

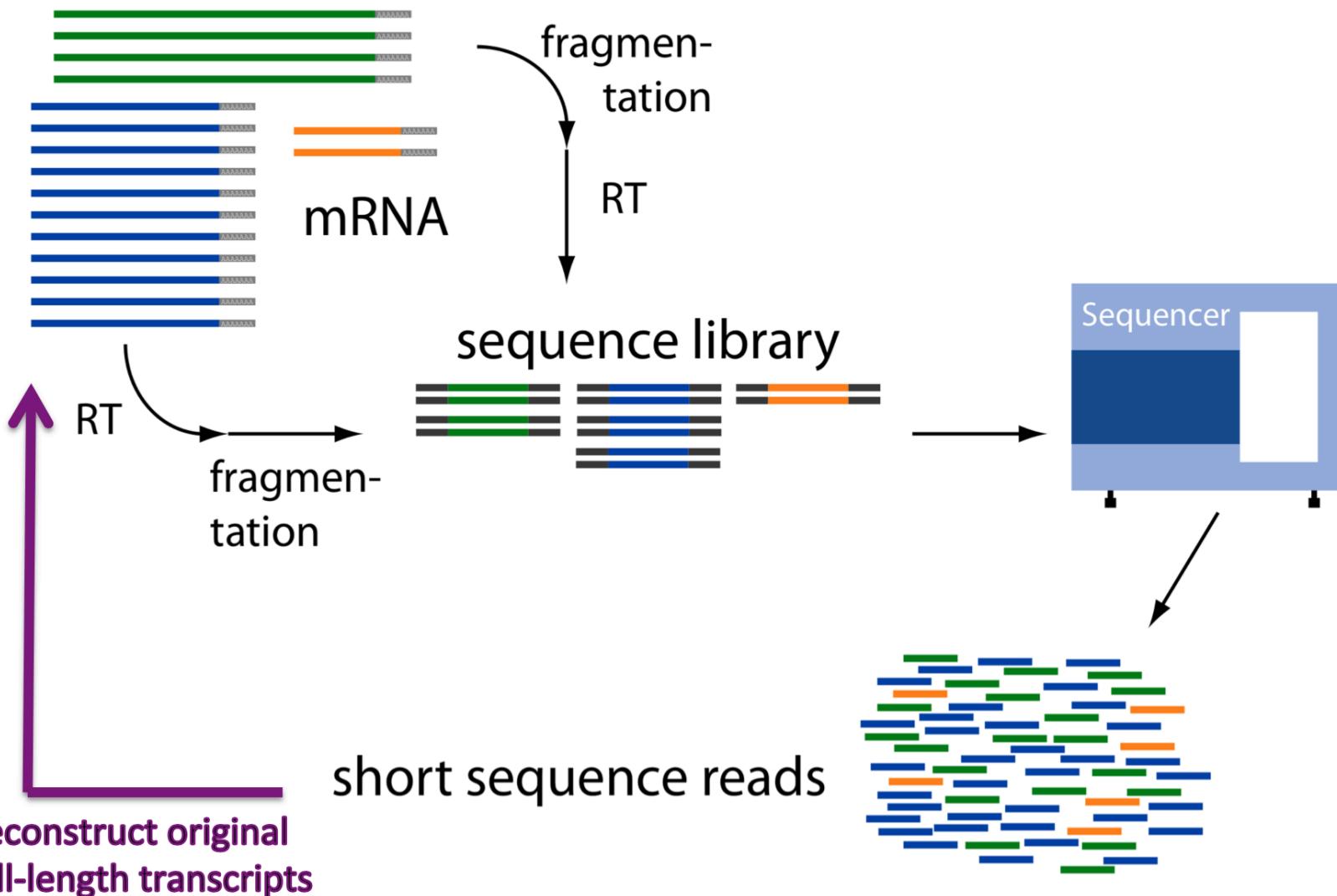
BI694

Bioinformatics & Phylogenetics

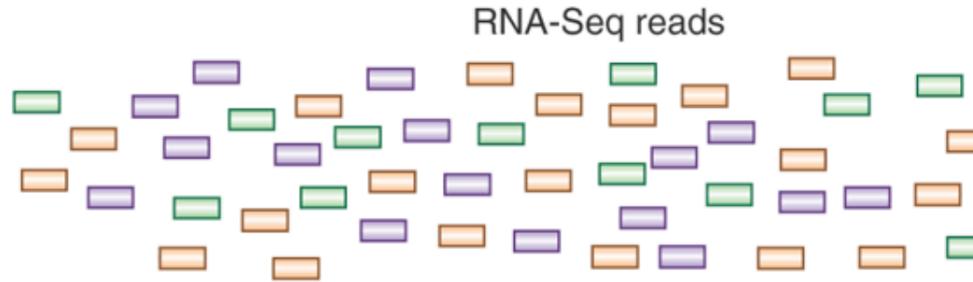
Winter Semester 2017

Read Mapping and DE Analysis Review

Overview of RNA-Seq



Transcript Reconstruction from RNA-Seq Reads



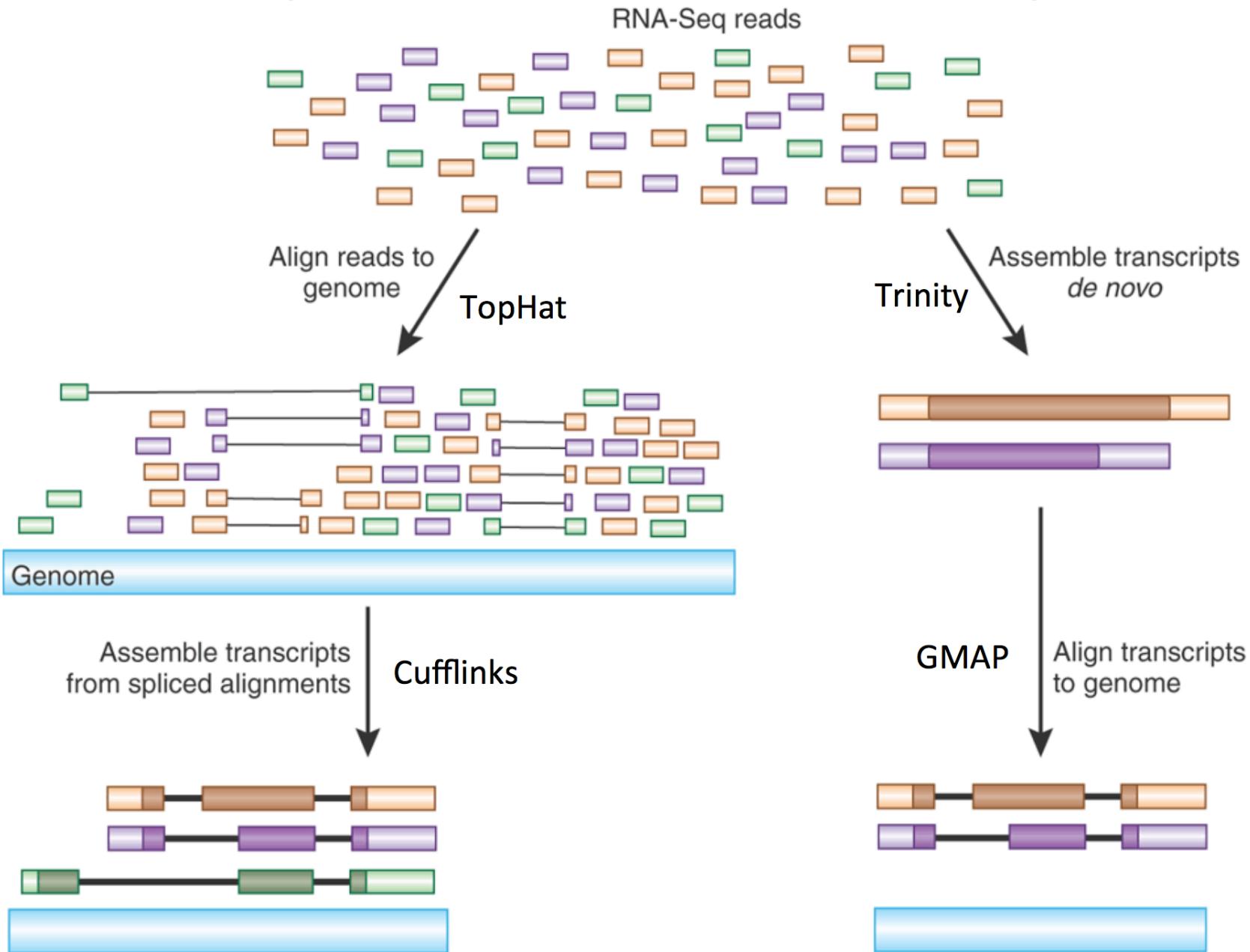
Advancing RNA-Seq analysis

Brian J Haas & Michael C Zody

Nature Biotech, 2010

New methods for analyzing RNA-Seq data enable *de novo* reconstruction of the transcriptome.

Transcript Reconstruction from RNA-Seq Reads



De novo transcriptome assembly

No genome required

Empower studies of non-model organisms

- expressed gene content
- transcript abundance
- differential expression

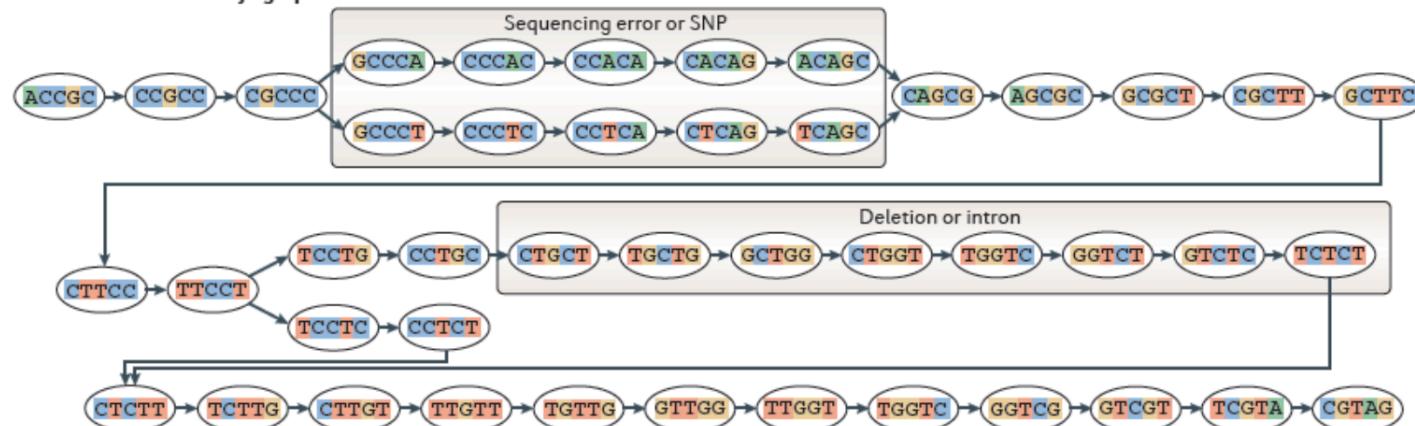
The General Approach to
De novo RNA-Seq Assembly
Using De Bruijn Graphs

Sequence Assembly via De Bruijn Graphs

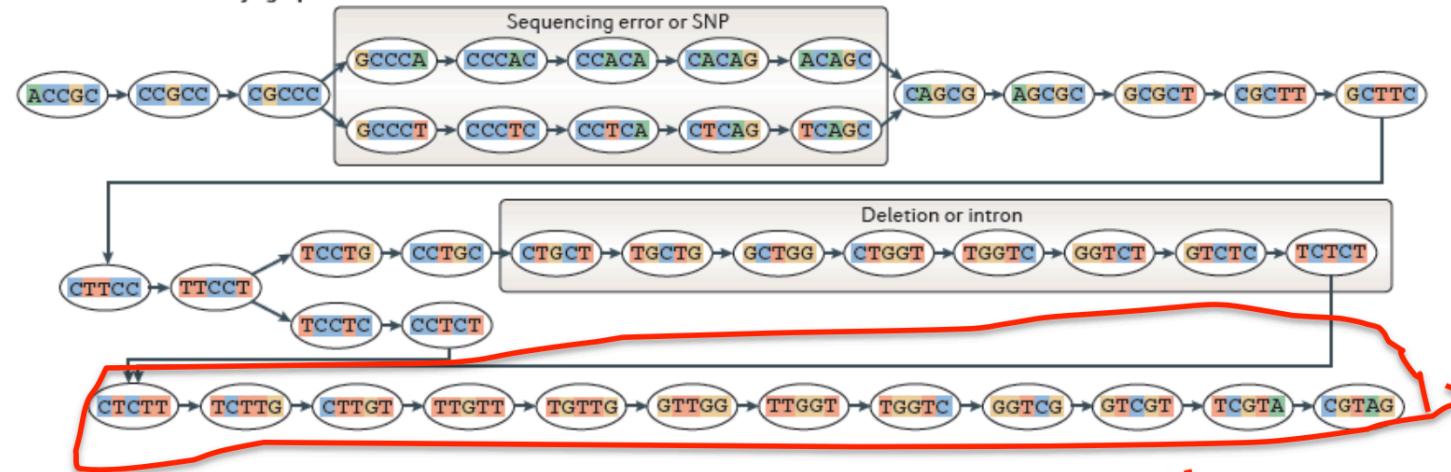
a Generate all substrings of length k from the reads

ACAGC	TCCTG	GTCTC		AGCGC	CTCTT	GGTCG	k-mers (k=5)	
CACAG	TTCCCT	GGTCT		CAGCG	CCTCT	TGGTC		
CCACA	CTTCC	TGGTC	TGTTG	TCAGC	TCCTC	TTGGT		
CCCAC	GCTTC	CTGGT	TTGTT	CTCAG	TTCCCT	GTTGG		
GCCCA	CGCTT	GCTGG	CTTGT	CCTCA	CTTCC	TGTTG		
CGCCC	GCGCT	TGCTG	TCTTG	CCCTC	GCTTC	TTGTT		CGTAG
CCGCC	AGCGC	CTGCT	CTCTT	GCCCT	CGCTT	CTTGT		TCGTA
ACCGC	CAGCG	CCTGC	TCTCT	CGCCC	GCGCT	TCTTG		GTCTG
ACCGCCCCACAGCGCTTCCCTGCTGGTCTCTTGGTG				CGCCCTCAGCGCTTCCCTTGGTGGTCGTAG				Reads

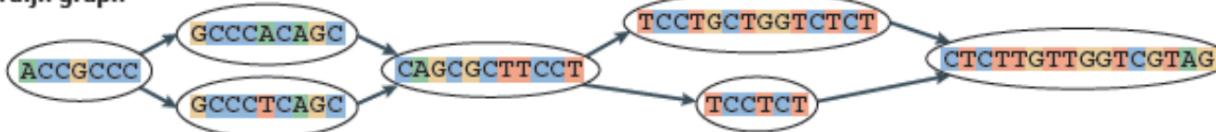
b Generate the De Bruijn graph



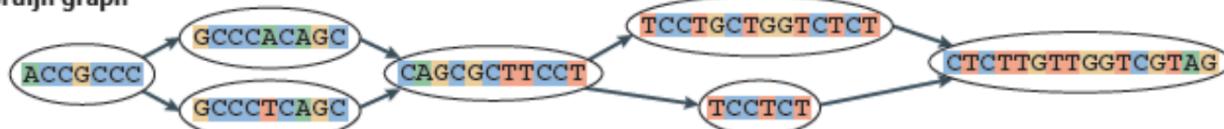
b Generate the De Bruijn graph



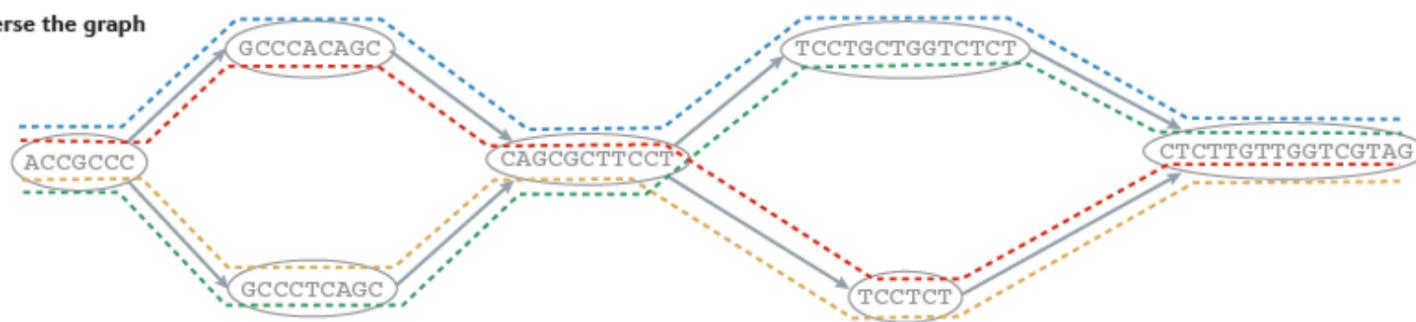
c Collapse the De Bruijn graph



c Collapse the De Bruijn graph



d Traverse the graph



e Assembled isoforms

----- ACCGCCCCAAGCGCTTCCTGCTGGTCTTTGGTGGTCGTAG
----- ACCGCCCCAAGCGCTTCCT ----- CTTGGTGGTCGTAG
---- ACCGGCCCTCAGCGCTTCCT ----- CTTGGTGGTCGTAG
---- ACCGGCCCTCAGCGCTTCCTGCTGGTCTTTGGTGGTCGTAG

Contrasting Genome and Transcriptome Assembly

Genome Assembly

- Uniform coverage
- Single contig per locus
- Double-stranded

Transcriptome Assembly

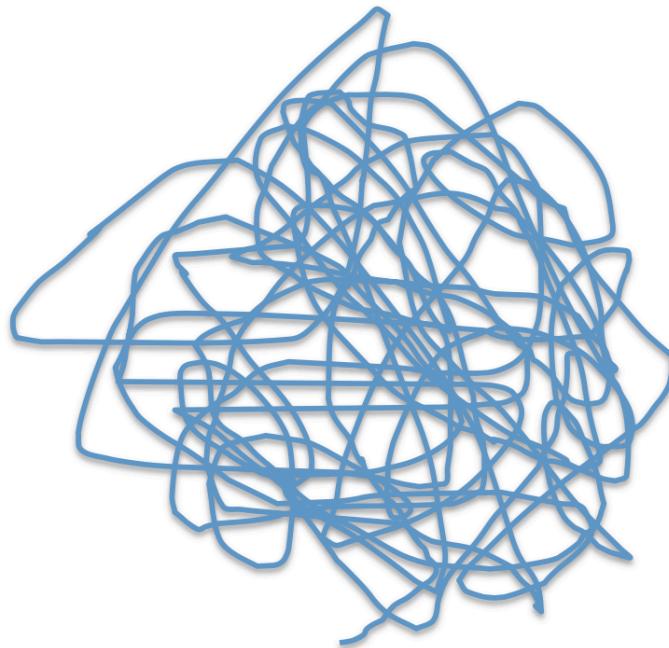
- Exponentially distributed coverage levels
- Multiple contigs per locus (alt splicing)
- Strand-specific



Trinity Aggregates Isolated Transcript Graphs

Genome Assembly

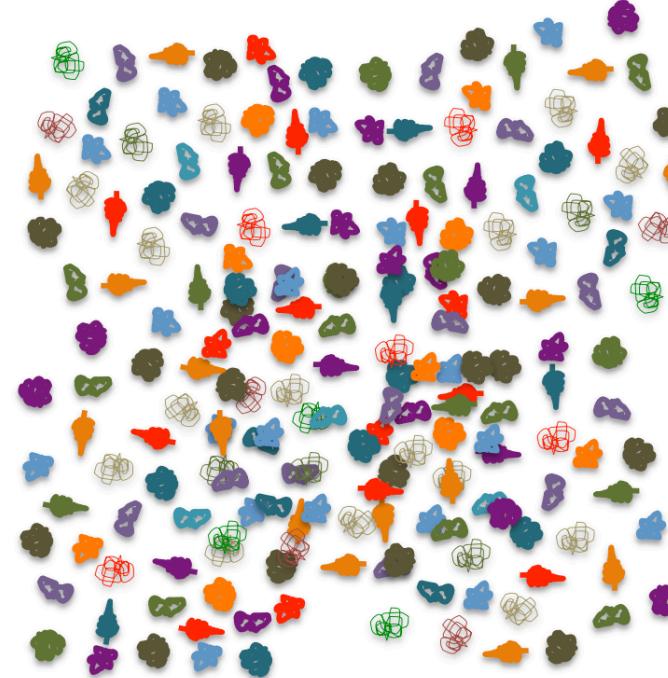
Single Massive Graph



Entire chromosomes represented.

Trinity Transcriptome Assembly

Many Thousands of Small Graphs



Ideally, one graph per expressed gene.

Trinity output: A multi-fasta file

Abundance Estimation

(Aka. Computing Expression Values)

How abundant is a transcript?

GGCGTCTATATCTGGCTCTAGGCCCTCATTTTT Transcript

GGCGTCTATATCT

GGCGTCTATATCTCG

TATCTCGGCTCTAGG

TATCTCAGCTCTAGGCC

TATCTCAGCTCTAGGCCCTCA

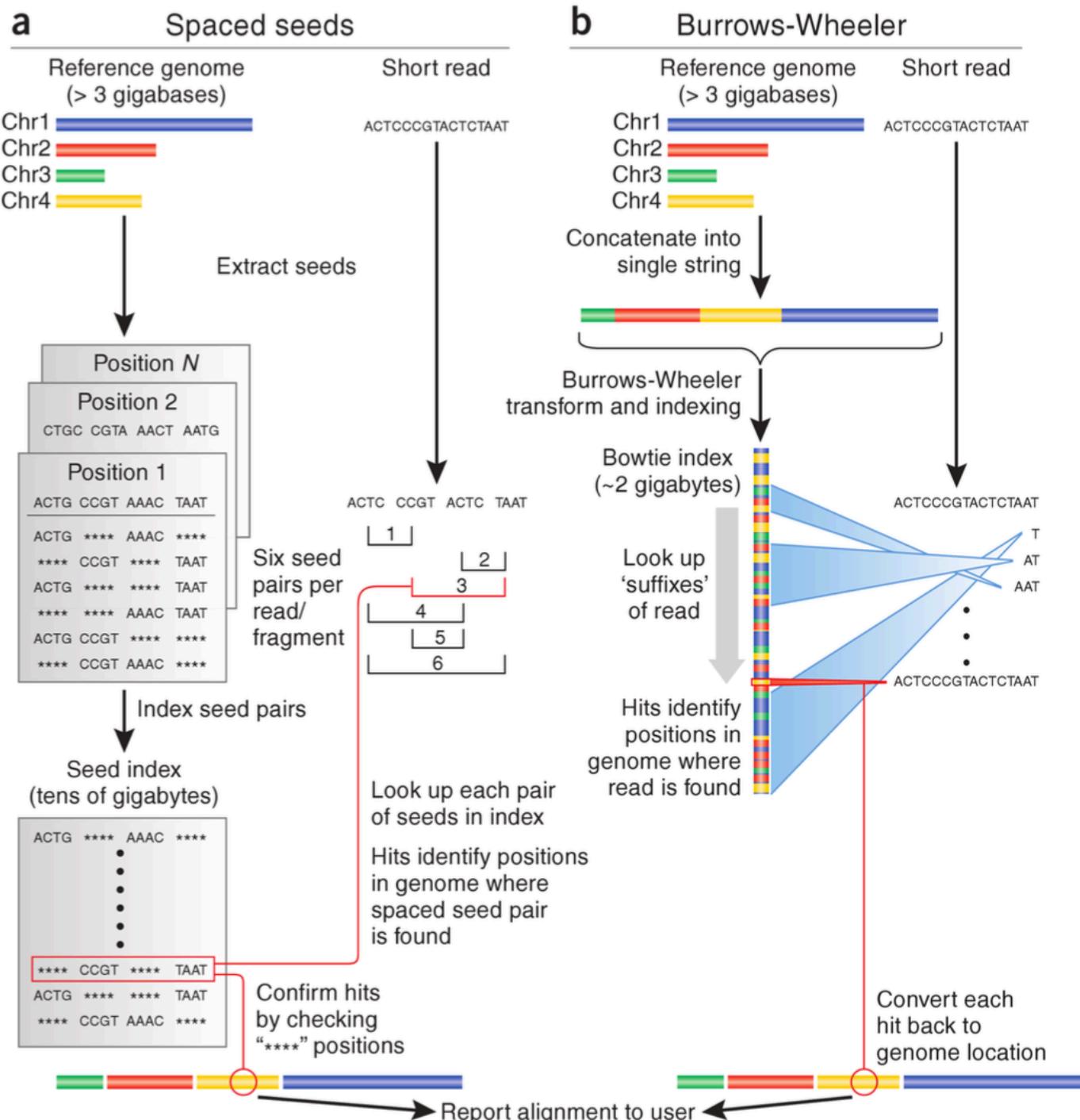
CTCGGCTCTAGGCCCTCATTTT

GGCTCTAGGCCCTCATTTTTT

CTCTAGGCCCTCATTTTTT

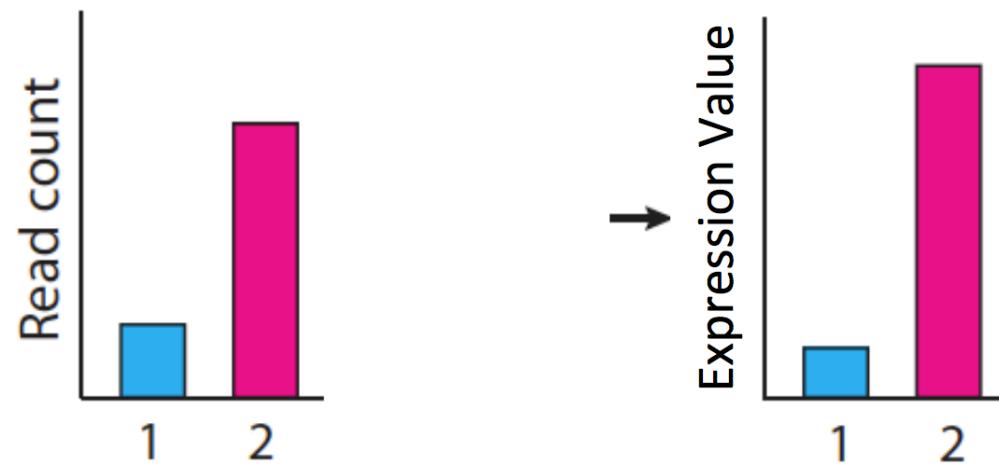
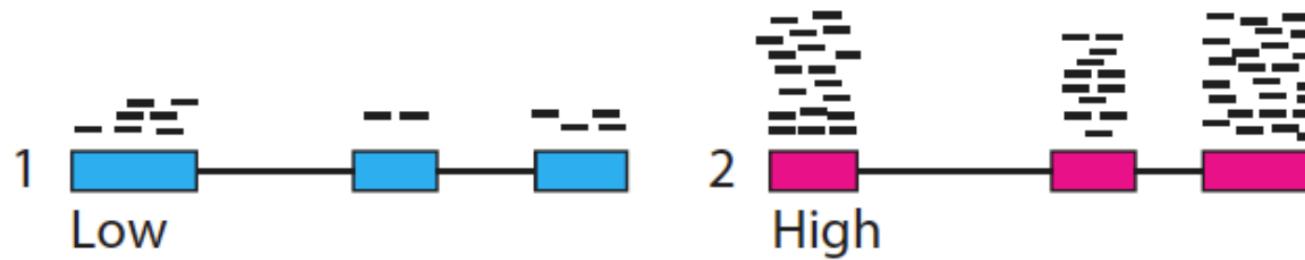
CTAGGCCCTCATTTTTT

Mapped reads

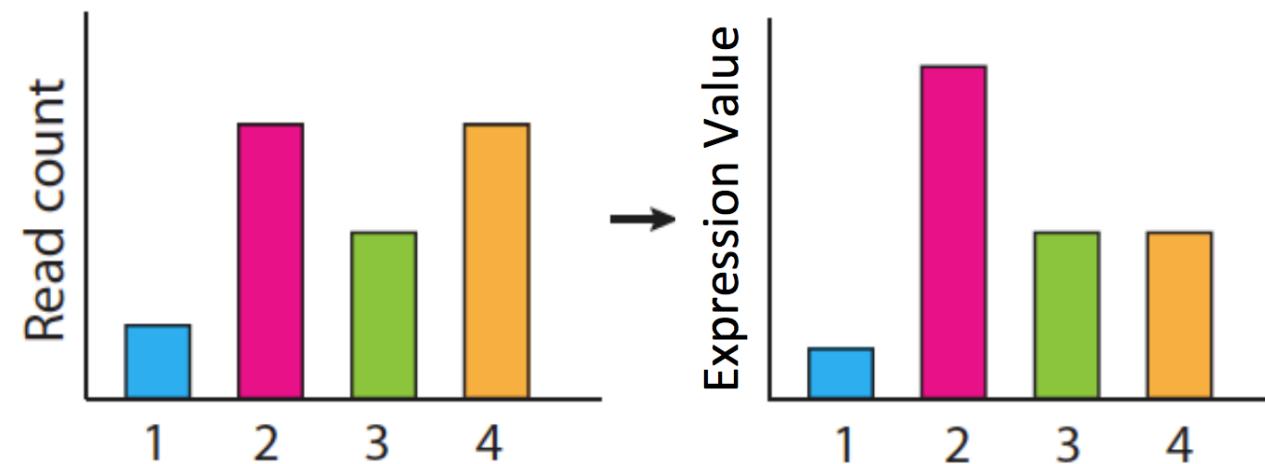
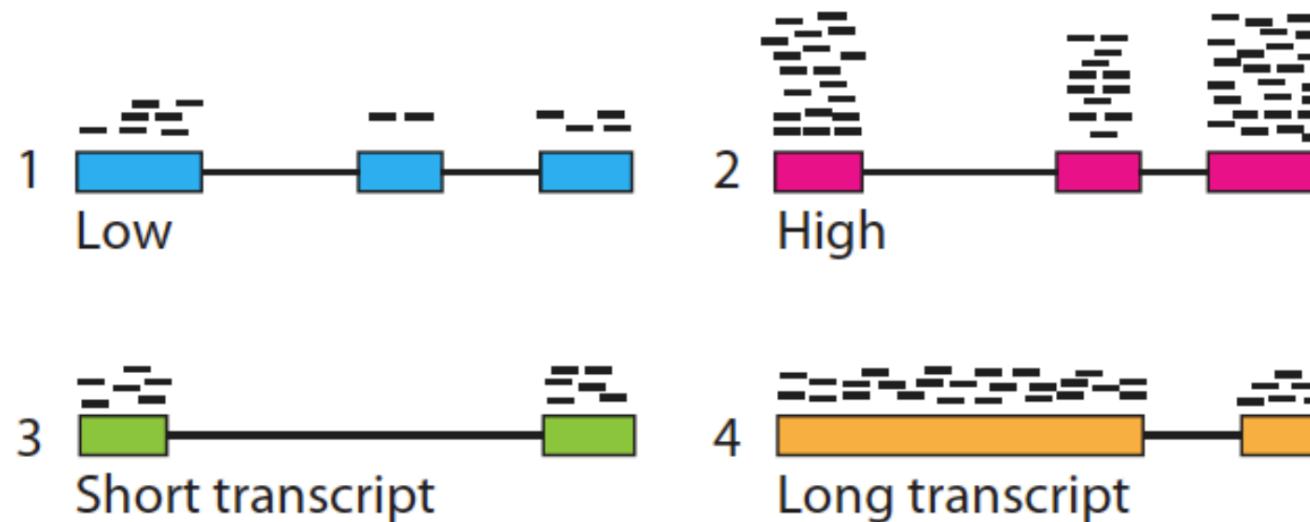


What is a Burrows Wheeler Transform?

Calculating expression of genes and transcripts



Calculating expression of genes and transcripts

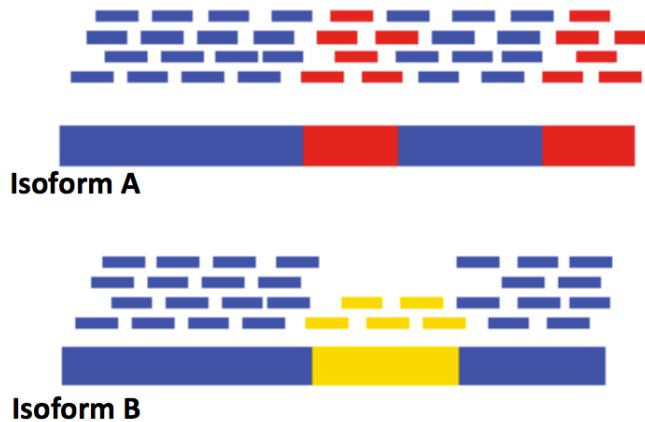


Normalized Expression Values

- Transcript-mapped read counts are normalized for both length of the transcript and total depth of sequencing.
- Reported as: Number of RNA-Seq **F**ragments
Per **K**ilobase of transcript
per total **M**illion fragments mapped

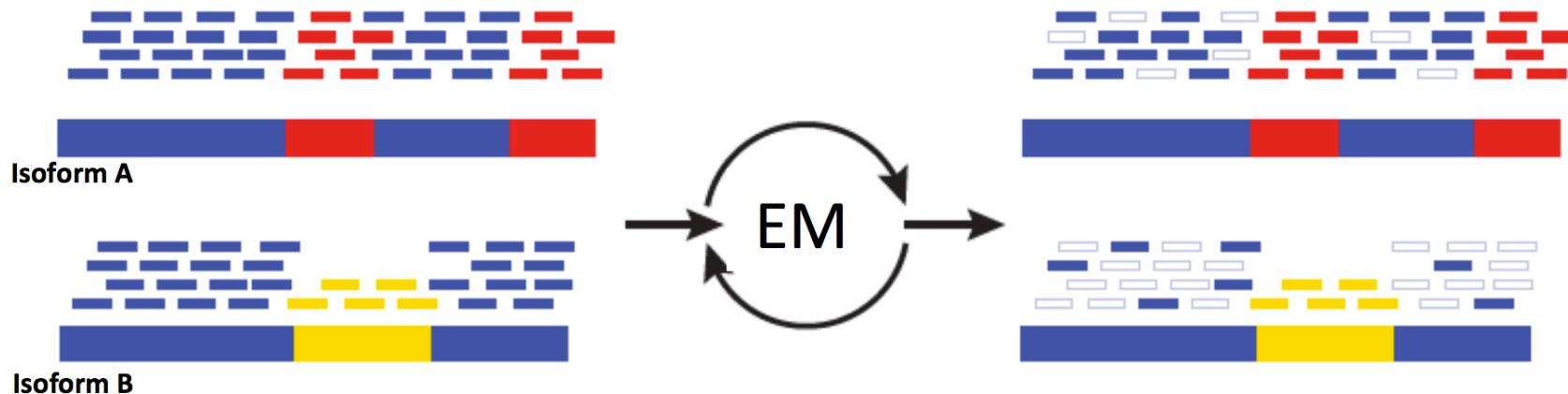
FPKM

Multiply-mapped Reads Confound Abundance Estimation



Blue = multiply-mapped reads
Red, Yellow = uniquely-mapped reads

Multiply-mapped Reads Confound Abundance Estimation



Blue = multiply-mapped reads
Red, Yellow = uniquely-mapped reads

Use Expectation Maximization (EM) to find the most likely assignment of reads to transcripts.

Performed by:

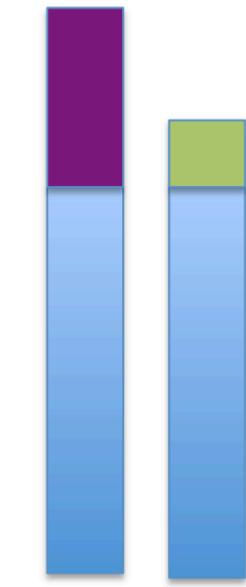
RSEM, eXpress, kallisto, salmon, ...

New fast alignment-free methods
now available! eg. Kallisto

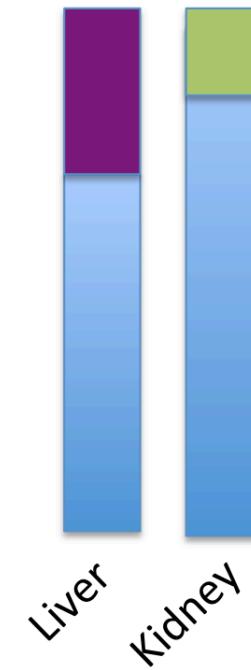
Differential Expression Analysis Using RNA-Seq

Why cross-sample normalization is important

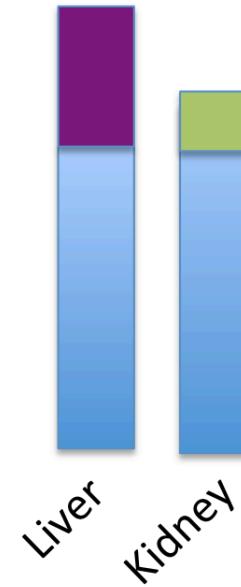
Absolute RNA
quantities per cell



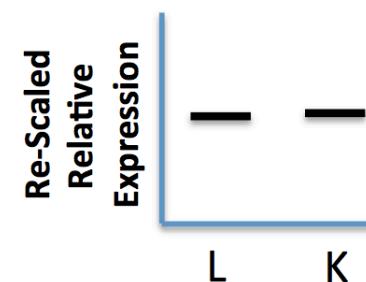
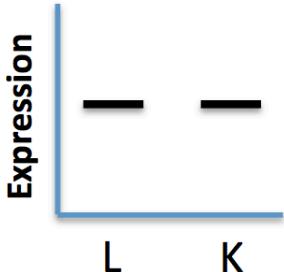
Measured relative
abundance via
RNA-Seq



Cross-sample
normalized
(rescaled) relative
abundance



e.g. Some housekeeping gene's expression level:



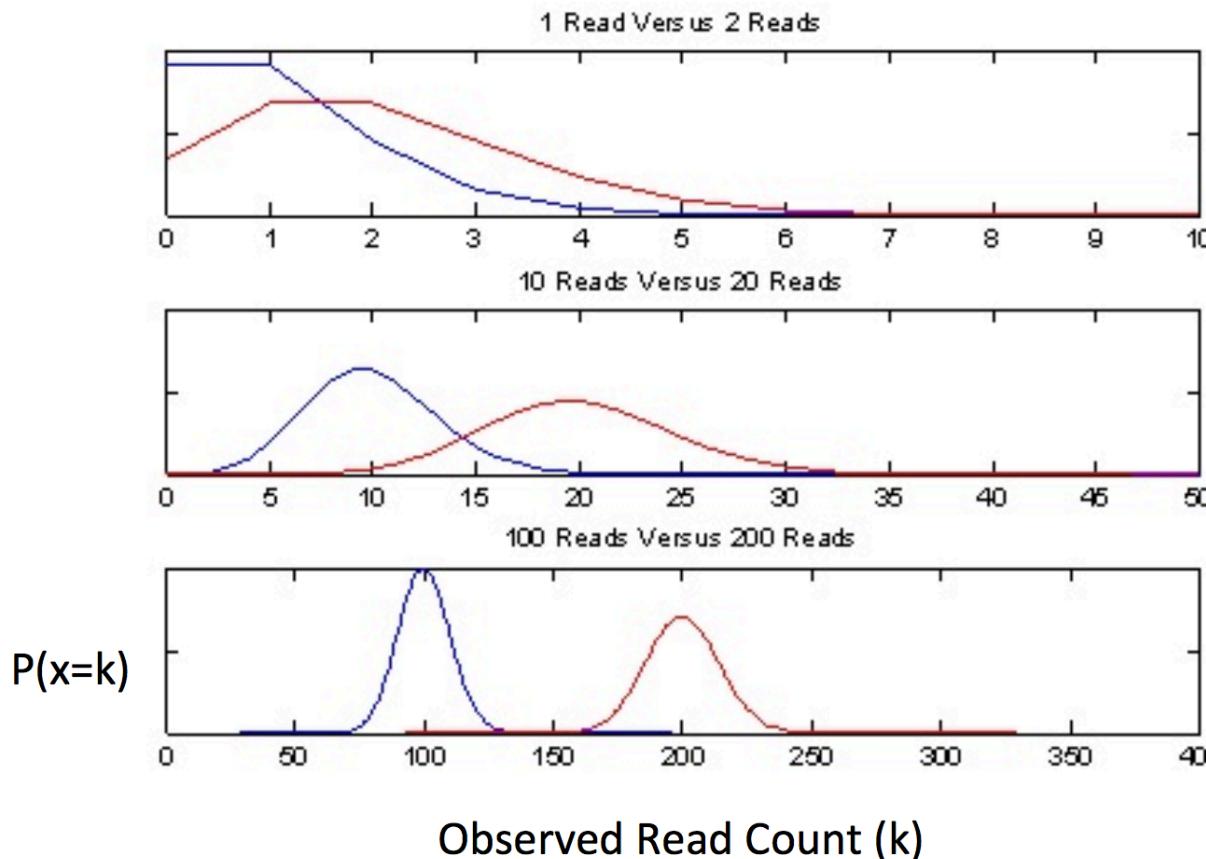
Diff. Expression Analysis Involves

- Counting reads
- Statistical significance testing

	Sample_A	Sample_B	Fold_Change	Significant?
Gene A	1	2	2-fold	No
Gene B	100	200	2-fold	Yes

Beware of concluding fold change from small numbers of counts

Poisson distributions for counts based on **2-fold** expression differences



No confidence in 2-fold difference. Likely observed by chance.

High confidence in 2-fold difference. Unlikely observed by chance.

Mapping can tell us something about
structural variation...

SNP Calling

GGCGTCTATATCTCGGCTCTAGGCCCTCATTTTT
GGCGTCTATATCT
GGCGTCTATATCTCG
TATCTCGGCTCTAGG
TATCTCAGCTCTAGGCC
TATCTCAGCTCTAGGCCCTCA
CTCGGCTCTAGGCCCTCATTTT
GGCTCTAGGCCCTCATTTTTT
CTCTAGGCCCTCATTTTTT
CTAGGCCCTCATTTTTT

Coverage = 5

SNP Calling

GGCGTCTATATCTCGGCTCTAGGCCCTCATT
GGCGTCTATATCT
GGCGTCTATATCTCG
TATCTCGGCTCTAGG
TATCTCAGCTCTAGGCC
TATCTCAGCTCTAGGCCCTCA
CTCGGCTCTAGGCCCTCATT
GGCTCTAGGCCCTCATT
CTCTAGGCCCTCATT
CTAGGCCCTCATT



SNP Calling

GGCGTCTATATCTGGCTCTAGGCCCTCATTTTT
GGCGTCTATATCT
GGCGTCTATATCTCG
TATCTCGGCTCTTGG
TATCTCAGCTCTGGCC
TATCTCAGCTCTGGCCCTCA
CTCGGCTCTGGCCCTCATTTT
GGCTCTGGCCCTCATTTTTT
CTCTGGCCCTCATTTTTT
CTTGCCCTCATTTTTT



SNP Calling

Basic Procedure:

- Map read to reference
- Read out differences

Considerations:

- Reads are short but genomes/transcriptomes long and complex
- Is the mapping position unique? Think about abundance est...
- Reads contain errors and errors are not random
- Base quality score? Error model?