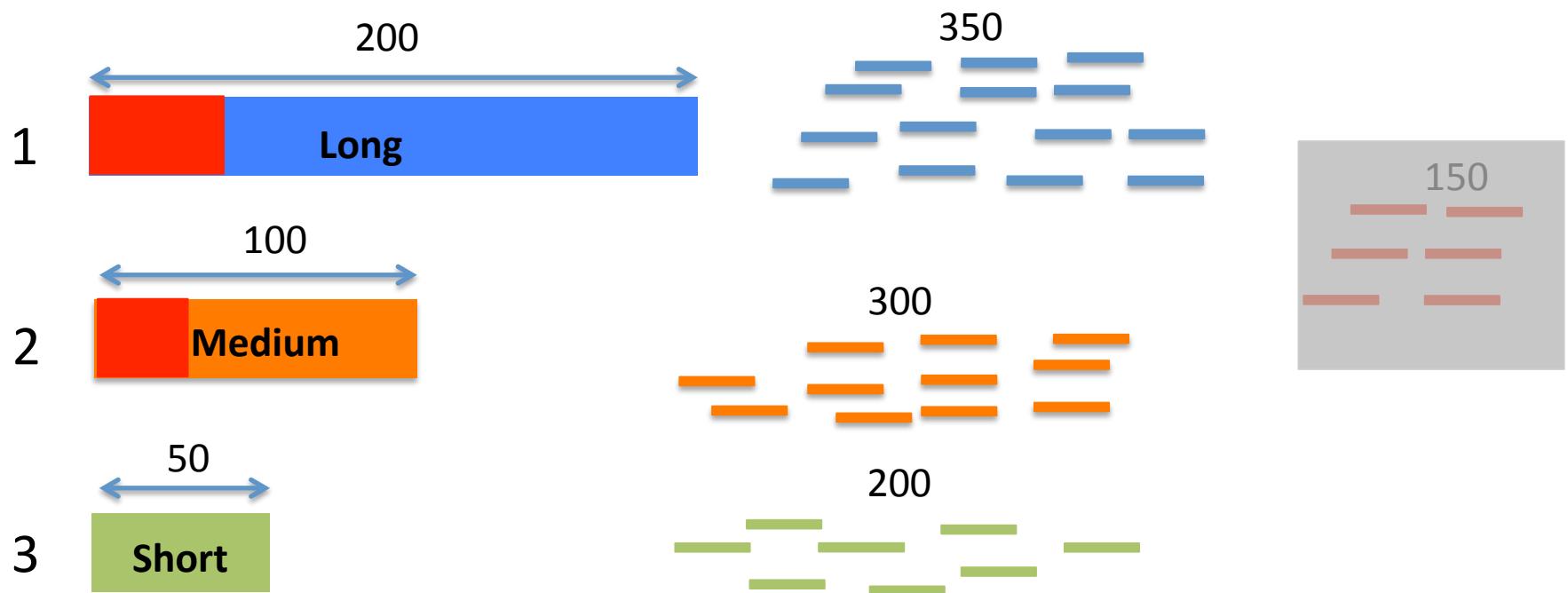
Relative abundance for these genes, f_1, f_2, f_3

Multireads: Reads Mapping to Multiple Genes/Transcripts



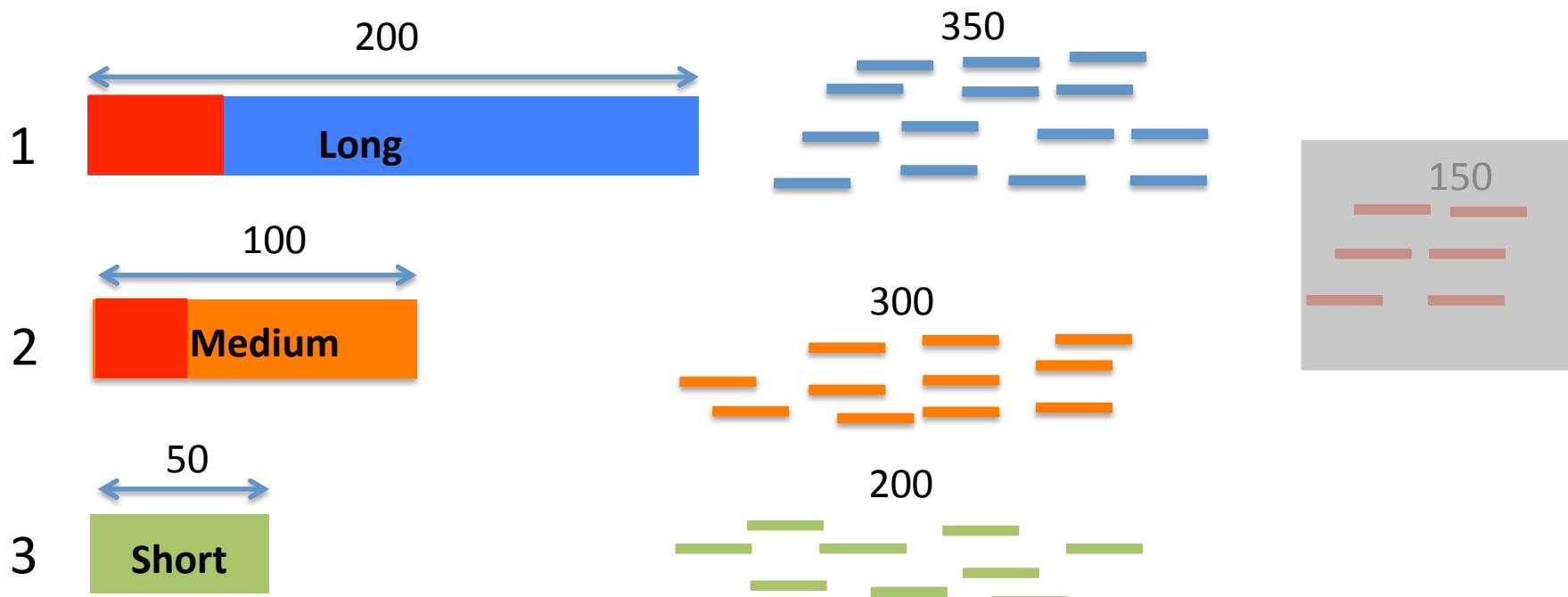
Relative abundance for these genes, f_1, f_2, f_3

Approach 1: Ignore Multireads



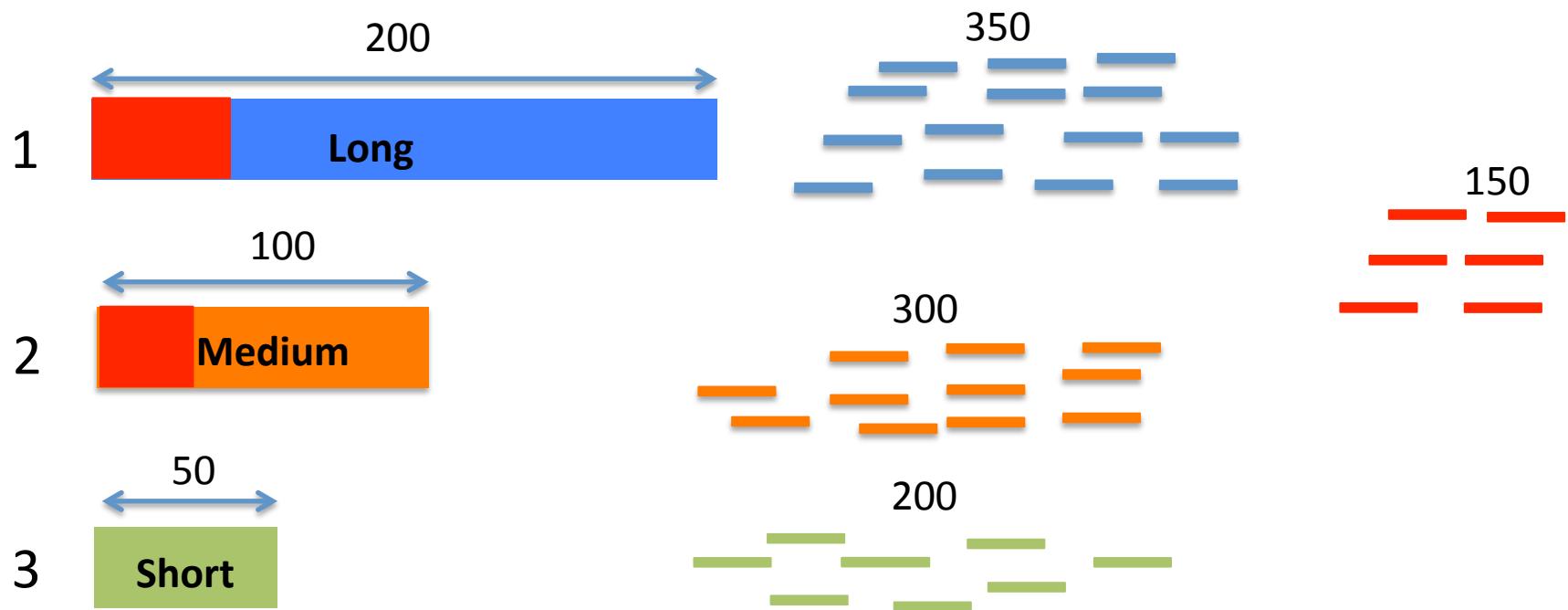
Relative abundance for these genes, f_1, f_2, f_3

Approach 1: Ignore Multireads



- Over-estimates the abundance of genes with unique reads
- Under-estimates the abundance of genes with multireads
- Not an option at all, if interested in isoform expression

Approach 2: EM algorithm based allocation of Multireads



Relative abundance for these genes, f_1, f_2, f_3

Approach 2: EM algorithm based allocation of Multireads



—

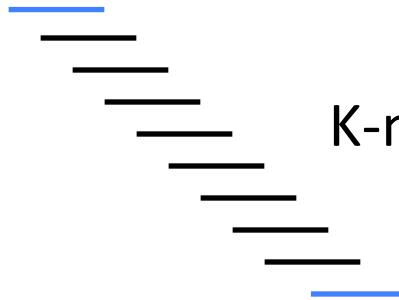
Approach 2: EM algorithm based allocation of Multireads



The rise of Pseduo-alignment a.k.a alignment-free methods

100bp Read

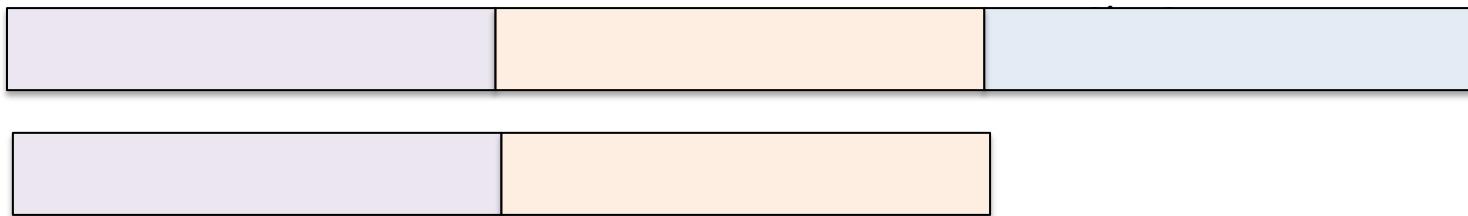
ACATGCTGCGGA



Sailfish, Salmon, and Kallisto

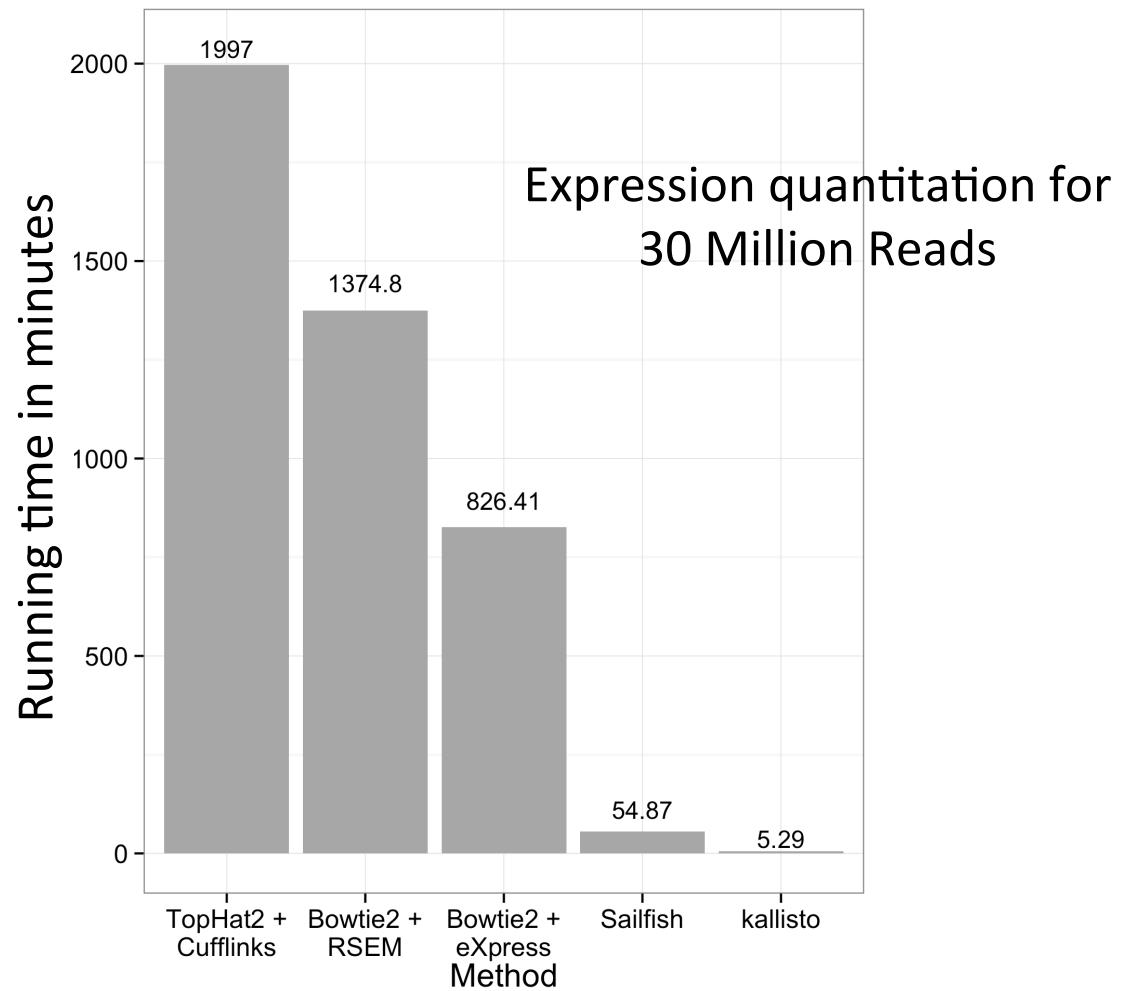
K-mers

Transcriptome



Kallisto: K-mer based pseudo-alignment

100bp Read

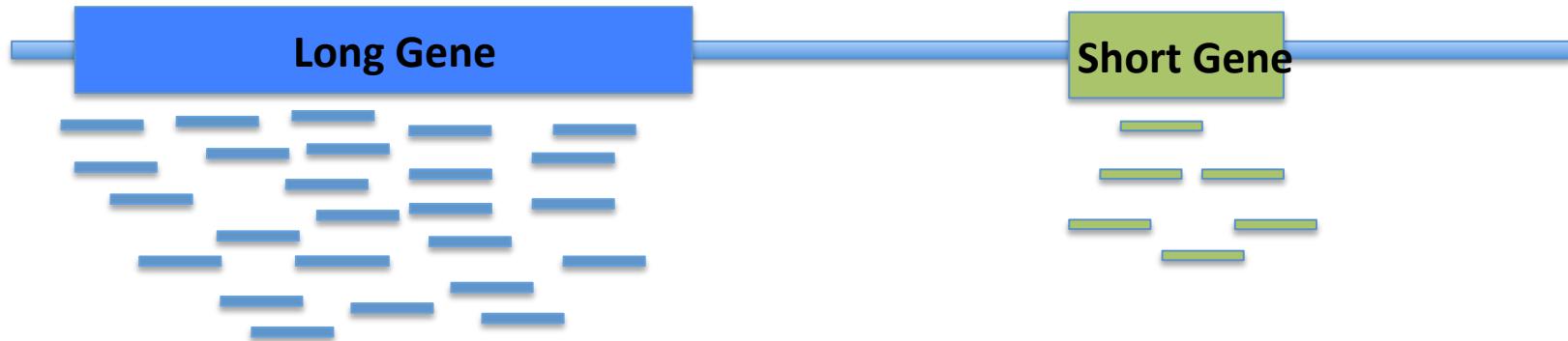


Conclusions for quantitation

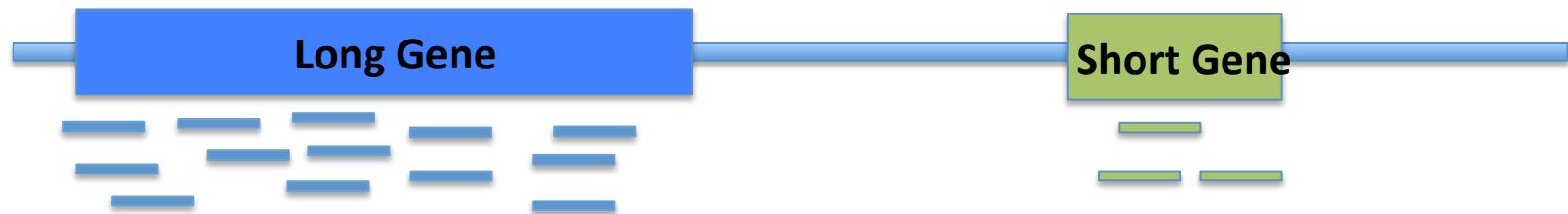
- EM approaches are currently the best option.
- Isoform-level estimates are still challenging and will become easier as read length increases.
- K-mer counting methods (Salmon, Kallisto) are very fast – they can be run easily on your own PC – and are reasonably accurate.

Expression Abundance: Counts, RPKM/FPKM, TPM

Sample 1



Sample 2



$$\text{FPKM} = \frac{\text{Number of Fragments Matched to a Gene / Kilo base}}{\text{Total matched reads in Millions}}$$

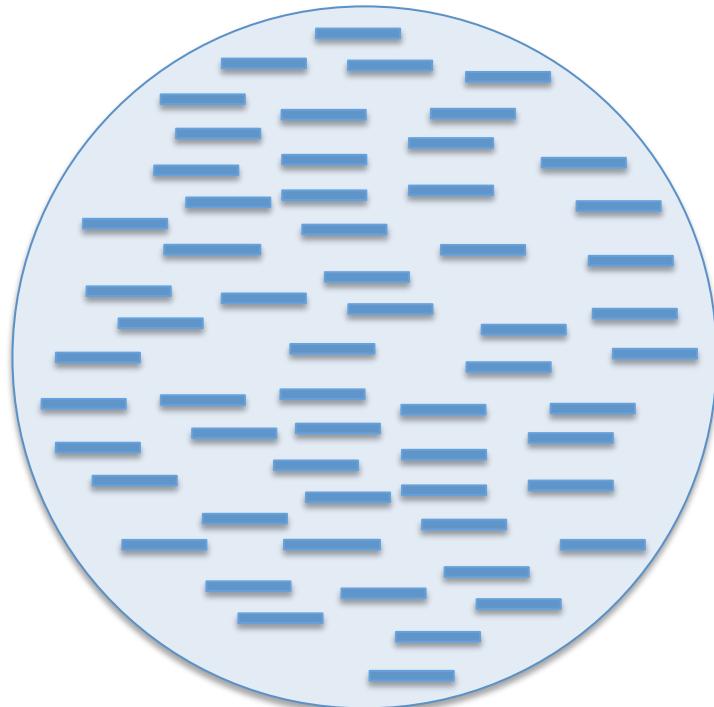
A speed bump on the road from raw counts to differential expression.

NORMALIZATION

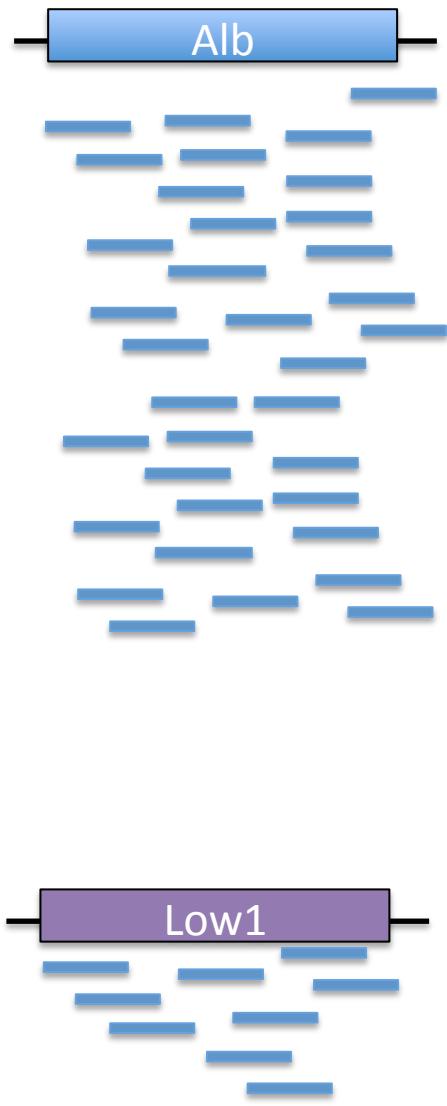
Large pool, small sample problem

- Typical RNA library estimated to contain 2.4×10^{12} molecules. McIntyre et al 2011
- Typical sequencing run = 25 million reads/sample.
- This means that only 0.00001 (1/1000th of a percent) of RNA molecules are sampled in a given run.
- High abundance transcripts are sampled more frequently.
Example: Albumin = 13% of all reads in liver RNA-seq samples.
- Sampling errors affect low-abundance transcripts most.

A finite pool of reads.

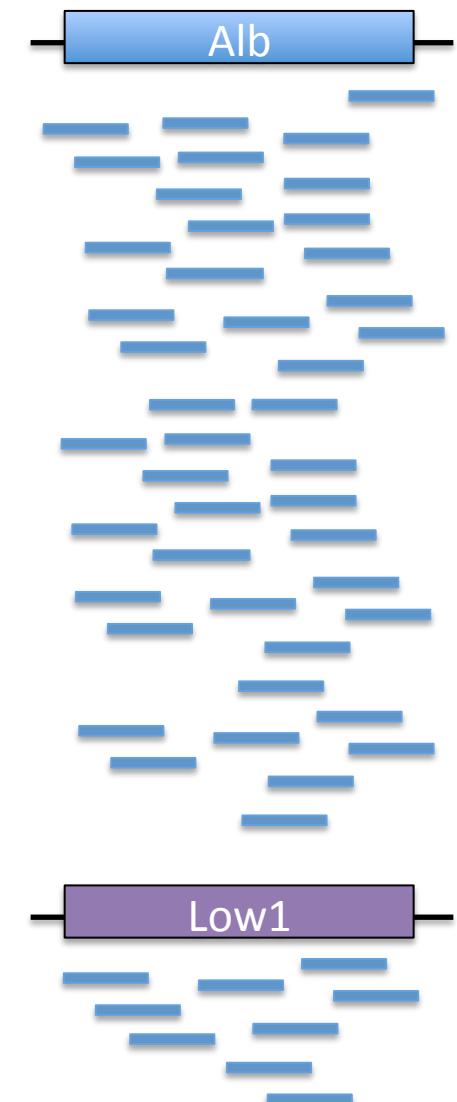


Sample 1

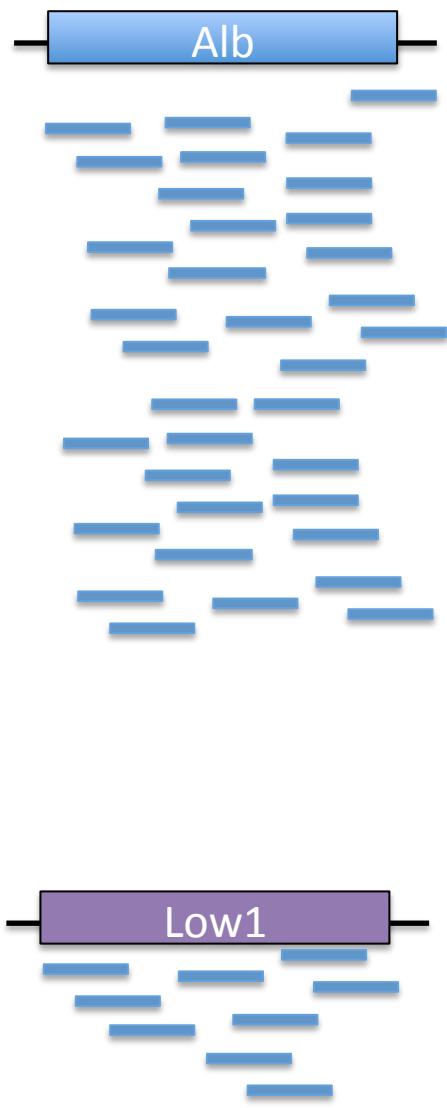


Perfect world:
All transcripts
counted.

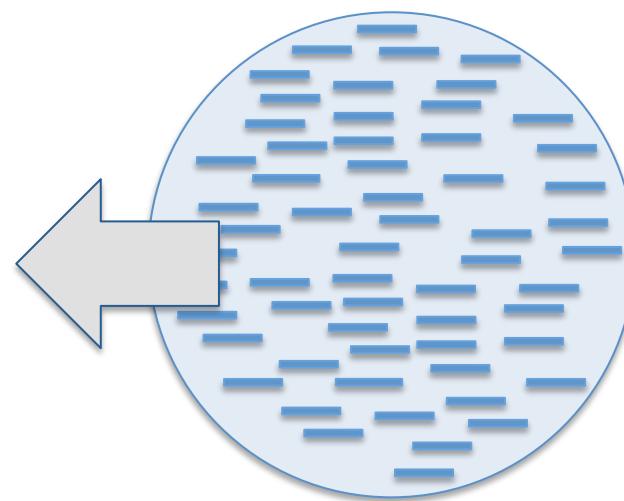
Sample 2



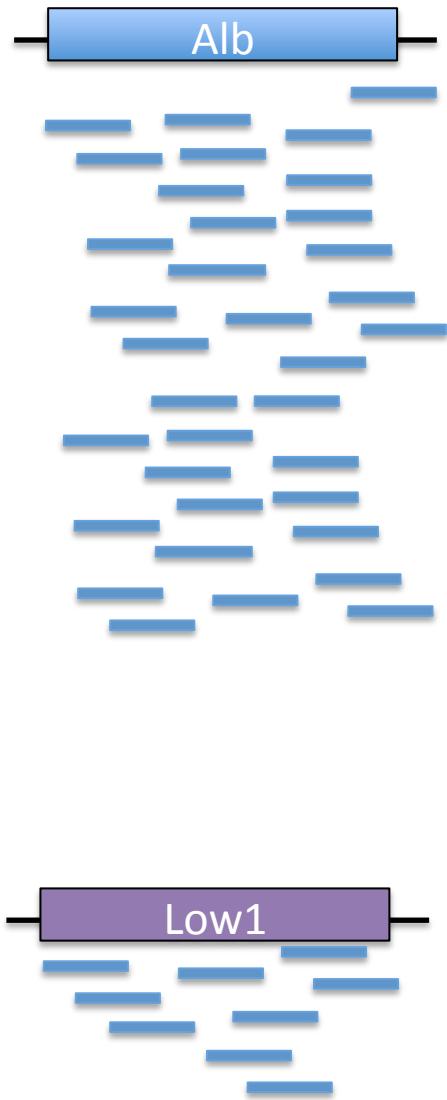
Sample 1



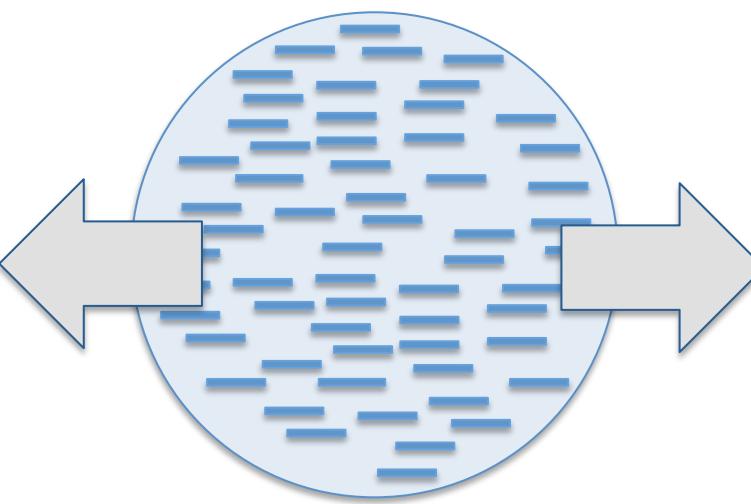
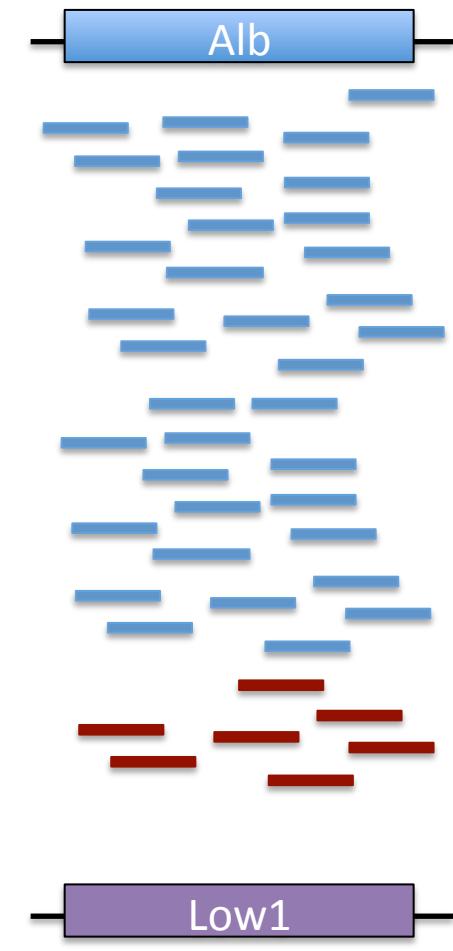
Real world: More reads taken up by highly expressed genes means less reads available for lowly expressed genes.



Sample 1



Sample 2

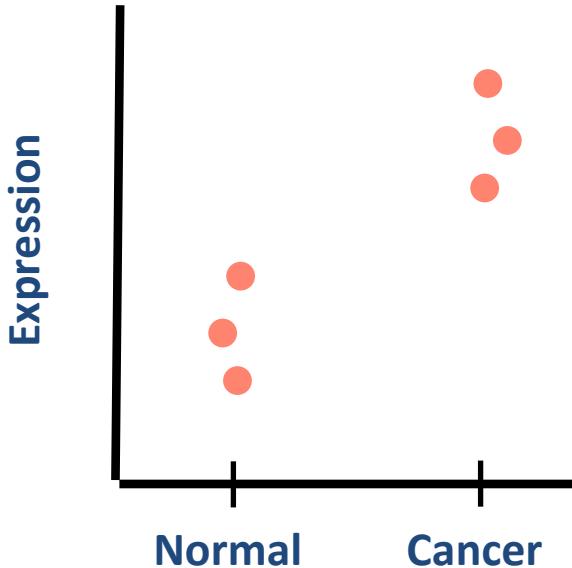


Highly expressed genes that are differentially expressed can cause lowly expressed genes that are not actually differentially expressed to appear that way.

Normalization of raw counts

- Wrong way to normalize data
 - Normalizing to the total number of mapped reads (e.g. FPKM). Top 10 highly expressed genes soak up 20% of reads in the liver. FPKM is widely used, and problematic.
- Better ways to measure data
 - Normalize to upper quartile (75th %) of non-zero counts, median of scaled counts (DESeq), or the weighted trimmed mean of the log expression ratios (EdgeR).

Differential Expression Analysis



T-test

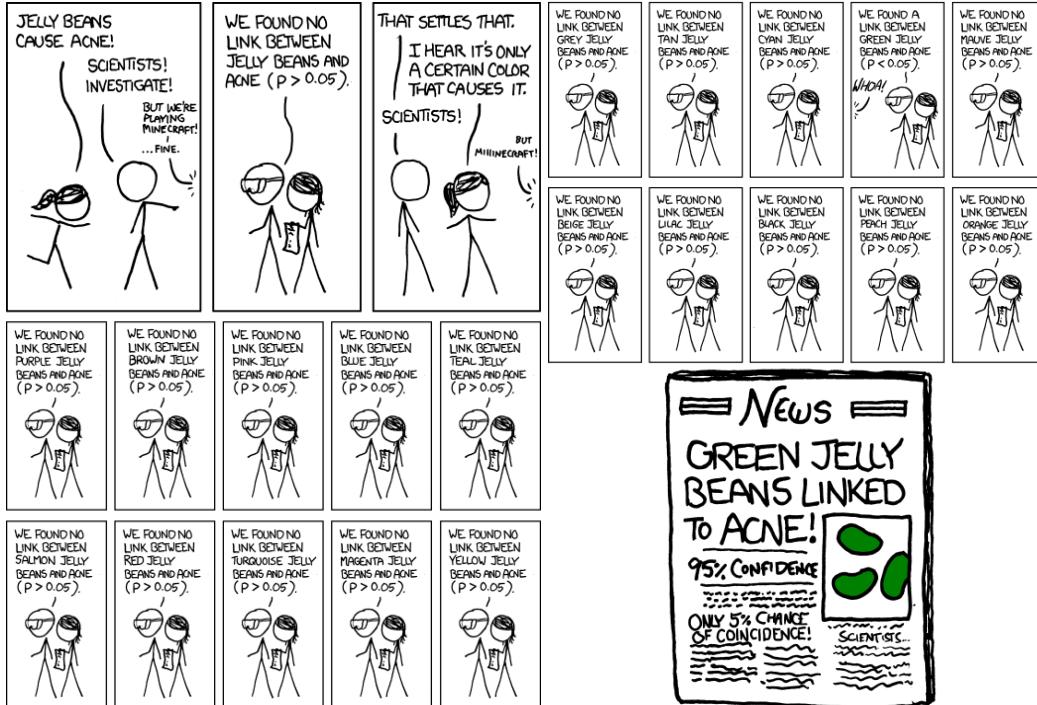
$$t_g = \frac{\hat{\mu}_{g,1} - \hat{\mu}_{g,2}}{\sqrt{\frac{\hat{\sigma}_{g,1}^2}{N_1} + \frac{\hat{\sigma}_{g,2}^2}{N_2}}}$$

Over-estimation of $\hat{\sigma}_g^2$ Too conservative

Under-estimation of $\hat{\sigma}_g^2$ Too sensitive
(Many false positives)

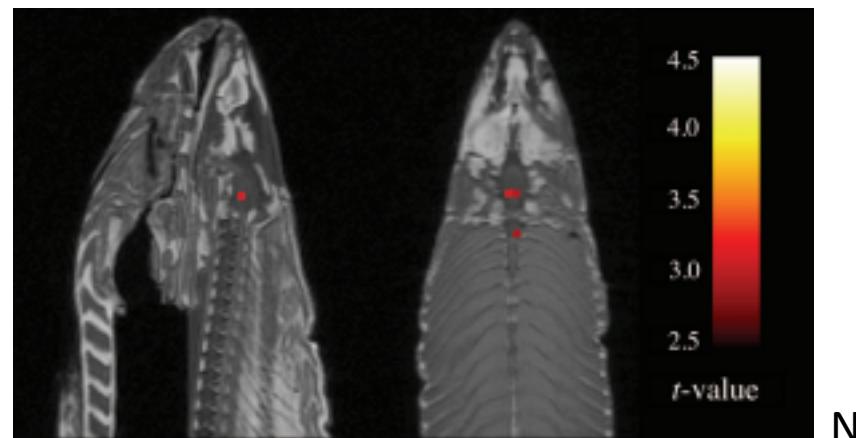
DESEQ2, edgeR, Voom, & CuffDiff

Multiple Testing Correction and False Discovery rate

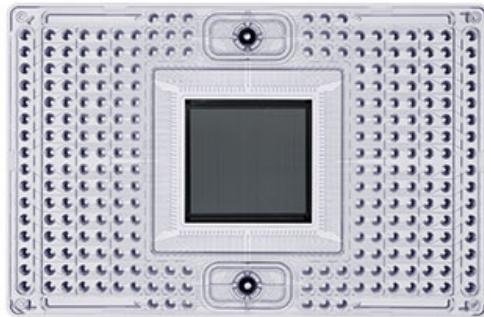


XKCD Significant

2012 IgNobel prize in
Neuroscience for “finding
Brain activity signal in dead salmon using fMRI”

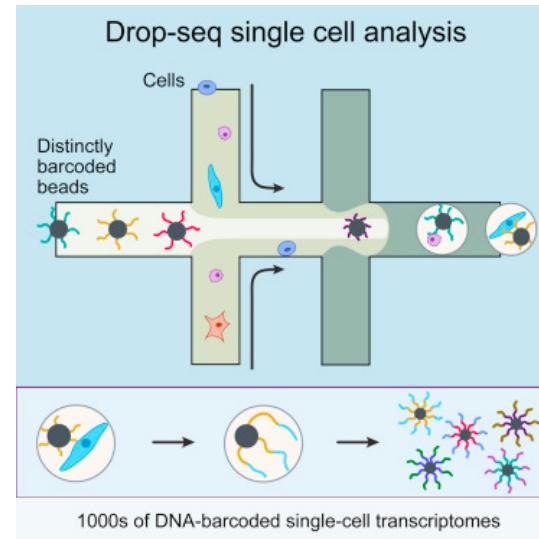


Single Cell RNA-seq Technologies



Fluidigm C1 Chip
96 cells / 800 Cells

DropSeq: 40,000 cells



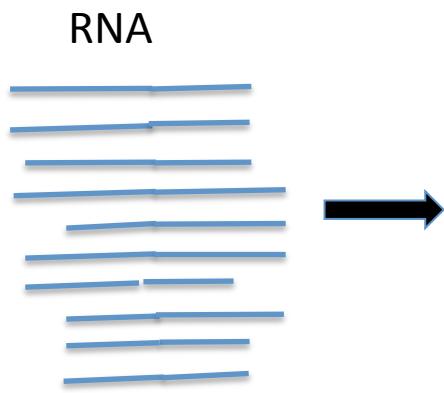
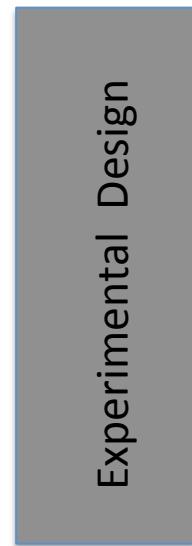
10X Genomics:
48,000 cells



Summary



Summary



ATGCTCA AGCTA
TAGATGCTCAAGCTA
ATGCTCAAGCTAAC
ATGCTCAAGCTA
AGTAGATGCTCAAGCTA
ATGCTCAAGCTA
ATGCTCA AGCTA
ATGCTCAAGCTA
TAGATGCTCAAGCTAAC
CTCAAGCTAACCTAG

RNA-seq analysis pipeline

As sequences get longer, alignment and isoform quantitation becomes easier!

Resources

Aligner

- Bowtie 2 <http://bowtie-bio.sourceforge.net/bowtie2/index.shtml>
- GSNAP <http://research-pub.gene.com/gmap/>

Transcript Discovery/Annotation

- STAR <https://github.com/alexdobin/STAR/releases>
- Tophat <http://tophat.cbcn.umd.edu/>

Transcript Abundance

- Kallisto <http://pachterlab.github.io/kallisto/>
- RSEM <http://deweylab.biostat.wisc.edu/rsem/>
- EMASE <https://github.com/churchill-lab/emase>

Differential Expression

- DESeq <http://www-huber.embl.de/users/anders/DESeq/>
- edgeR <http://bioconductor.org/packages/release/bioc/html/edgeR.html>
- EBSeq <https://www.biostat.wisc.edu/~kendzior/EBSEQ/>

Example 1

Differential expression in my mutant mouse compared to wild type. What genes are up- or down-regulated?

Things to consider...

- Differential expression of highly expressed and well annotated genes?
 - Smaller sample depth; more biological replicates
 - No need for paired end reads; shorter reads (50bp) may be sufficient.
 - Better to have 20 million 50bp reads than 10 million 100bp reads.
- Looking for novel genes/splicing/isoforms?
 - More read depth, paired-end reads from longer fragments.

Example 2

- How to quantify gene expression in a species that has not been sequenced or annotated?
 - Multistep strategy using multiple sequencing technologies.

Example 3

- How to quantify single cell gene expression in a heterogeneous human tumor?

Any other applications you are
interested in?

Steve Munger
steven.munger@jax.org

Narayanan Raghupathy
narayanan.raghupathy@jax.org

Acknowledgements

- KB Choi
- Gary Churchill
- Ron Korstanje/ Karen Svenson/ Elissa Chesler
- Joel Graber
- Doug Hinerfeld
- Anuj Srivastava
- Churchill Lab – Dan Gatti
- Al Simons and Matt Hibbs
- Lisa Somes, Steve Ciciotte, mouse room staff at JAX
- Gene Expression Technologies group at JAX