



# RNA-Seq Analysis Workshop

## Tuxedo and Trinity for Next-Gen Transcriptome Studies

CSHL 2012-10

Brian Haas

Broad Institute

# Next-gen Sequencing Transforming Modern Science

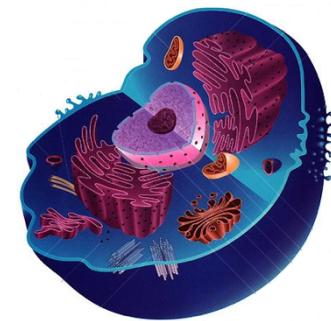
## Molecular Biology of the Cell

Chromatin structure

Histone occupancy

Transcription factor binding

DNA 3D topology



Genes and transcripts

gene content

alternative splicing

expression

RNA-editing

## Evolution



## Population Genetics

## Sequencing Methods

DNA-Seq

ChIP-Seq

RNA-Seq

Methyl-Seq

## Algorithms and Software Tools

Sequence assembly

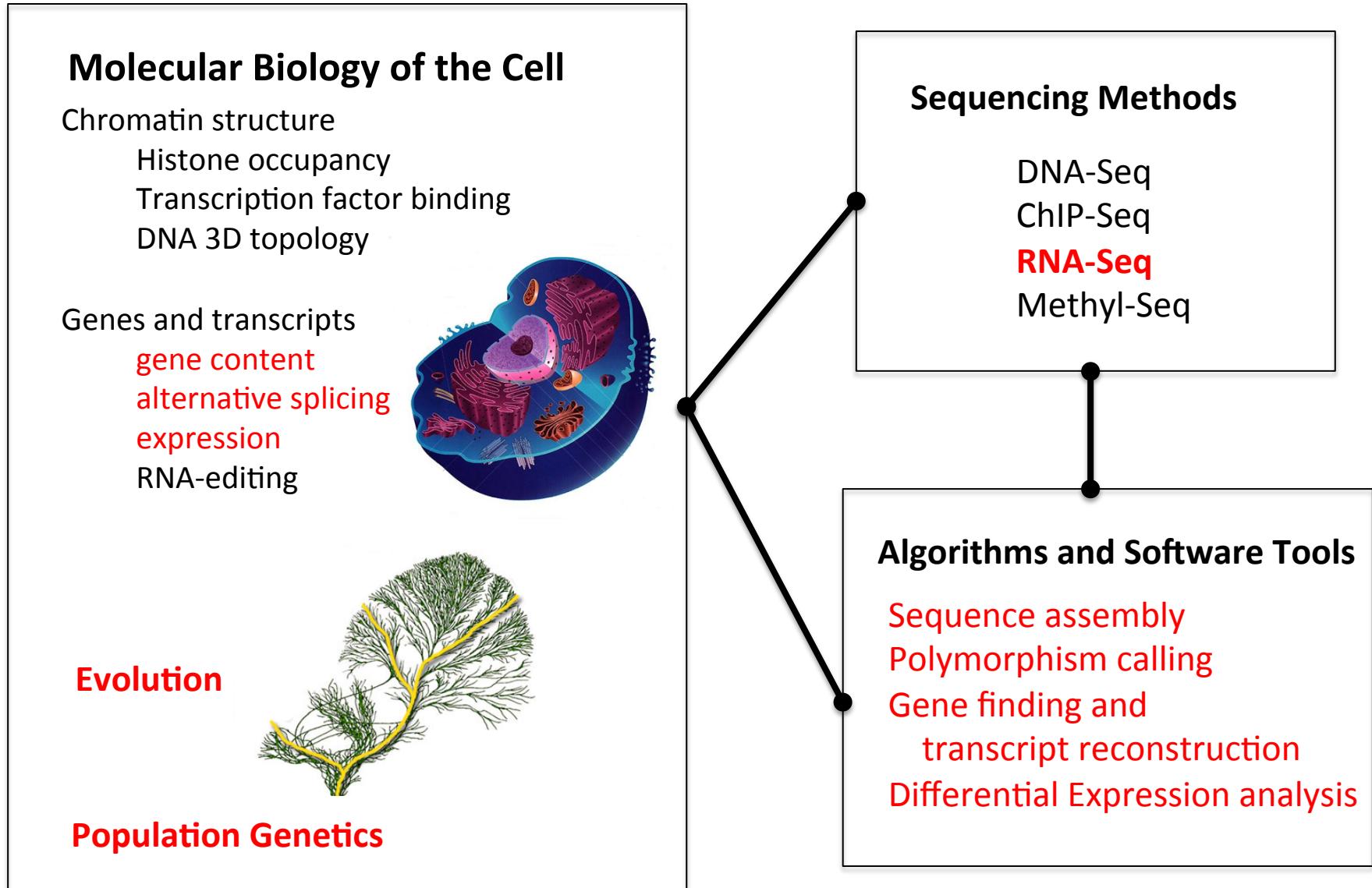
Polymorphism calling

Gene finding and

transcript reconstruction

Differential Expression analysis

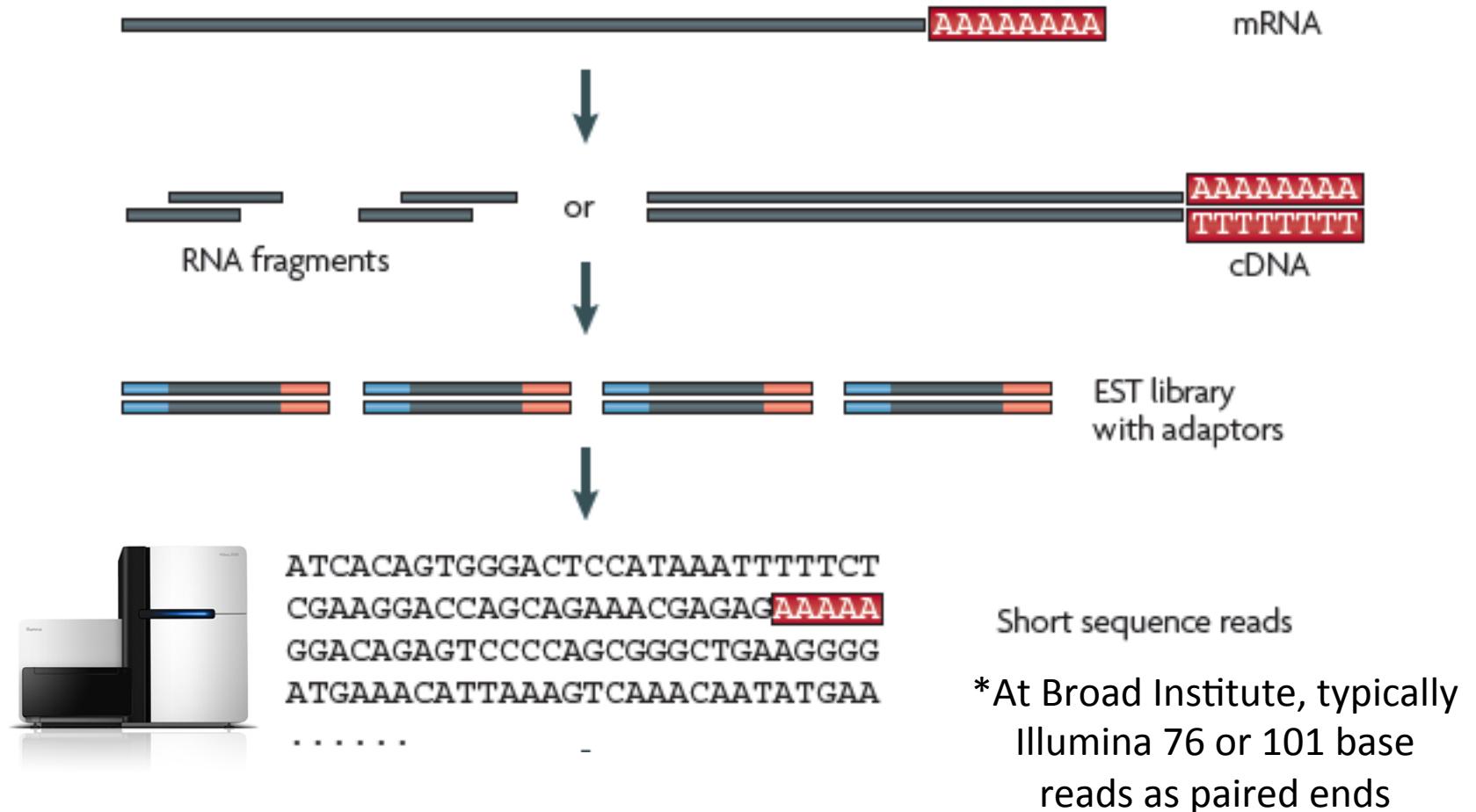
# RNA-Seq as a Driving Technology



# Outline of what's to follow:

- RNA-Seq basics
- Analysis paradigms
- Genome-based rna-seq studies
- Data formats and visualization
- De novo transcriptome-based rna-seq studies
- Transcript Quantification
- Differential Expression Analysis

# RNA-Sequencing Methodology



\*Adapted from Wang, Gerstein, and Snyder, Nature Reviews Genetics, 2009

# Common Data Formats for RNA-Seq

FastA format:

```
>61DFRAAXX100204:1:100:10494:3070/1  
AAACAAACAGGGCACATTGTCACTCTTGTATTGAAAAAACACTTCCGGCCAT
```

FastQ format:

```
@61DFRAAXX100204:1:100:10494:3070/1  
AAACAAACAGGGCACATTGTCACTCTTGTATTGAAAAAACACTTCCGGCCAT  
+  
ACCCCCCCCCCCCCCCCCCCCCCCCCCCCCCBC?CCCCCCCC@CACCCCCA
```

Read

Quality values

$$\text{AsciiEncodedQual} = -10 * \log_{10}(\text{Pwrong}) + 30$$

↑  
 $\text{Ascii } ('C') = 64$

$$\text{So, } \text{Pwrong}('C') = 10^{(64-30)/(-10)} = 10^{-3.4} = 0.0004$$

# Paired-end Sequences

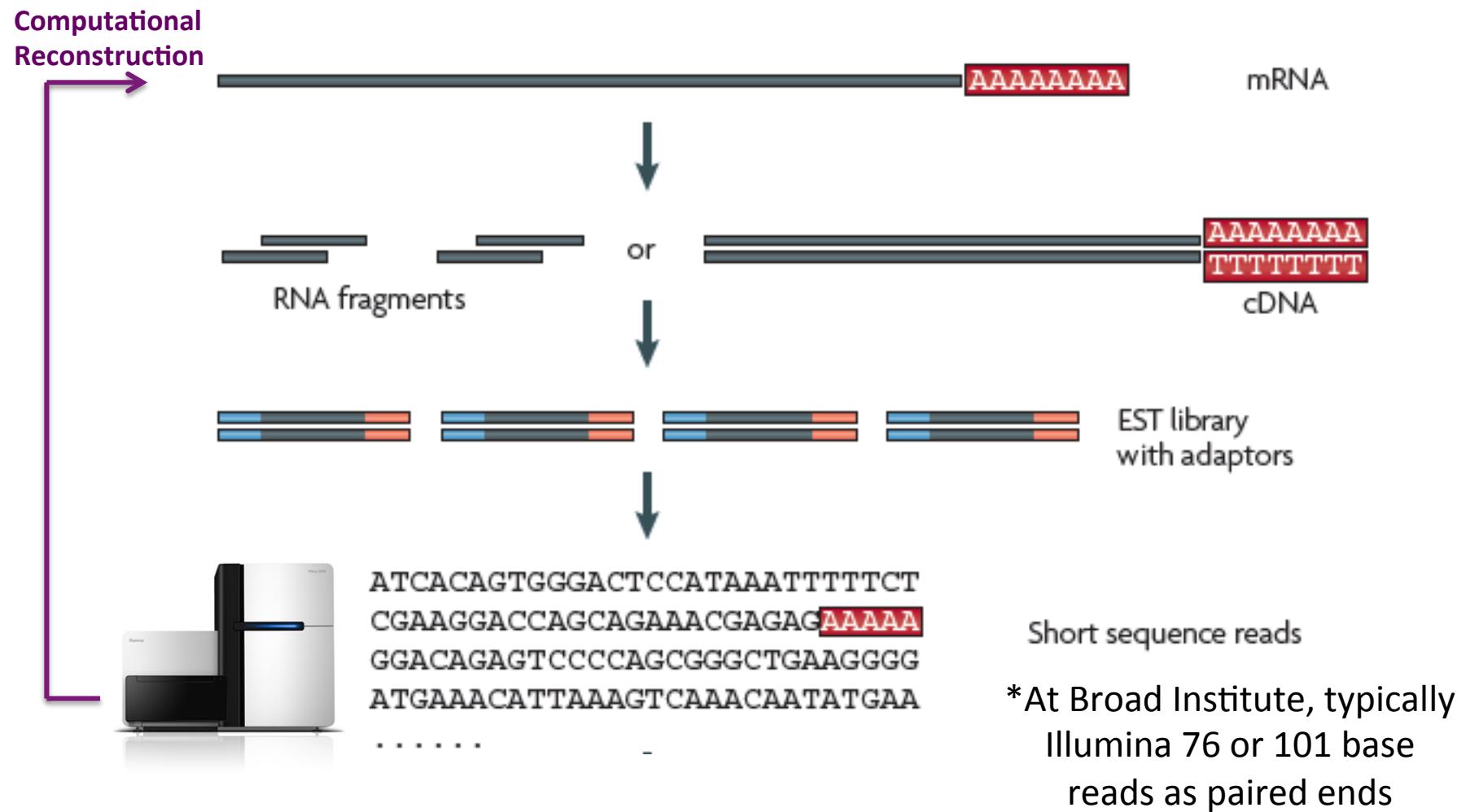


Two FastQ files, read name indicates  
left (/1) or right (/2) read of paired-end

```
@61DFRAAXX100204:1:100:10494:3070/1
AAACAAACAGGGCACATTGTCACTCTTGTATTGAAAAACACTTCCGGCCAT
+
ACCCCCCCCCCCCCCCCCCCCCCCCCCCCCBC?CCCCCCCC@CACCCCA
```

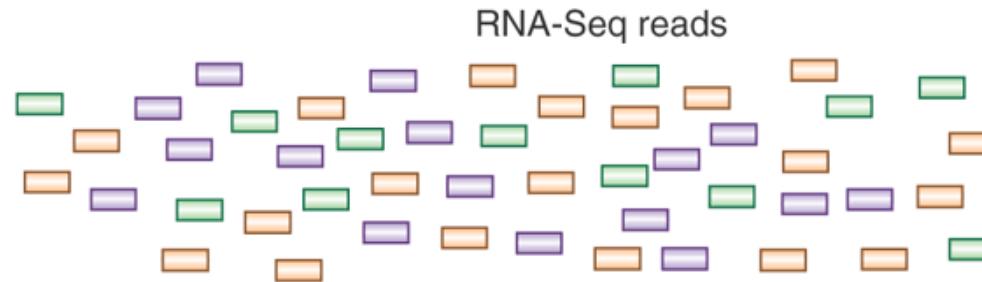
```
@61DFRAAXX100204:1:100:10494:3070/2
CTCAAATGGTTAATTCTCAGGCTGCAAATATTGTTAGGATGGAAGAAC
+
C<CCCCCCCCACCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCBCCCC
```

# RNA-Sequencing Methodology



\*Adapted from Wang, Gerstein, and Snyder, Nature Reviews Genetics, 2009

# Transcript Reconstruction from RNA-Seq Reads



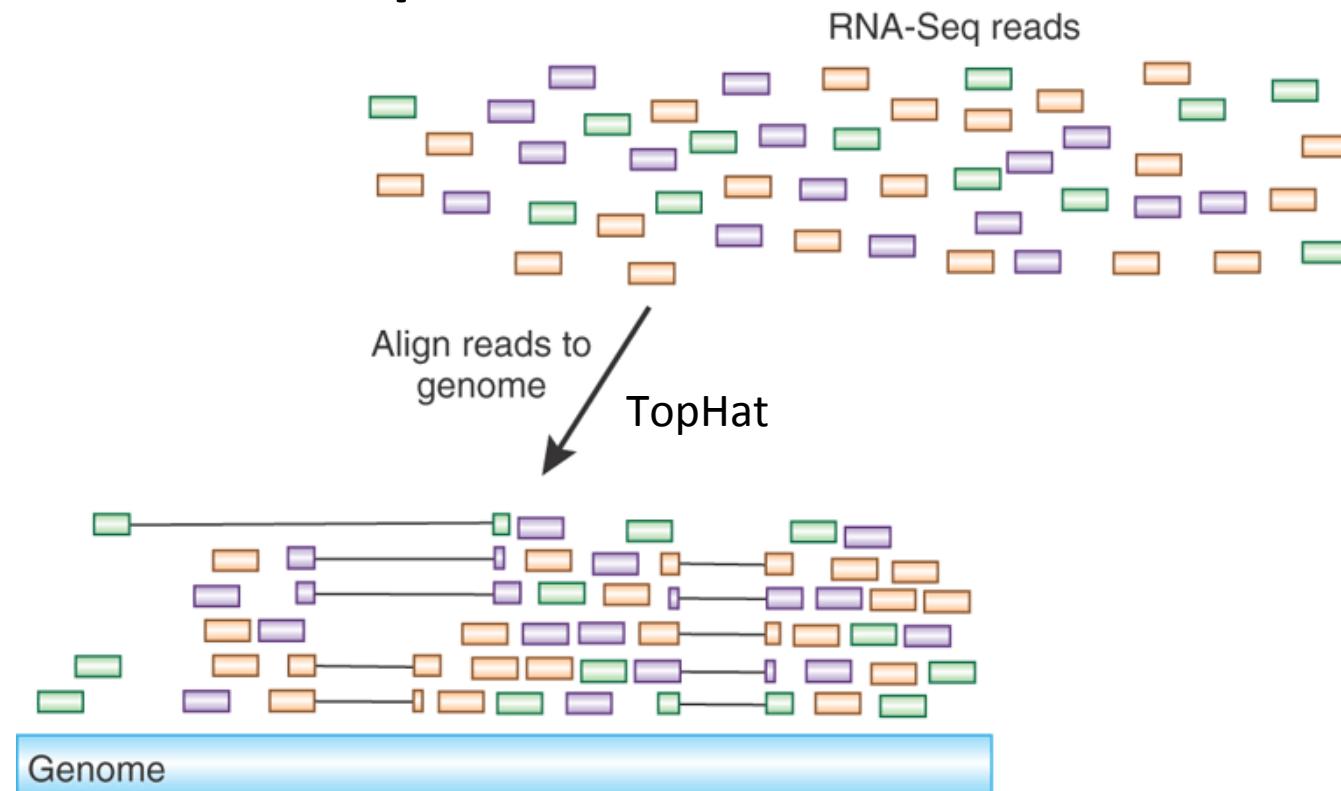
## Advancing RNA-Seq analysis

Brian J Haas & Michael C Zody

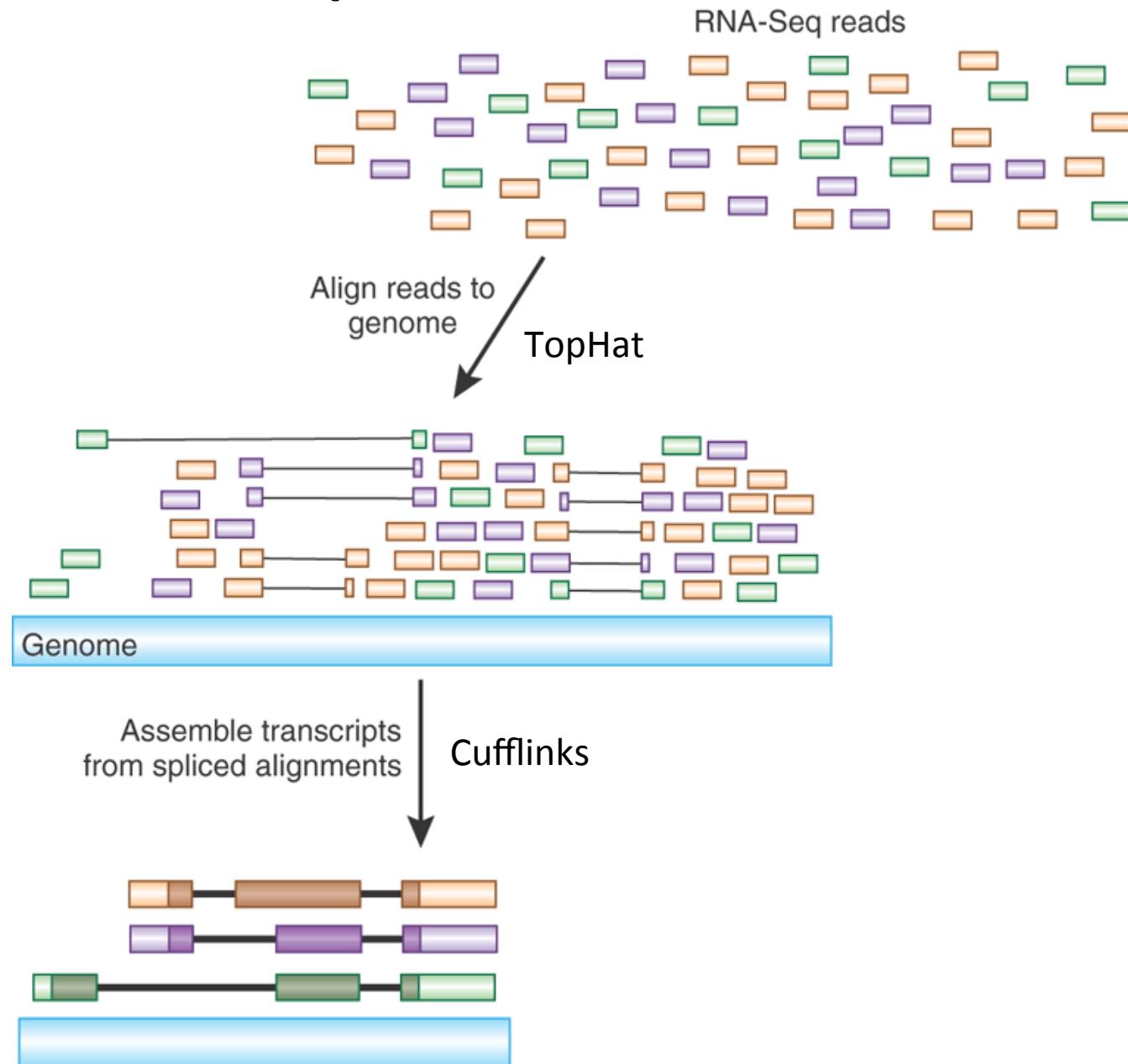
Nature Biotech, 2010

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

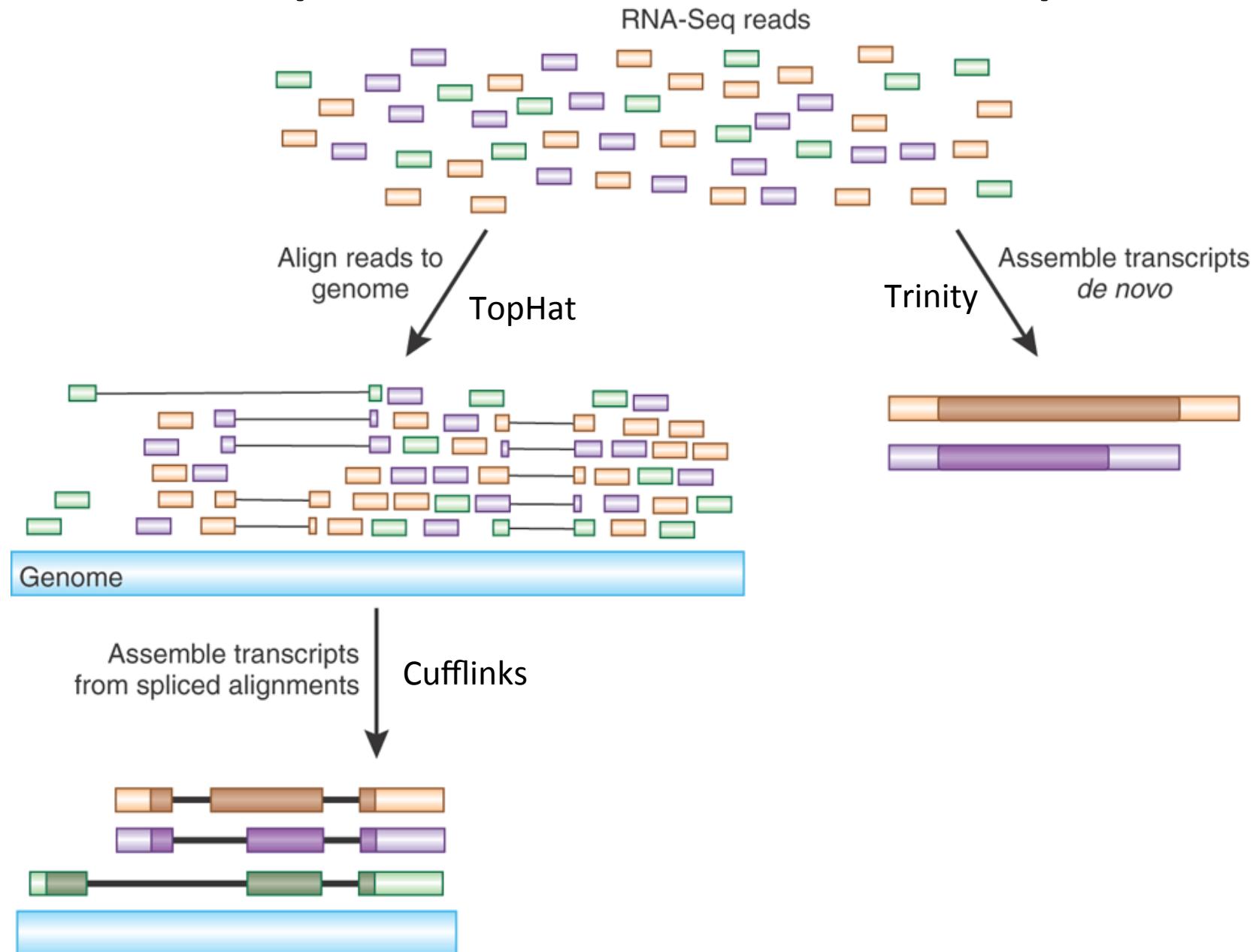
# Transcript Reconstruction from RNA-Seq Reads



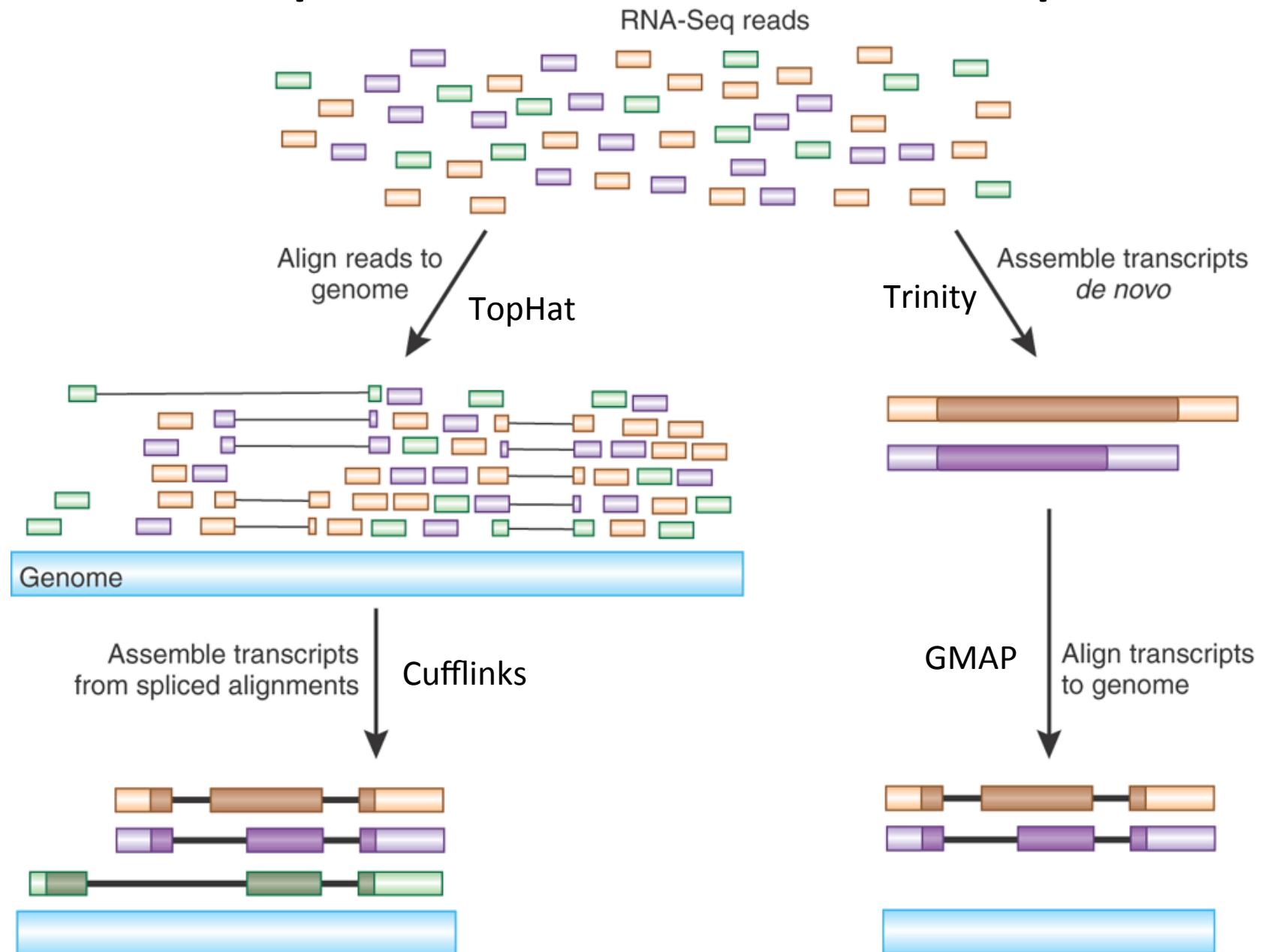
# Transcript Reconstruction from RNA-Seq Reads



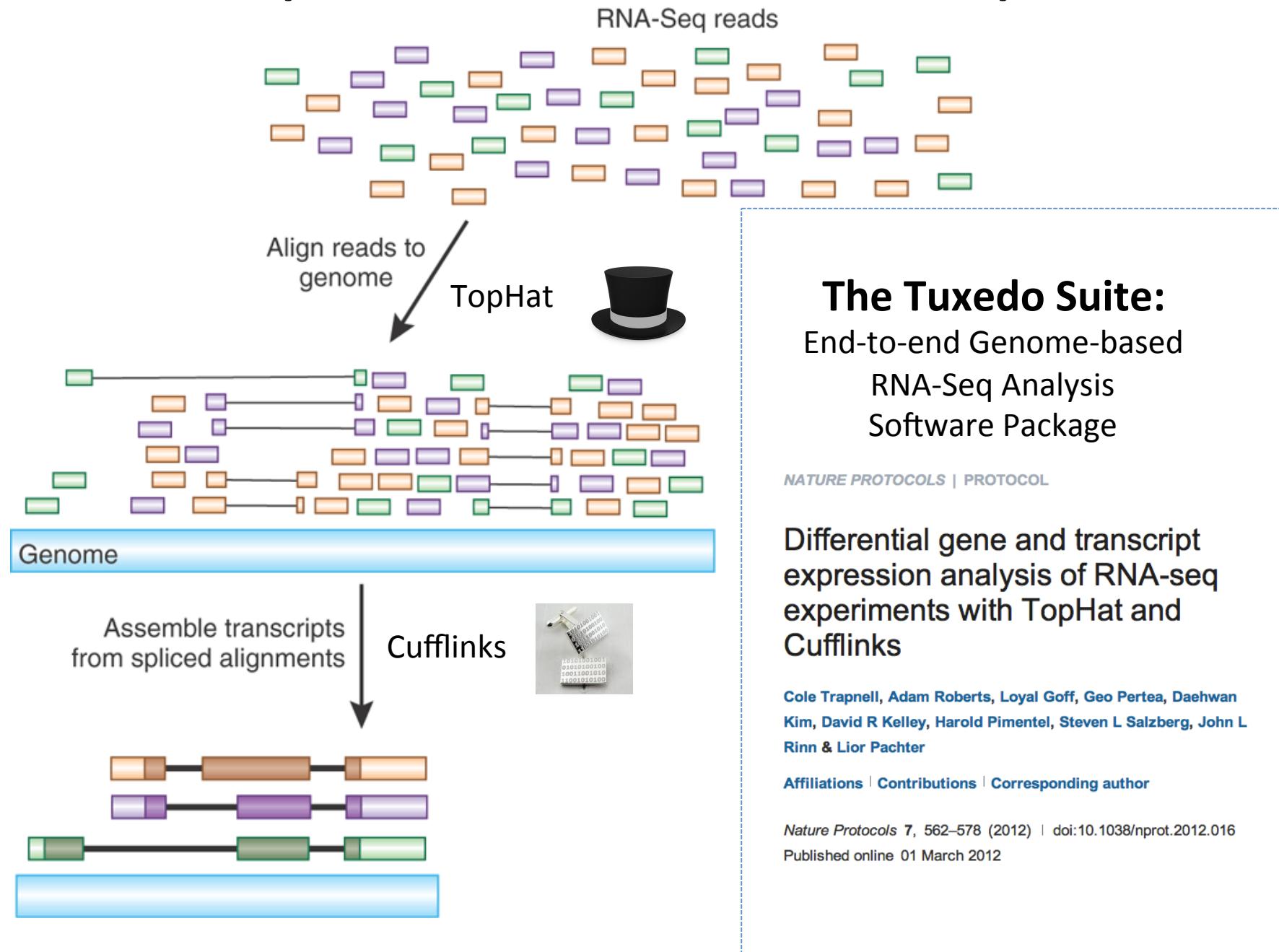
# Transcript Reconstruction from RNA-Seq Reads



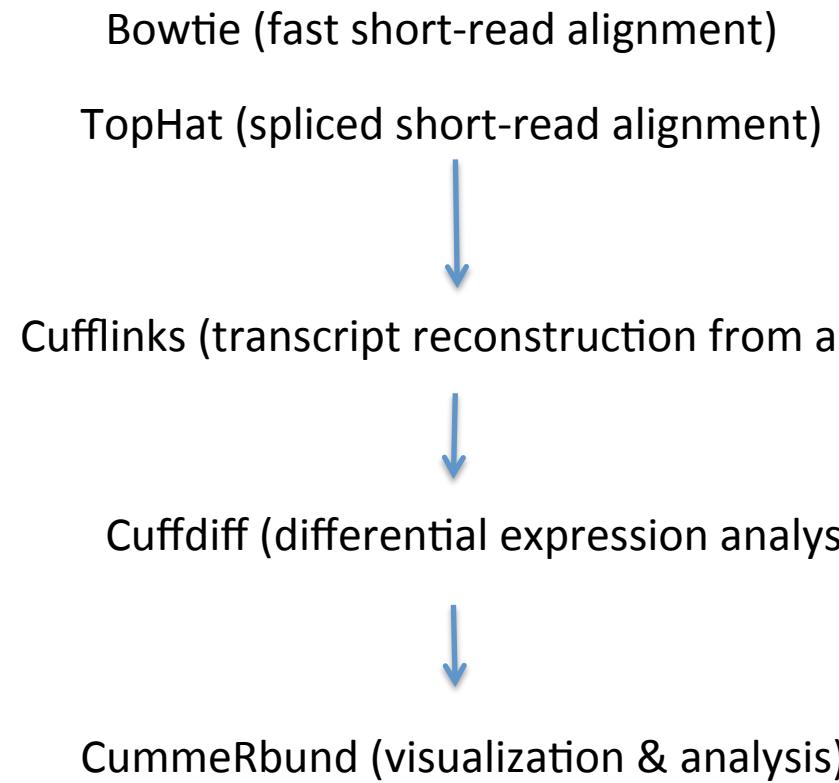
# Transcript Reconstruction from RNA-Seq Reads



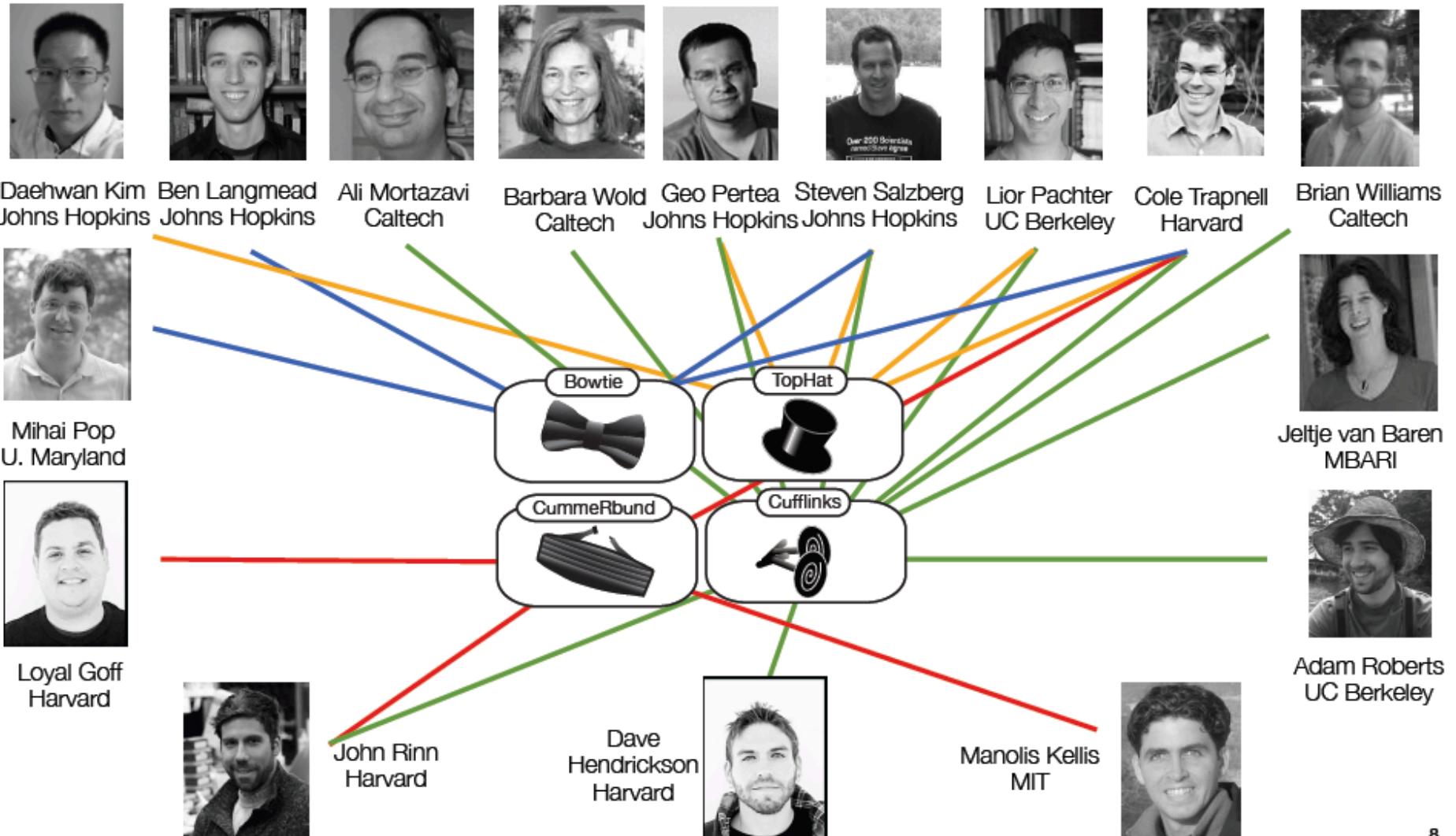
# Transcript Reconstruction from RNA-Seq Reads



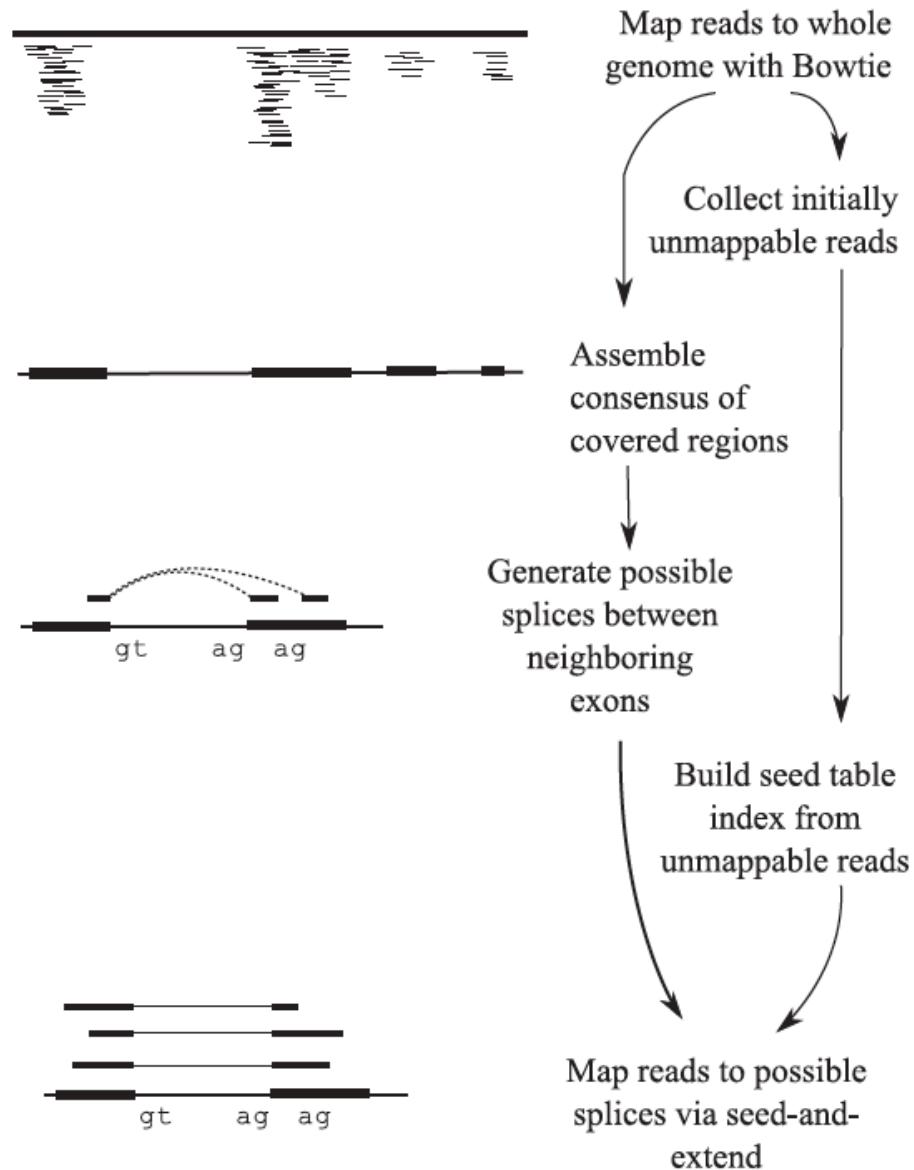
# Overview of the Tuxedo Software Suite



# Tuxedo development team



# The TopHat Pipeline



From Trapnell, Pachter, & Salzberg. Bioinformatics. 2009

## Alignments are reported in a compact representation: SAM format

```
0      61G9EAAXX100520:5:100:10095:16477
1      83
2      chr1
3      51986
4      38
5      46M
6      =
7      51789
8      -264
9      CCCAAACAAGCCGAACTAGCTGATTTGGCTCGTAAAGACCCGGAAA
10     ###CB?=ADDBCBCDEFFDEFFFDEFFGDBEFGEDGCFGFGGGGG
11     MD:Z:67
12     NH:i:1
13     HI:i:1
14     NM:i:0
15     SM:i:38
16     XQ:i:40
17     X2:i:0
```

[Link to SAM format description](#)

## Alignments are reported in a compact representation: SAM format

```
0      61G9EAAXX100520:5:100:10095:16477 (read name)
1      83 (FLAGS stored as bit fields; 83 = 00001010011 )
2      chr1 (alignment target)
3      51986 (position alignment starts)
4      38
5      46M (Compact description of the alignment in CIGAR format)
6      =
7      51789
8      -264 → (read sequence, oriented according to the forward alignment)
9      CCCAAACAAGCCGAACTAGCTGATTGGCTCGTAAAGACCCGGAAA
10     ###CB?=ADDBCBCDEEFFDEFDEFFGDBEFGEDGCFGFGGGGG
11     MD:Z:67                                     → (base quality values)
12     NH:i:1
13     HI:i:1
14     NM:i:0
15     SM:i:38          (Metadata)
16     XQ:i:40
17     X2:i:0
```

[Link to SAM format description](#)

Alignments are reported in a compact representation: SAM format

```
0      61G9EAAXX100520:5:100:10095:16477 (read name)
1      83 (FLAGS stored as bit fields; 83 = 00001010011 )
2      chr1 (alignment target)
```

Still not compact enough...

Millions to billions of reads takes up a lot of space!!

Convert SAM to binary – BAM format.

```
15     SM:i:38      (read data)
16     XQ:i:40
17     X2:i:0
```

[Link to SAM format description](#)

# Samtools

- Tools for
  - converting SAM <-> BAM
  - Viewing BAM files (eg. samtools view file.bam | less )
  - Sorting BAM files, and lots more:

```
Program: samtools (Tools for alignments in the SAM format)
Version: 0.1.18 (r982:295)

Usage:   samtools <command> [options]

Command: view      SAM<->BAM conversion
          sort      sort alignment file
          mpileup   multi-way pileup
          depth     compute the depth
          faidx    index/extract FASTA
          tview     text alignment viewer
          index    index alignment
          idxstats  BAM index stats (r595 or later)
          fixmate   fix mate information
          flagstat  simple stats
          calmd     recalculate MD/NM tags and '=' bases
          merge     merge sorted alignments
          rmdup    remove PCR duplicates
          reheader  replace BAM header
          cat       concatenate BAMs
          targetcut cut fosmid regions (for fosmid pool only)
          phase    phase heterozygotes
```

# Visualizing Alignments of RNA-Seq reads

# Text-based Alignment Viewer

```
% samtools tview alignments.bam target.fasta
```

# IGV

www.broadinstitute.org/igv/

The screenshot shows the IGV website at [www.broadinstitute.org/igv/](http://www.broadinstitute.org/igv/). The page features a large central image of the IGV software interface, which displays a genomic track with multiple panels showing sequencing data, tracks, and annotations. To the left of the main content area is a sidebar with the IGV logo and a navigation menu. The menu includes links for Home, Downloads, Documents, Hosted Genomes, FAQ, IGV User Guide, File Formats, Release Notes, Credits, Contact, and a search bar. Below the menu is a section for the Broad Institute, including links to Broad Home and Cancer Program, and the text "© 2012 Broad Institute".

**Home**

# Integrative Genomics Viewer

**What's New**

**NEWS**  
July 3, 2012. Soybean (*Glycine max*) and Rat (*rn5*) genomes have been updated.

**April 20, 2012.** IGV 2.1 has been released.  
See the [release notes](#) for more details.

**April 19, 2012.** See our new [IGV paper](#) in *Briefings in Bioinformatics*.

**Overview**

**Citing IGV**

To cite your use of IGV in your publication:

James T. Robinson, Helga Thorvaldsdóttir, Wendy Winckler, Mitchell Guttman, Eric S. Lander, Gad Getz, Jill P. Mesirov. [Integrative Genomics Viewer](#). *Nature Biotechnology* 29, 24–26 (2011), or  
Helga Thorvaldsdóttir, James T. Robinson, Jill P. Mesirov. [Integrative Genomics Viewer \(IGV\): high-performance genomics data visualization and exploration](#).

# GenomeView

← → C genomeview.org ⭐ a

# GenomeView

Demos Plug-ins JAnnot API Join mailing list Support - Frequently asked questions Cite us

## Start Now!

Webstart:  
[Launch](#)

High-mem webstart

Applet:  
[Launch](#)

## Documentation

- Quick start guide
- Manual
- Advanced manual
- Tutorials

## Navigation

- Download
- Demos
- Plug-ins

GenomeView is a next-generation stand-alone genome browser and editor initiated in the BSB group at VIB and currently developed at Broad Institute. It provides interactive visualization of sequences, annotation, multiple alignments, syntetic mappings, short read alignments and more. Many standard file formats are supported and new functionality can be added using a plugin system.

### Getting started



Get started with a five minute quick-start guide that will get up and running in no time

### Web start

[Launch](#) Click the launch button to start GenomeView

### Download



Download the current release. You can also start GenomeView

### Support

If you experience any issues, head over to the [support section](#), we like to help you.

### Recent questions

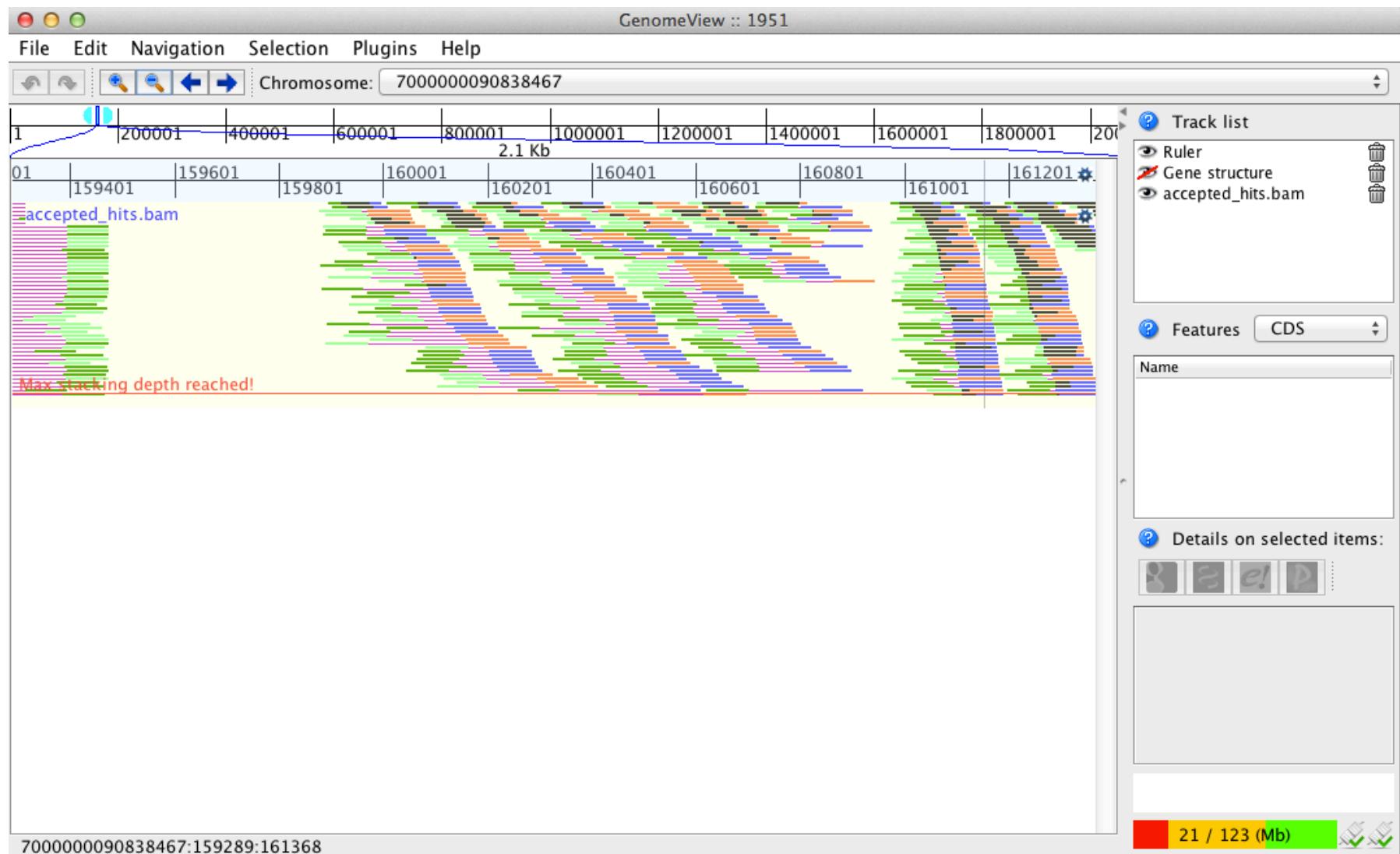
- How do I show annotation on a multiple alignment?
- Why does my multiple alignment load as reference sequences?
- Where do I find documentation?
- Why doesn't GenomeView correctly detect my BED file?
- How do I fix the order of the tracks in an integrated GenomeView instance?
- How do I integrate GenomeView in my

Most Creative Visualization  
IDEA Challenge 2011  
Academic

Most creative visualization award @ Illumina iDEA challenge 2011

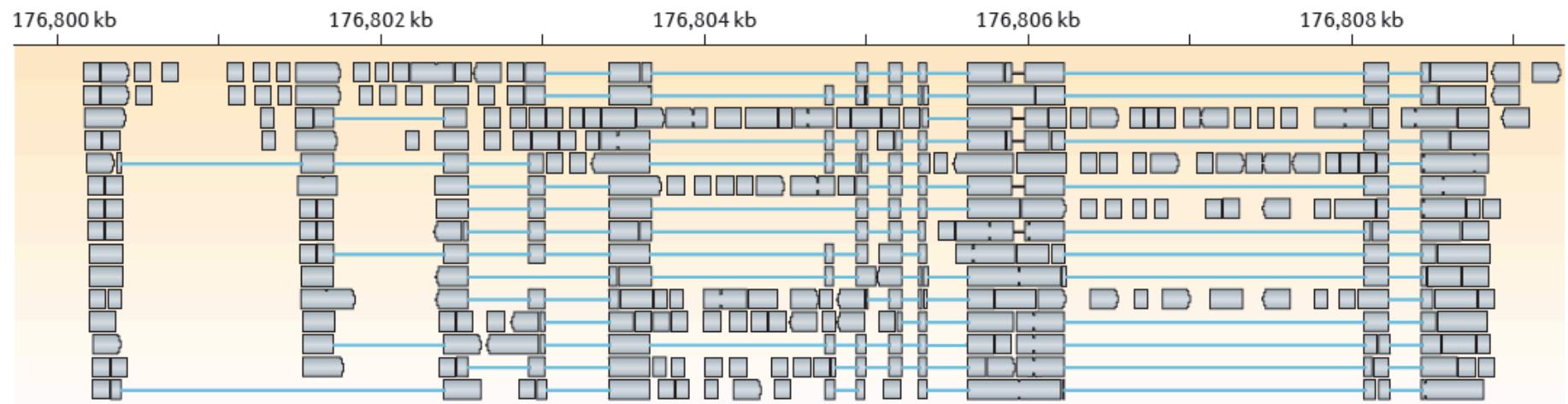
KILLER APP AWARD

# GenomeView: viewing TopHat alignments



# Transcript Reconstruction Using Cufflinks

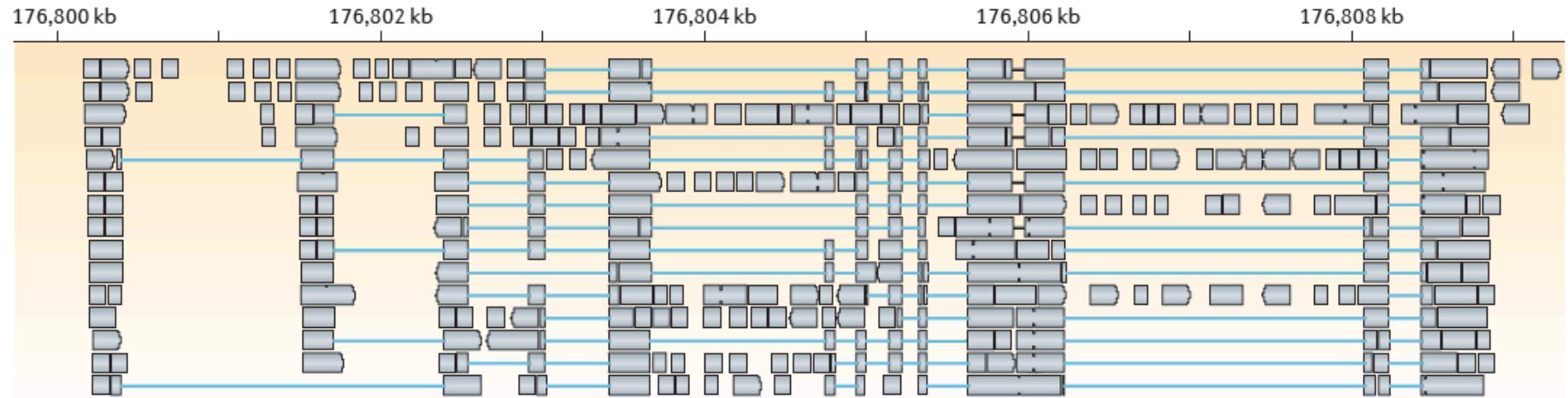
a Splice-align reads to the genome



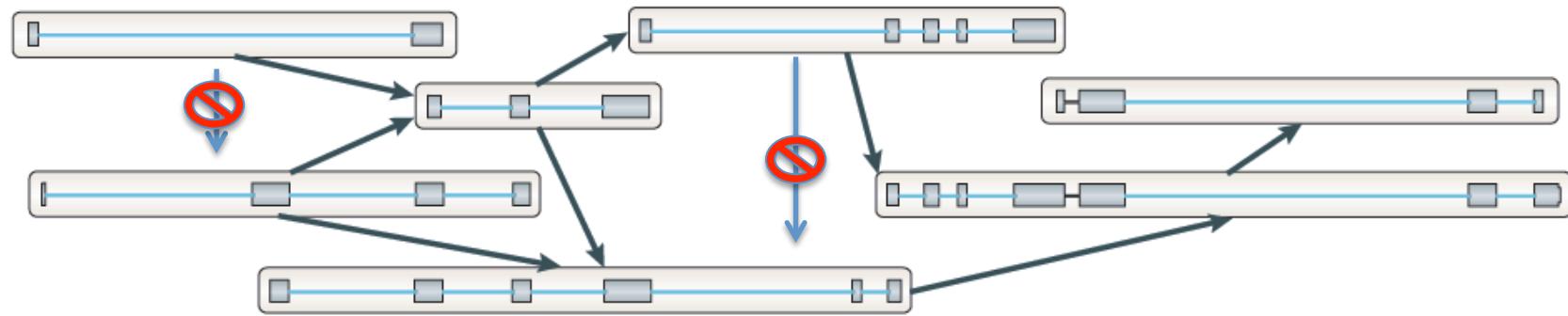
From Martin & Wang. Nature Reviews in Genetics. 2011

# Transcript Reconstruction Using Cufflinks

## a Splice-align reads to the genome



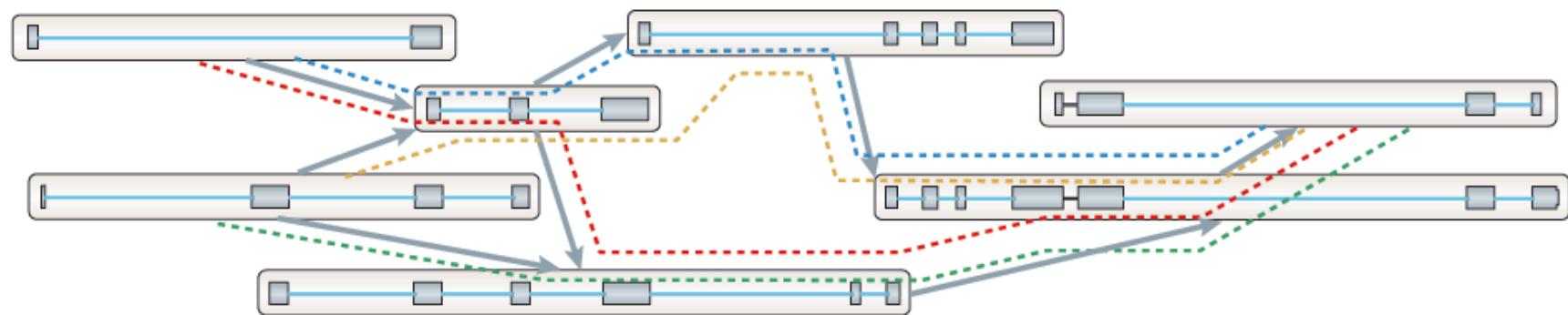
## b Build a graph representing alternative splicing events



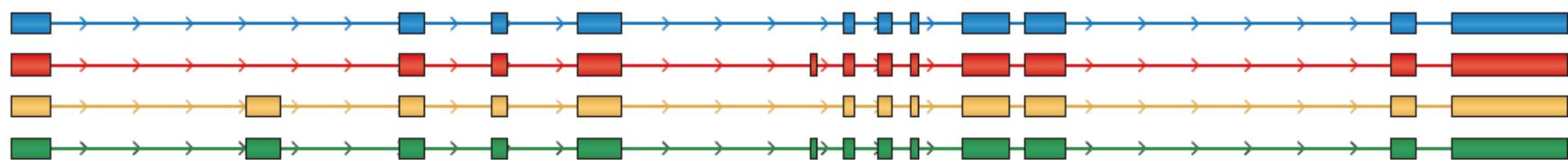
From Martin & Wang. Nature Reviews in Genetics. 2011

# Transcript Reconstruction Using Cufflinks

## C Traverse the graph to assemble variants



## d Assembled isoforms



# Transcript Structures in GTF Format

(tab-delimited fields per line shown transposed to a column format here)

```
0 7000000090838467  (genomic contig identifier)
1 Cufflinks
2 transcript
3 101  (left coordinate)
4 5716 (right coordinate)
5 1000
6 +      (strand)
7 .
8 gene_id "CUFF.1"; transcript_id "CUFF.1.1"; FPKM "378.0239937260"  (annotations)
```

```
0 7000000090838467
1 Cufflinks
2 exon
3 101
4 5716
5 1000
6 +
7 .
8 gene_id "CUFF.1"; transcript_id "CUFF.1.1"; exon_number "1"; FPKM "378.0239937260"
```

# De novo transcriptome assembly

No genome required

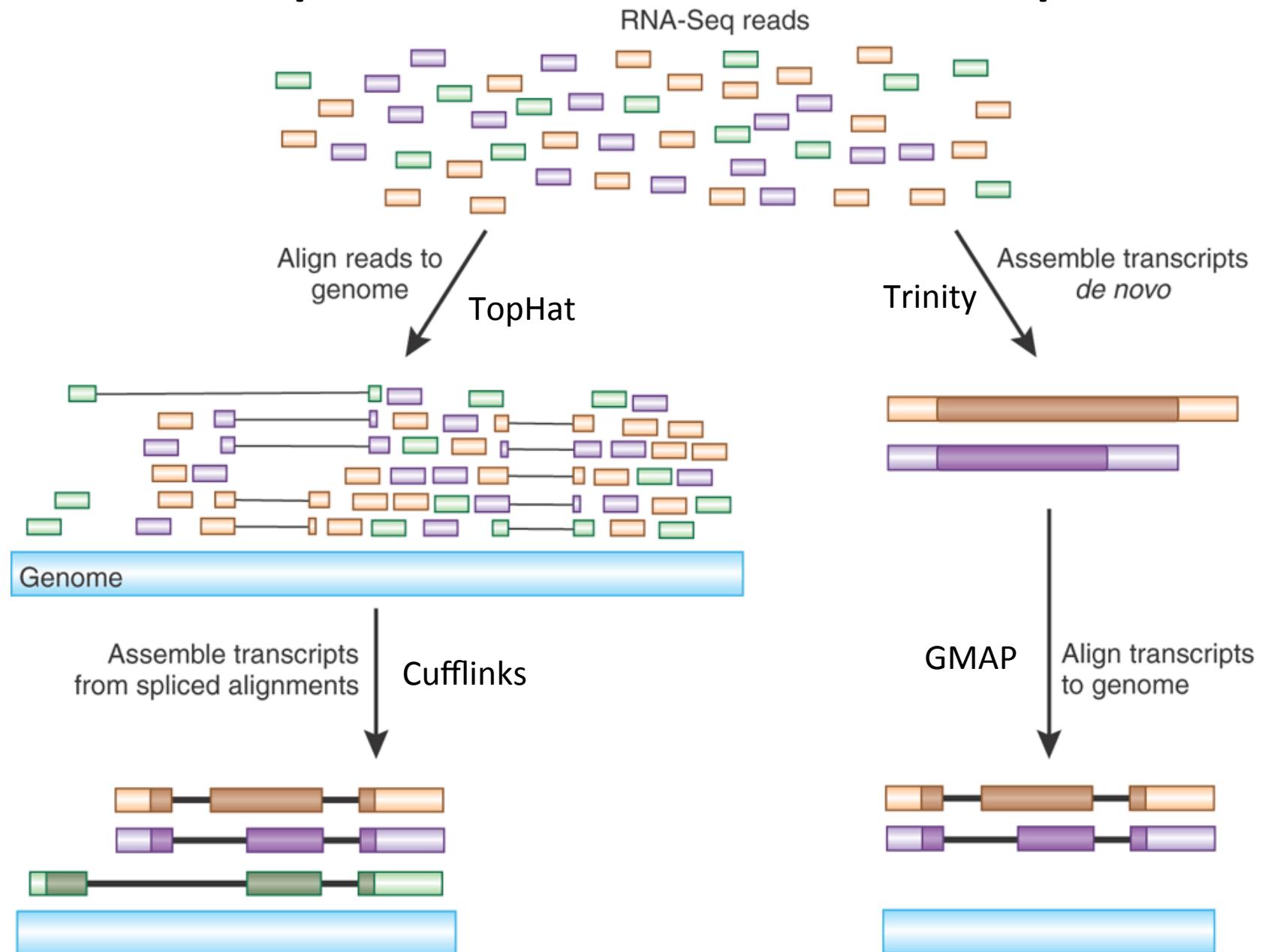
Empower studies of non-model organisms

expressed gene content

transcript abundance

differential expression

# Transcript Reconstruction from RNA-Seq Reads



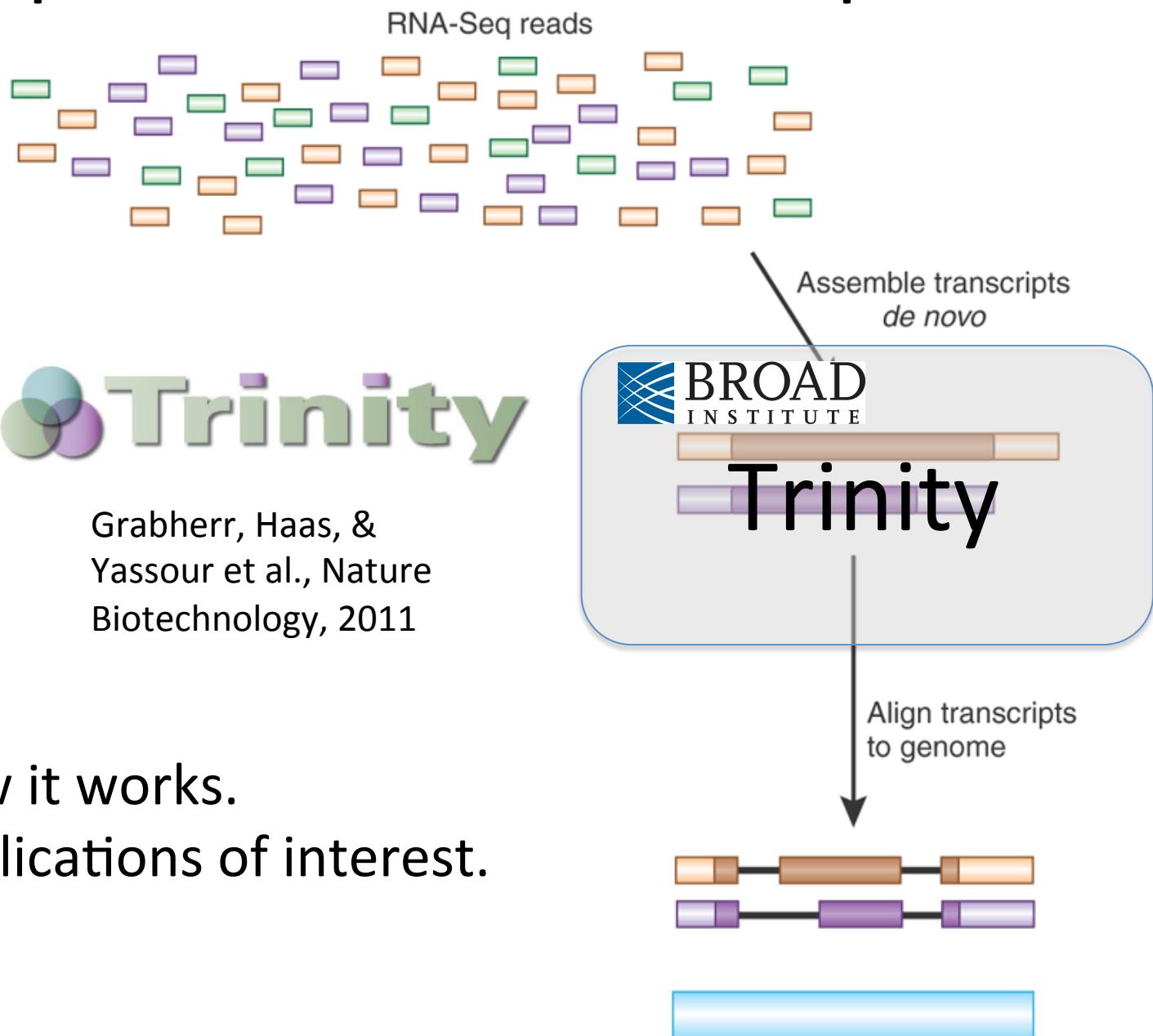
# Transcript Reconstruction from RNA-Seq Reads



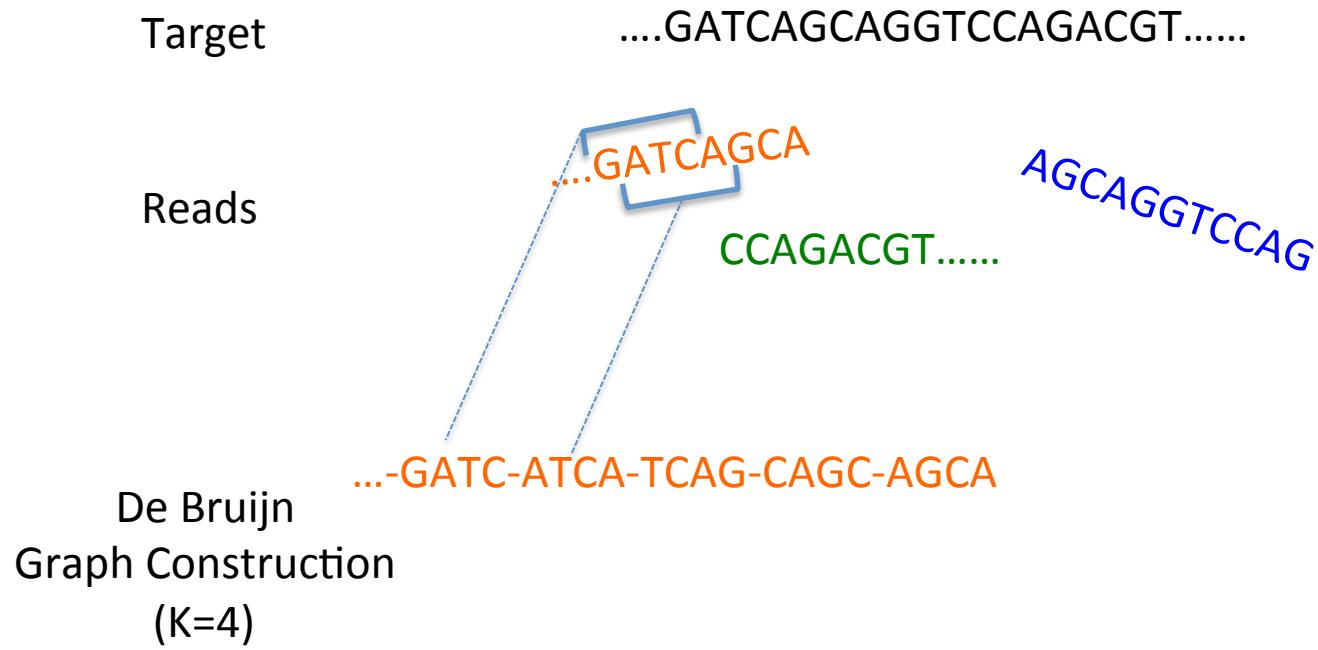
 **Trinity**

Grabherr, Haas, &  
Yassour et al., Nature  
Biotechnology, 2011

- How it works.
- Applications of interest.



# Short Read Assembly Using de Bruijn Graphs



# Short Read Assembly Using de Bruijn Graphs

Target

....GATCAGCAGGTCCAGACGT.....

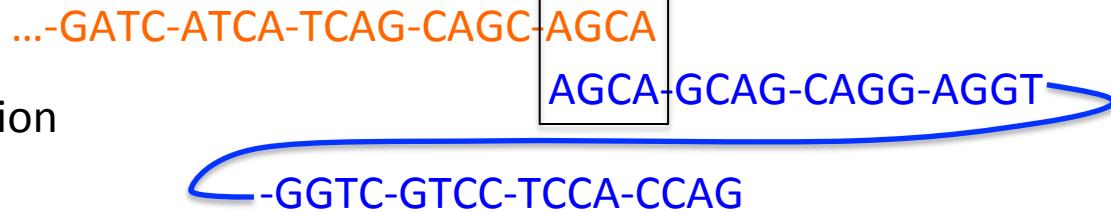
Reads

....GATCAGCA

CCAGACGT.....

AGCAGGTCCAG

De Bruijn  
Graph Construction  
(K=4)



# Short Read Assembly Using de Bruijn Graphs

Target

....GATCAGCAGGTCCAGACGT.....

Reads

....GATCAGCA

CCAGACGT.....

AGCAGGTCCAG

De Bruijn  
Graph Construction  
(K=4)

...-GATC-ATCA-TCAG-CAGC-

AGCA

AGCA-GCAG-CAGG-AGGT

-GGTC-GTCC-TCCA-CCAG

CCAG

CAGA-AGAC-GACG-ACGT-....

# Short Read Assembly Using de Bruijn Graphs

Target

....GATCAGCAGGTCCAGACGT.....

Reads

....GATCAGCA

CCAGACGT.....

AGCAGGTCCAG

De Bruijn  
Graph Construction  
(K=4)

...-GATC-ATCA-TCAG-CAGC-AGCA

AGCA-GCAG-CAGG-AGGT

-GGTC-GTCC-TCCA-CCAG

CCAG-CAGA-AGAC-GACG-ACGT-....

Sequence Reconstruction  
By Path Traversal

....GATCAGCAGGTCCAGACGT.....

# Contrasting Genome and Transcriptome Assembly

## Genome Assembly

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

## Transcriptome Assembly

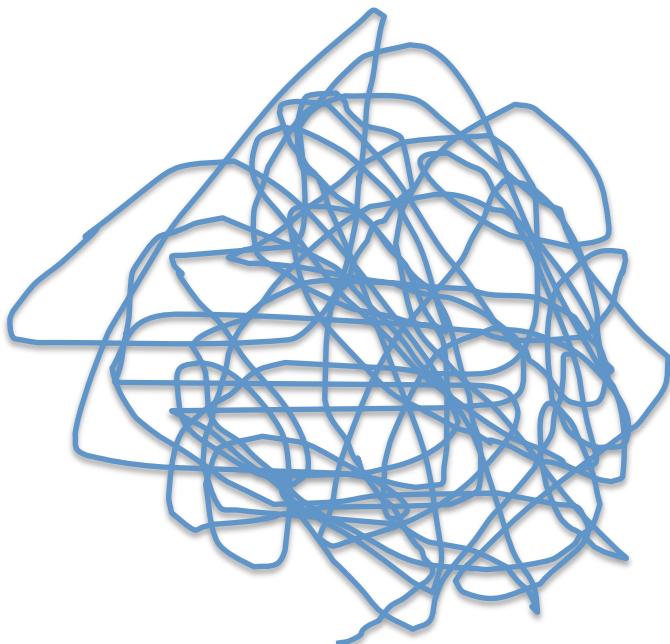
- Exponentially distributed coverage levels
- Multiple contigs per locus (alt splicing)
- Strand-specific



# Trinity Aggregates Isolated Transcript Graphs

## Genome Assembly

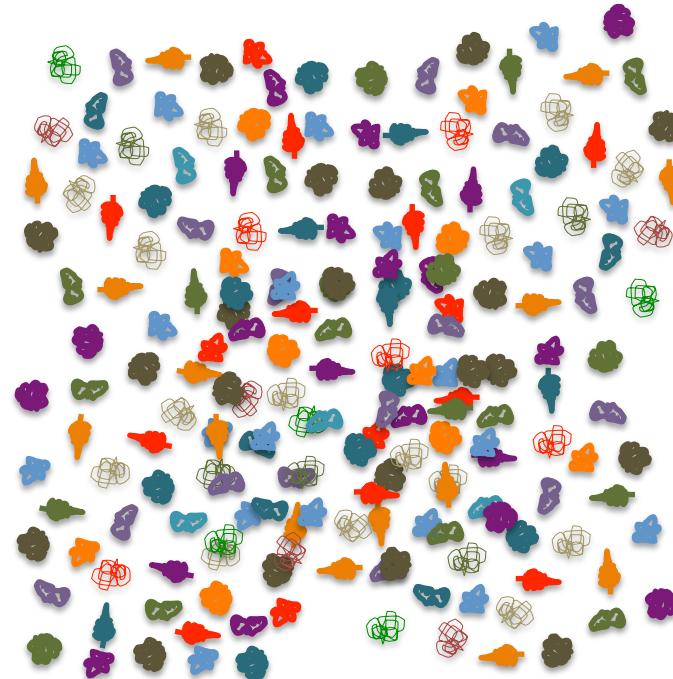
Single Massive Graph



Entire chromosomes represented.

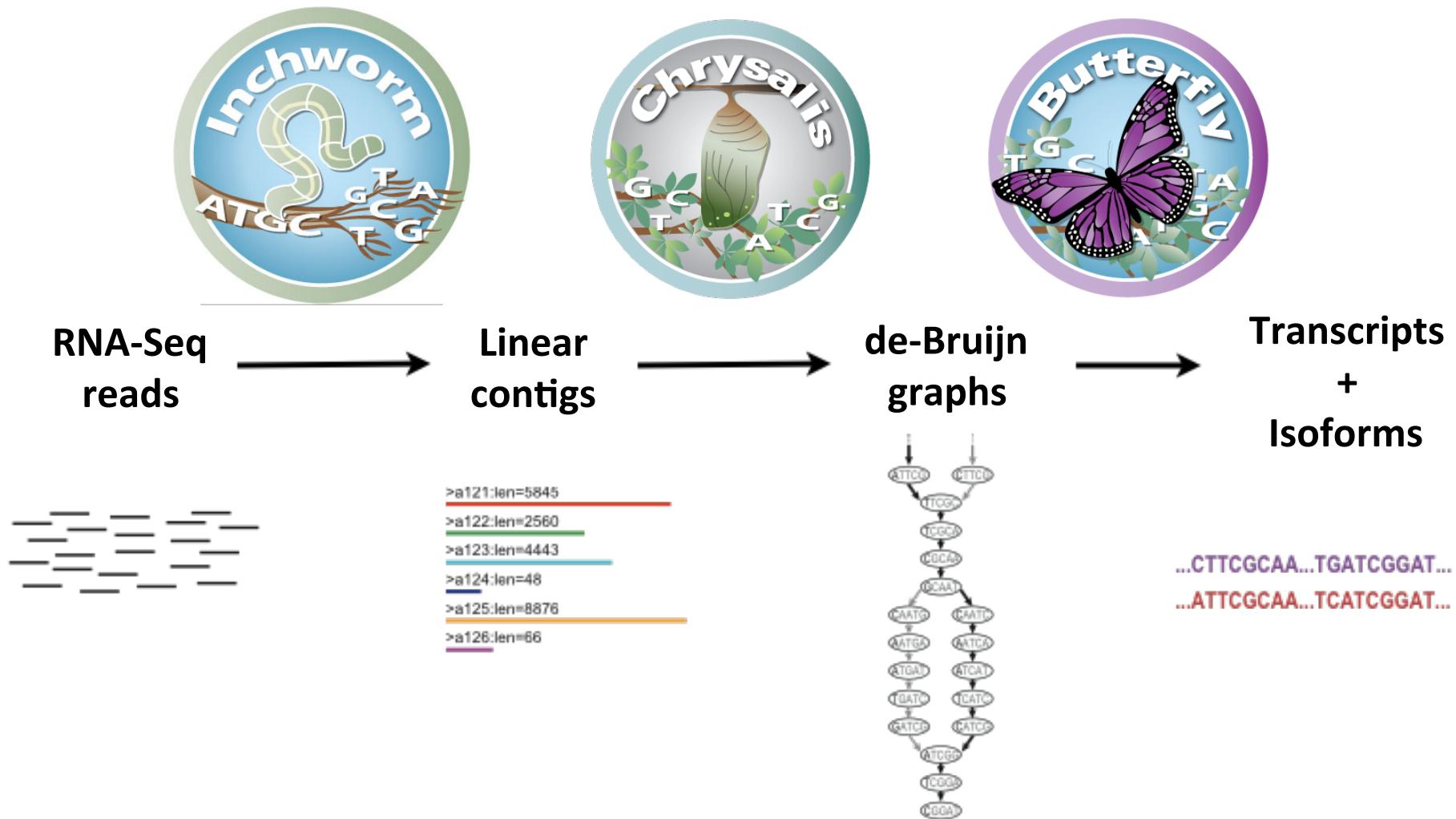
## Trinity Transcriptome Assembly

Many Thousands of Small Graphs



Ideally, one graph per expressed gene.

# Trinity



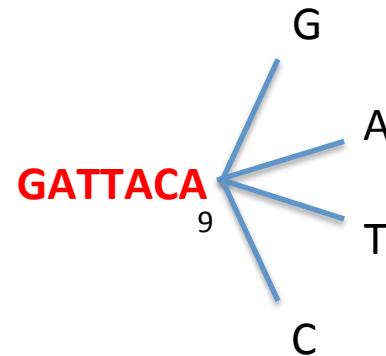


# Inchworm Algorithm

Decompose all reads into overlapping Kmers (25-mers)

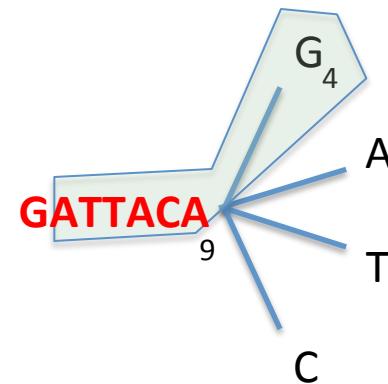
Identify seed kmer as most abundant Kmer, ignoring low-complexity kmers.

Extend kmer at 3' end, guided by coverage.



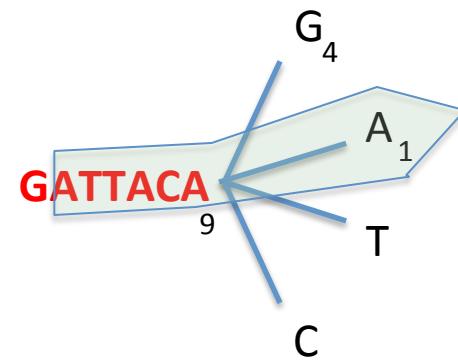


# Inchworm Algorithm



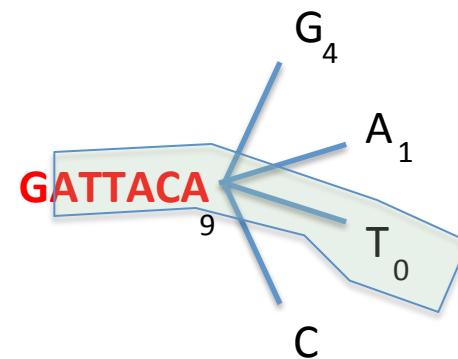


# Inchworm Algorithm



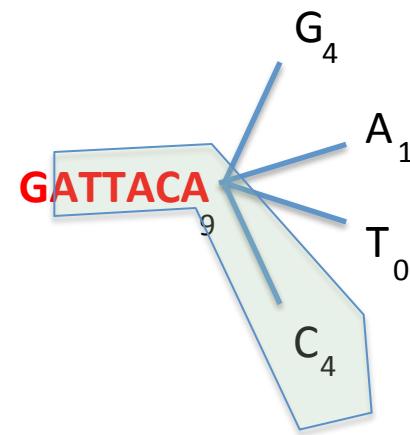


# Inchworm Algorithm



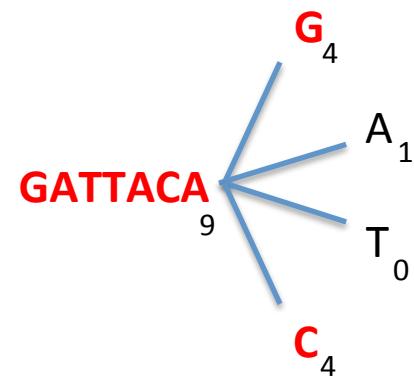


# Inchworm Algorithm



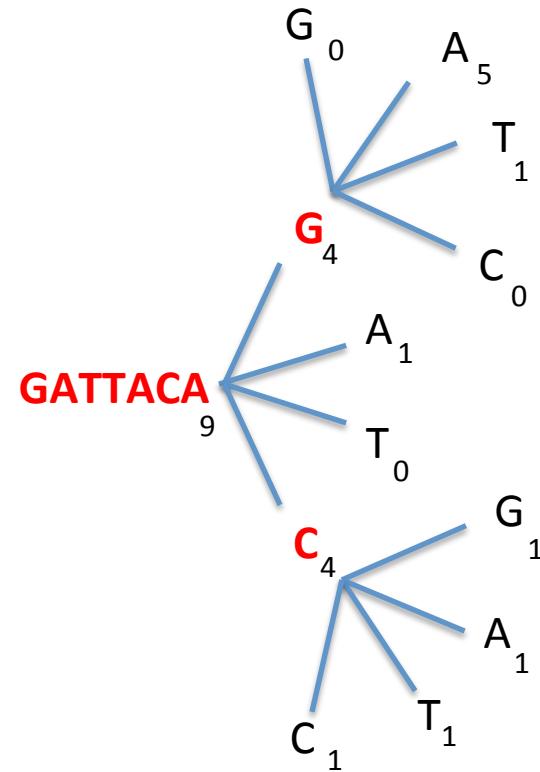


# Inchworm Algorithm



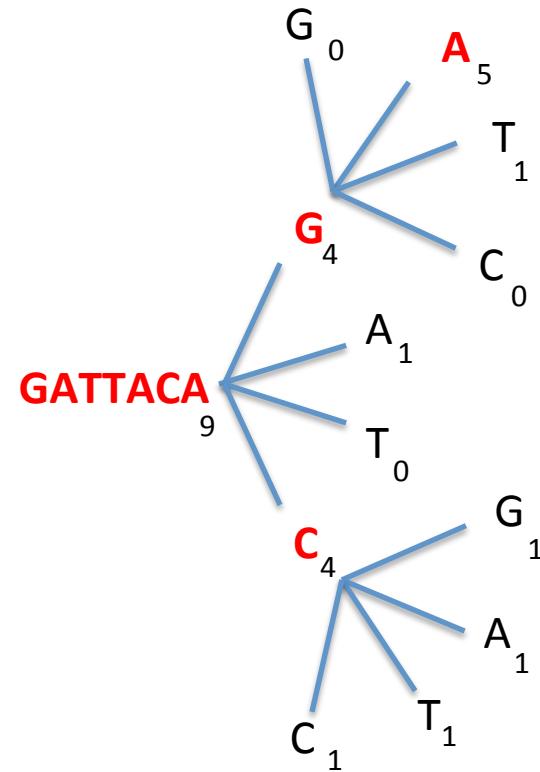


# Inchworm Algorithm



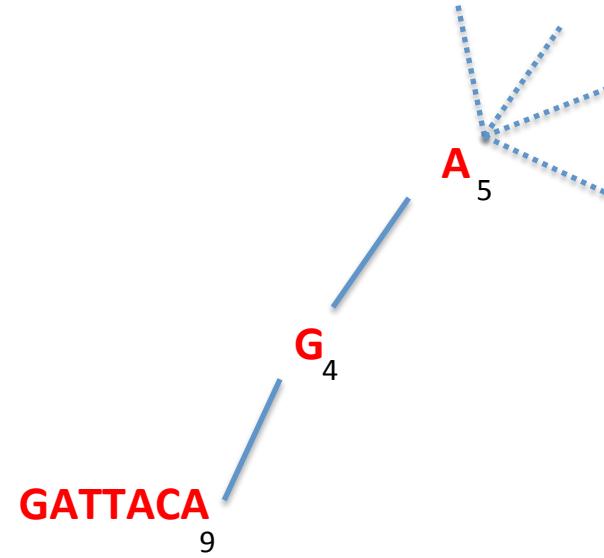


# Inchworm Algorithm



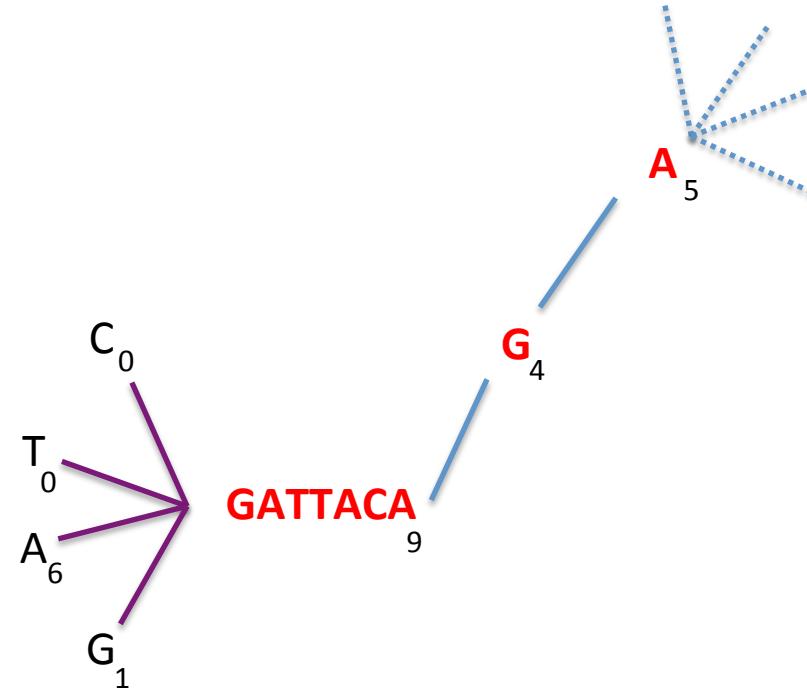


# Inchworm Algorithm



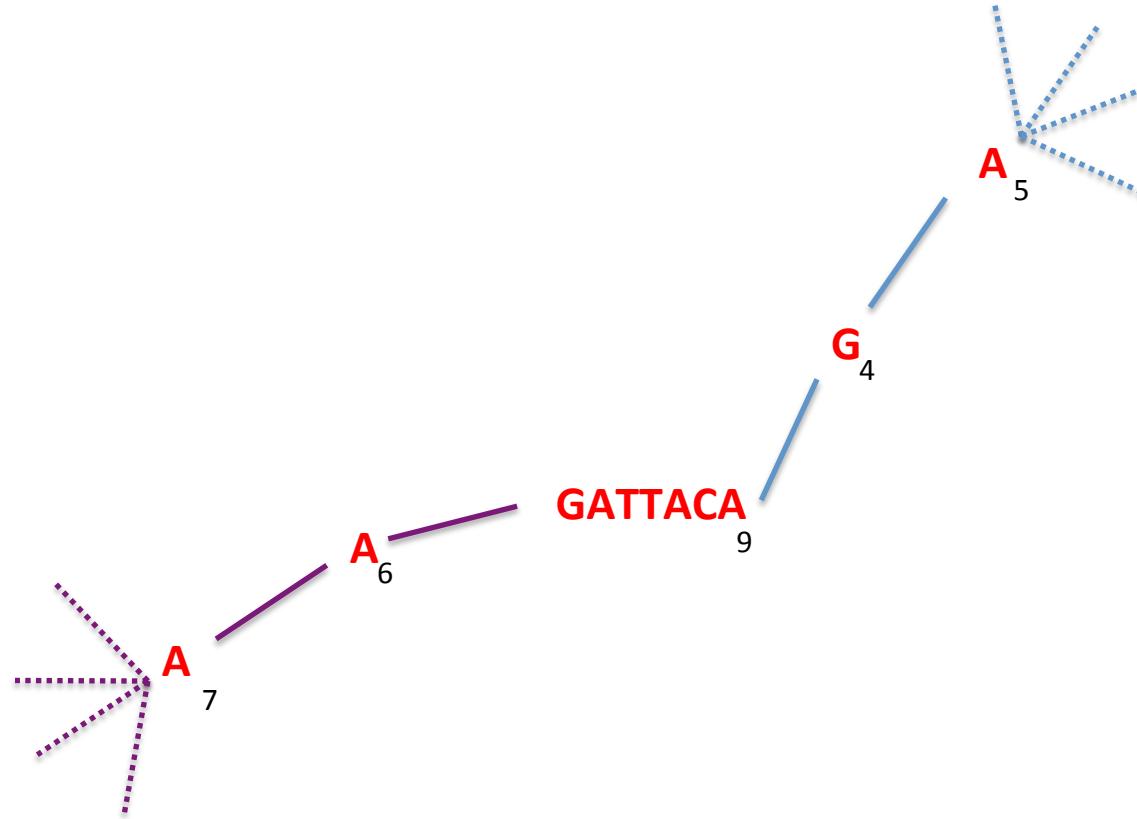


# Inchworm Algorithm





# Inchworm Algorithm

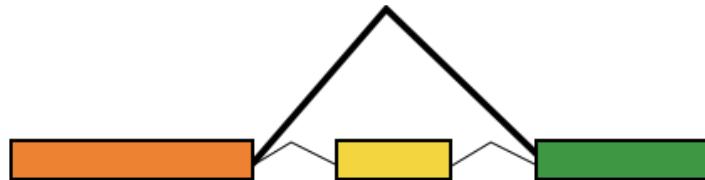


Report contig: ....**AAGATTACAGA**....

Remove assembled kmers from catalog, then repeat the entire process.

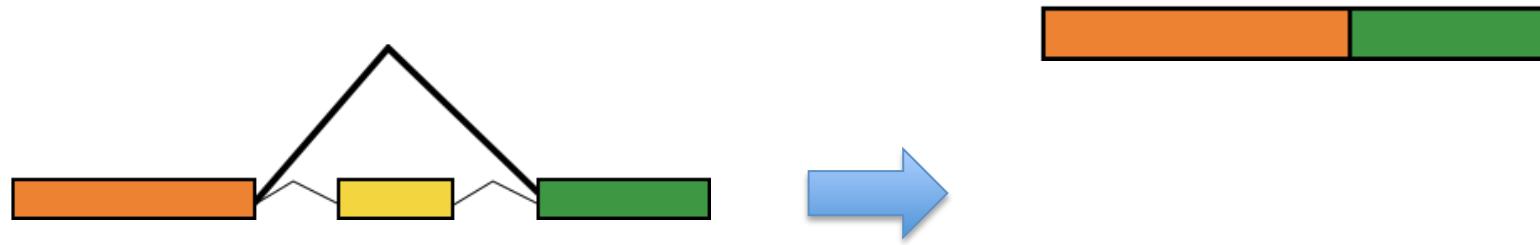


# Inchworm Contigs from Alt-Spliced Transcripts



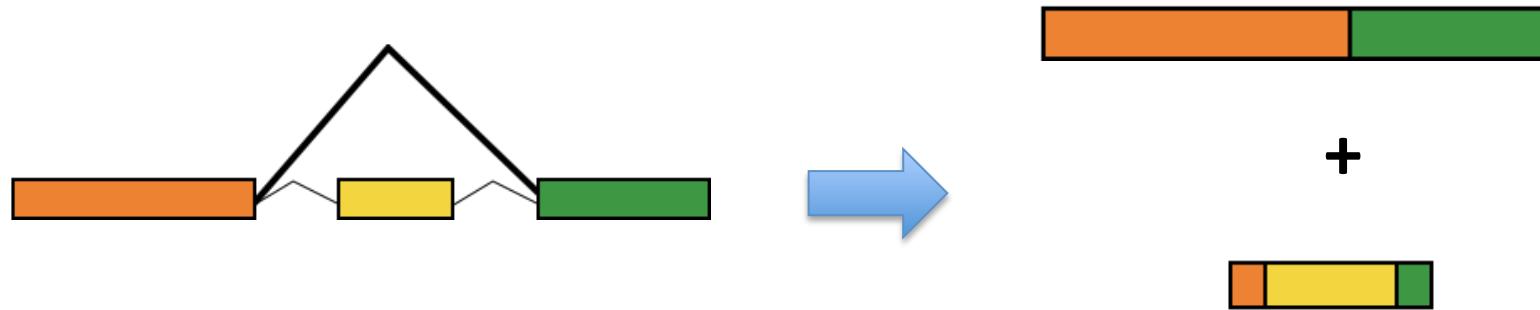


# Inchworm Contigs from Alt-Spliced Transcripts



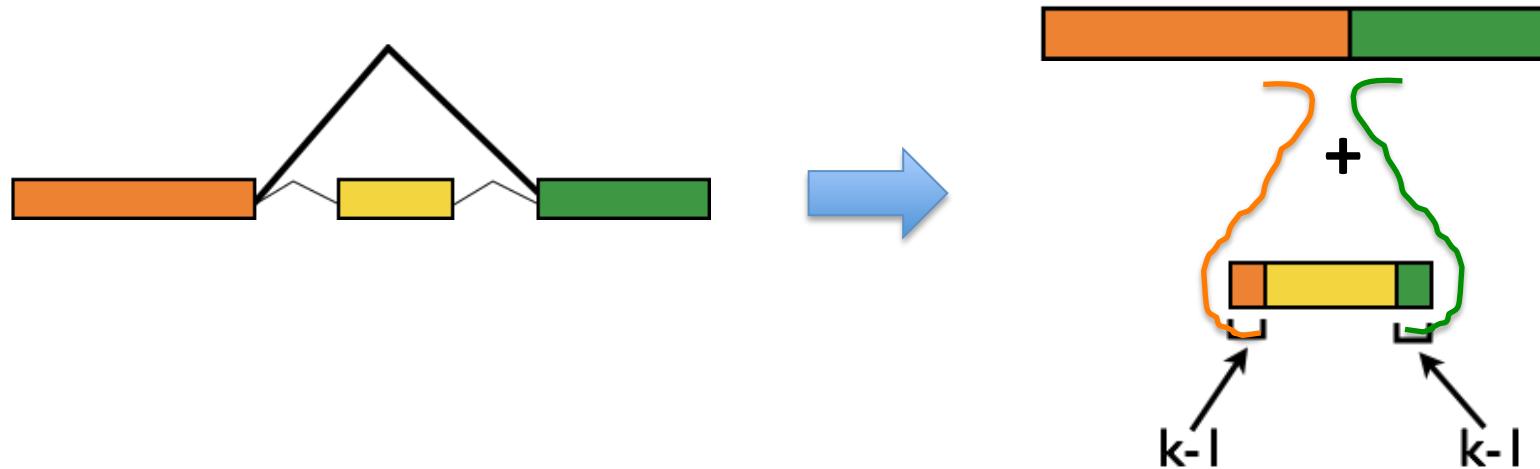


# Inchworm Contigs from Alt-Spliced Transcripts





# Inchworm Contigs from Alt-Spliced Transcripts



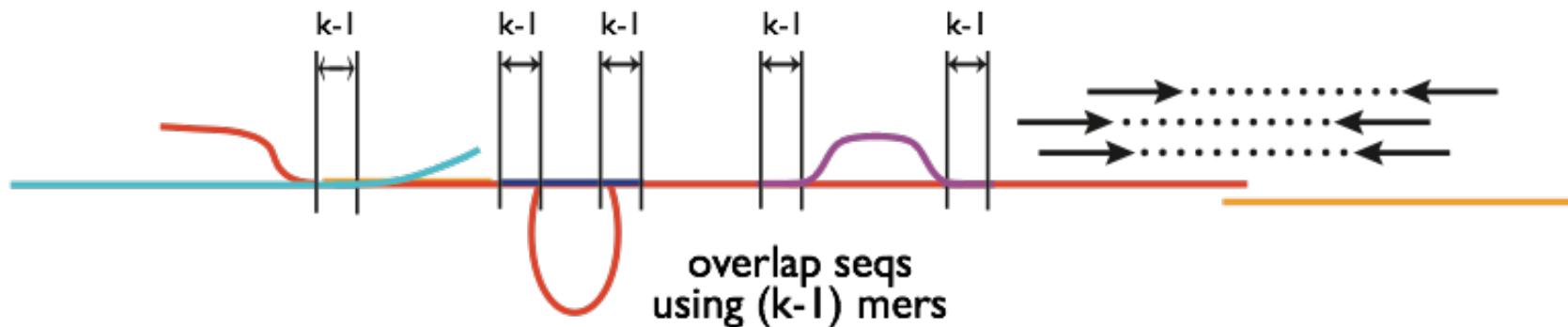
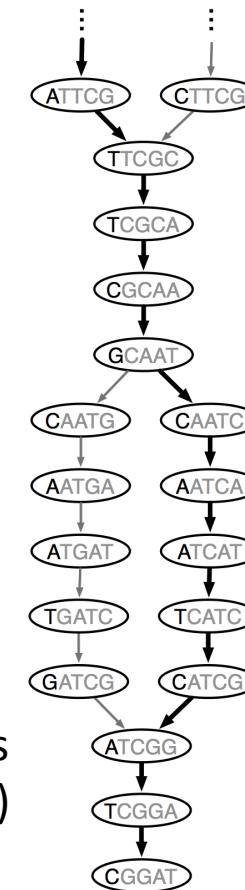
# Chrysalis

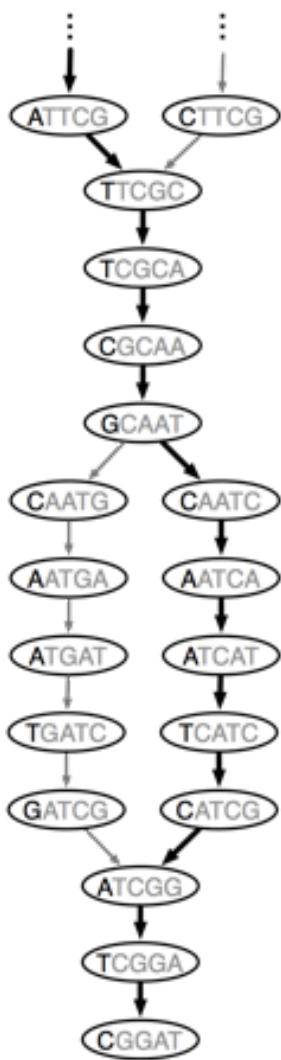
>a121:len=5845  
  |  
>a122:len=2560  
  |  
>a123:len=4443  
  |  
>a124:len=48  
  |  
>a125:len=8876  
  |  
>a126:len=66

Integrate isoforms  
via k-1 overlaps



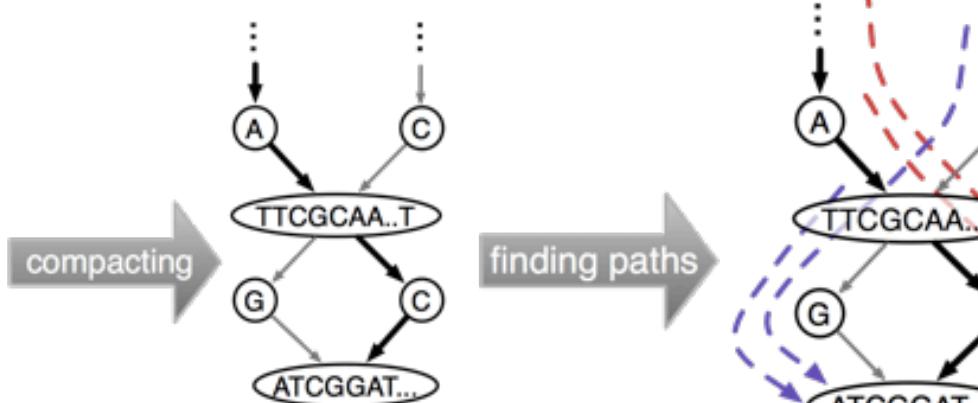
Build de Bruijn Graphs  
(ideally, one per gene)



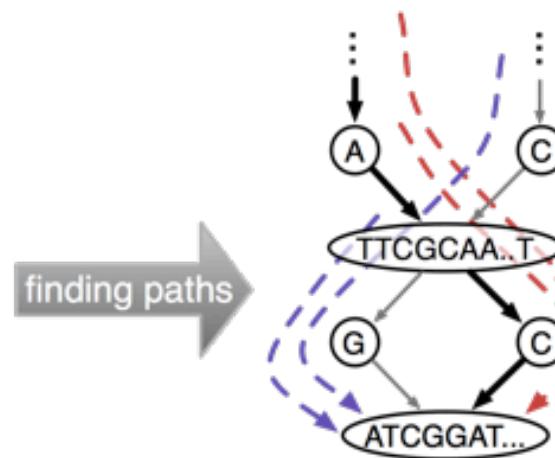


de Bruijn  
graph

# Butterfly



compact  
graph



compact  
graph with  
reads

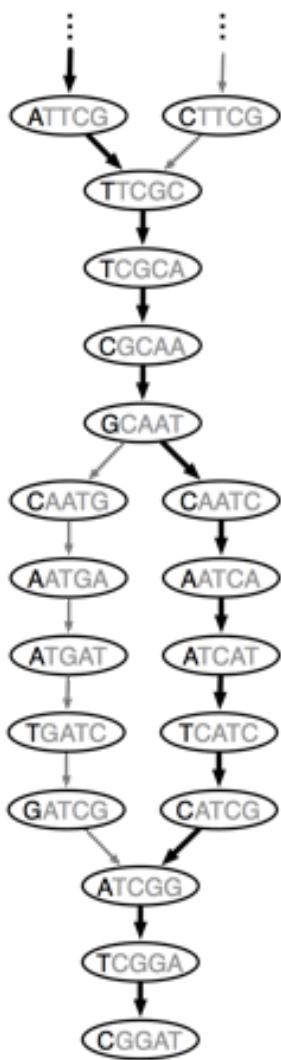


extracting  
sequences  
sequences  
(isoforms and paralogs)

..CTTCGCAA..TGATCGGAT...

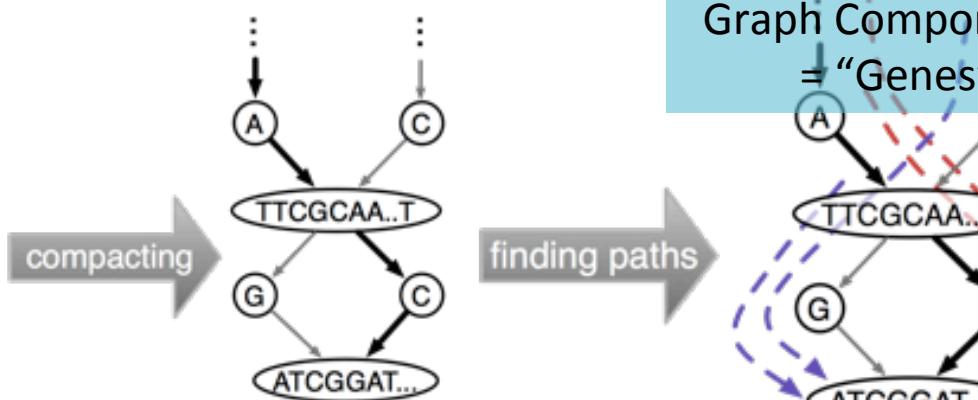
..ATTCGCAA..TCATCGGAT...

sequences  
(isoforms and paralogs)



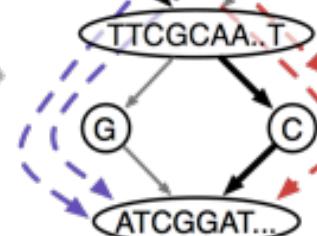
de Bruijn  
graph

# Butterfly



compact  
graph

Graph Components  
= "Genes"



finding paths

ATCGGAT...



extracting  
sequences

..CTTCGCAA..TGATCGGAT...

..ATTTCGCAA..TCATCGGAT...

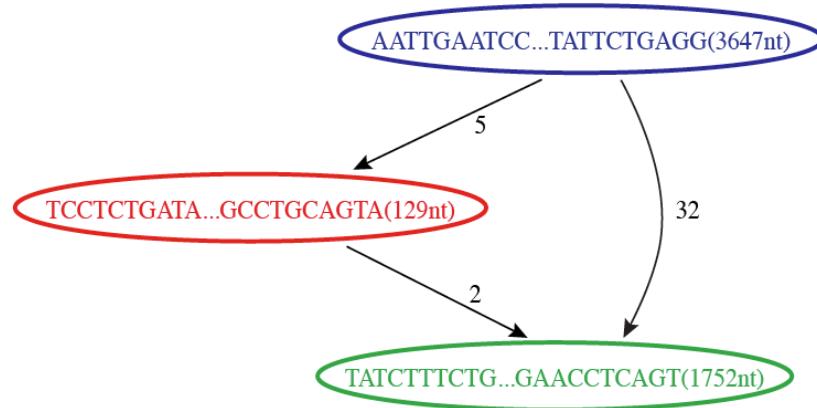
Sequences  
= "isoforms"

compact  
graph with  
reads

sequences  
(isoforms and paralogs)

# Reconstruction of Alternatively Spliced Transcripts

Butterfly's Compacted Sequence Graph

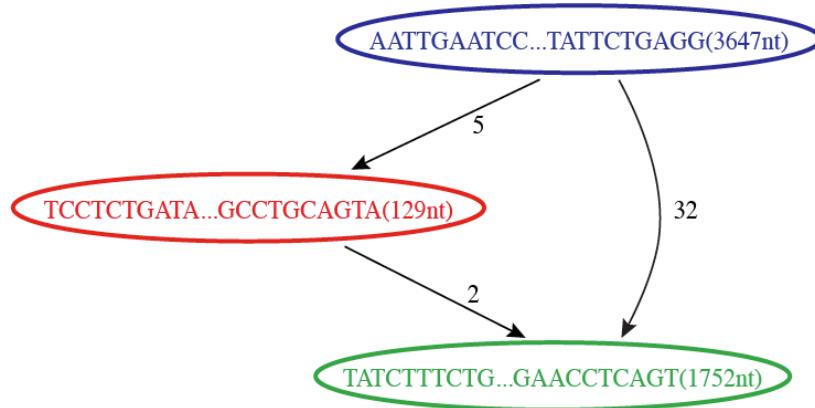


Reconstructed Transcripts



# Reconstruction of Alternatively Spliced Transcripts

Butterfly's Compacted Sequence Graph

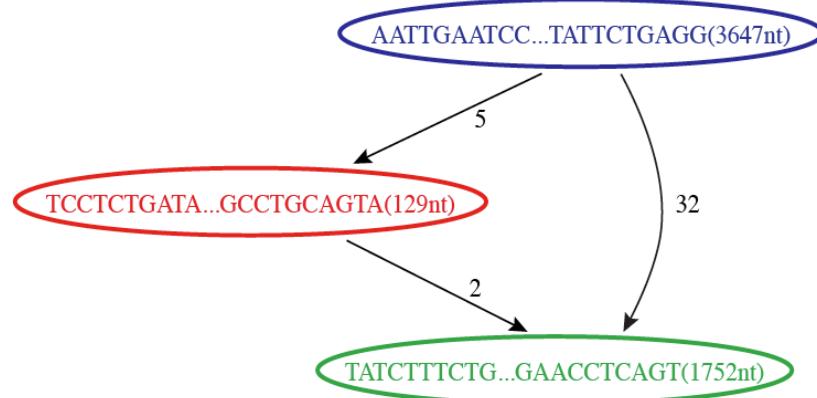


Reconstructed Transcripts



# Reconstruction of Alternatively Spliced Transcripts

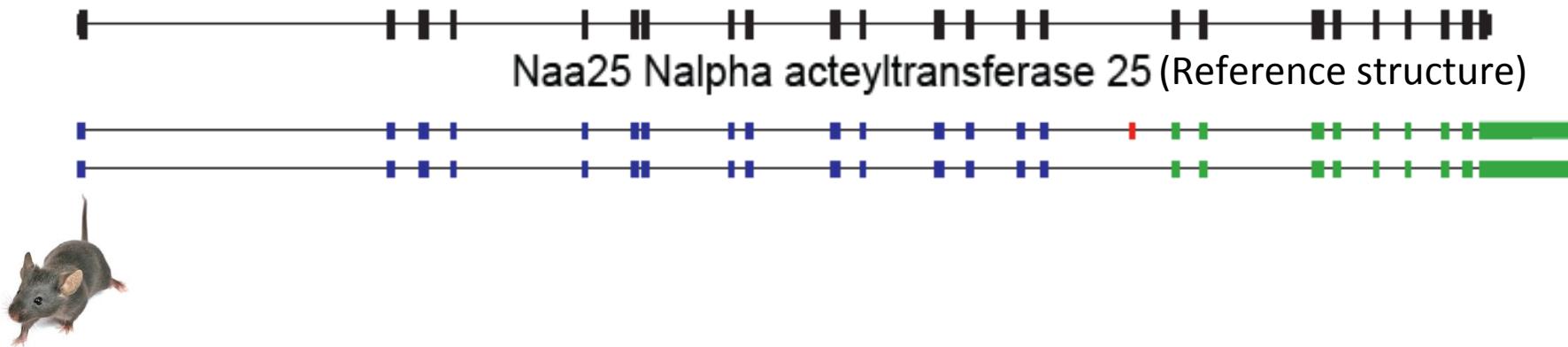
Butterfly's Compacted Sequence Graph



Reconstructed Transcripts



Aligned to Mouse Genome



# Teasing Apart Transcripts of Paralogous Genes



# Teasing Apart Transcripts of Paralogous Genes

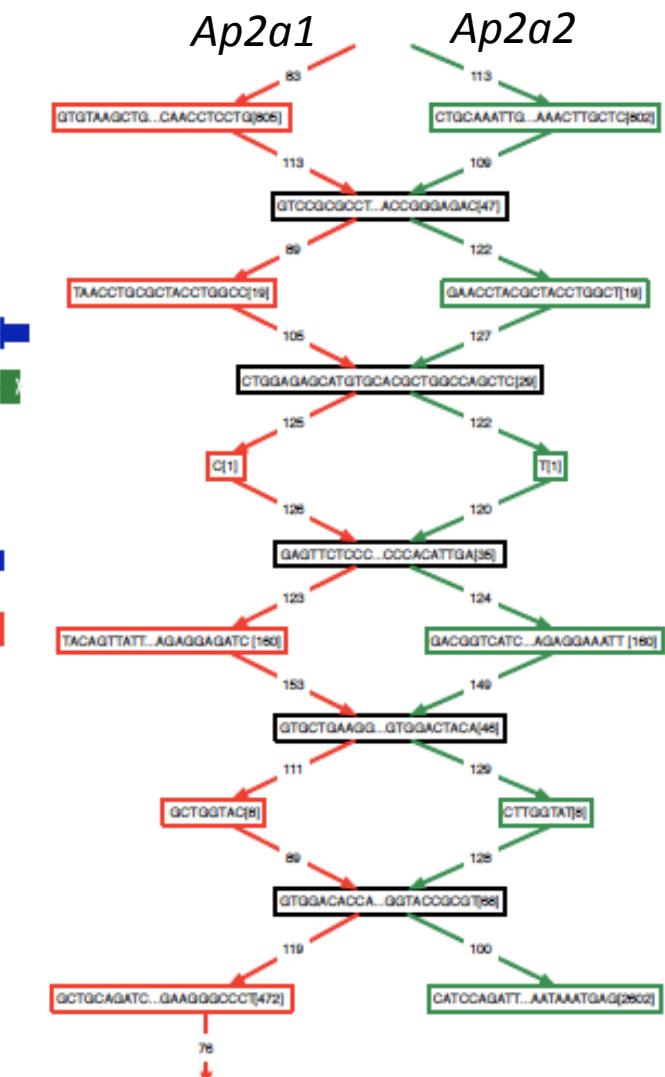
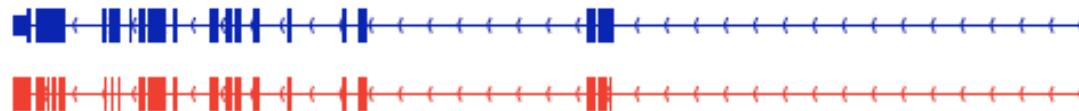
chr7:148,744,197–148,821,437

NM\_007459; Ap2a2 adaptor protein complex AP-2, alpha 2 subunit



chr7:52,150,889–52,189,508

NM\_001077264; Ap2a1 adaptor protein complex AP-2, alpha 1 subunit

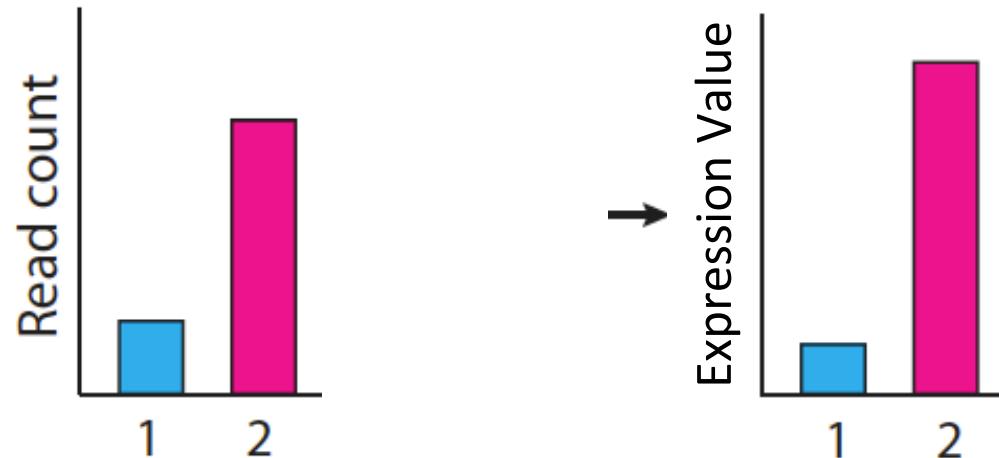
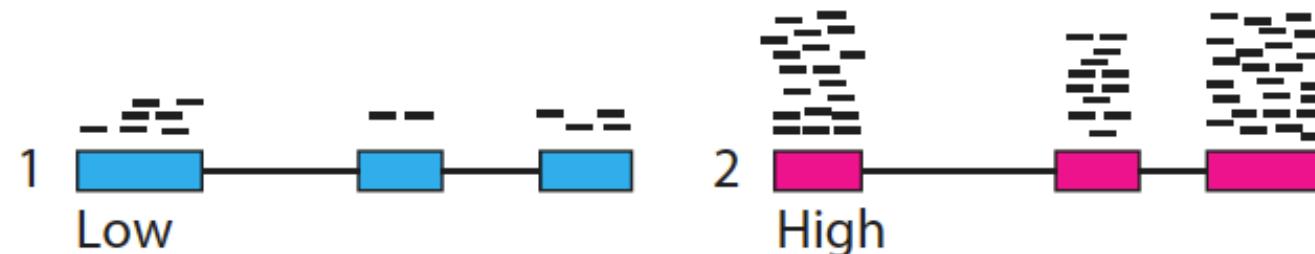


# Trinity output: A multi-fasta file

# Abundance Estimation

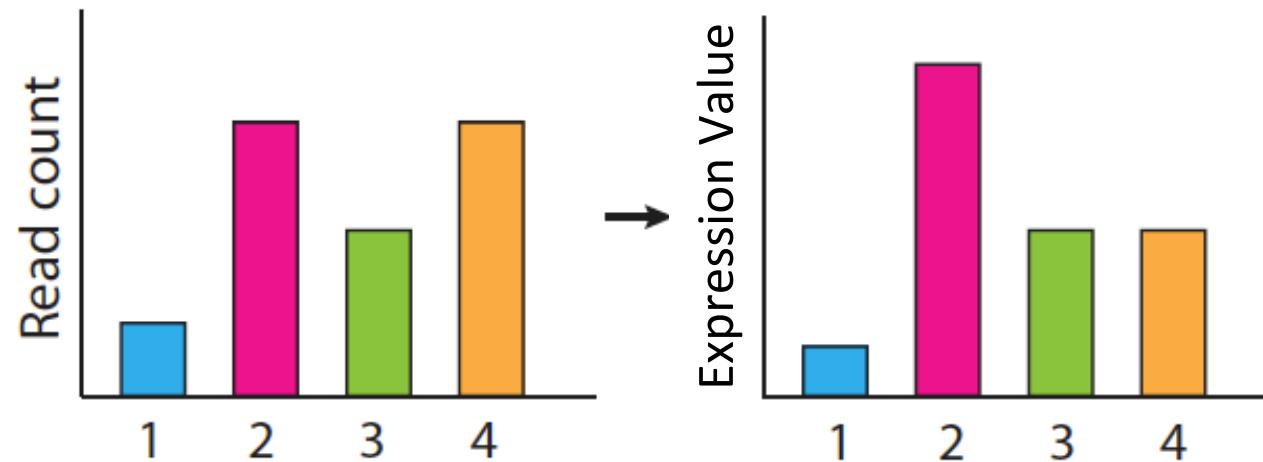
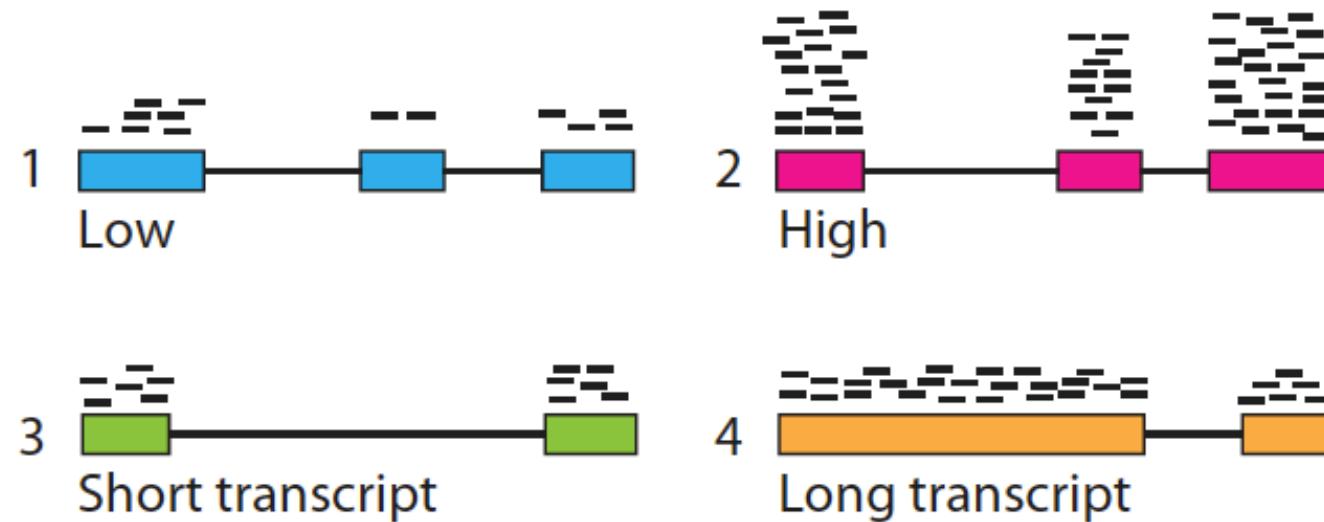
(Aka. Computing Expression Values)

# Calculating expression of genes and transcripts



Slide courtesy of Cole Trapnell

# Calculating expression of genes and transcripts



Slide courtesy of Cole Trapnell

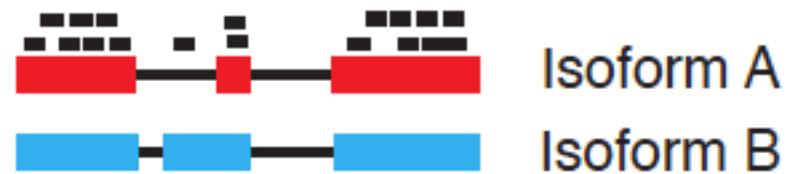
# Normalized Expression Values

- Normalized for both length of the transcript and total depth of sequencing.
- Number of RNA-Seq Fragments  
Per Kilobase of transcript  
per total Million fragments mapped

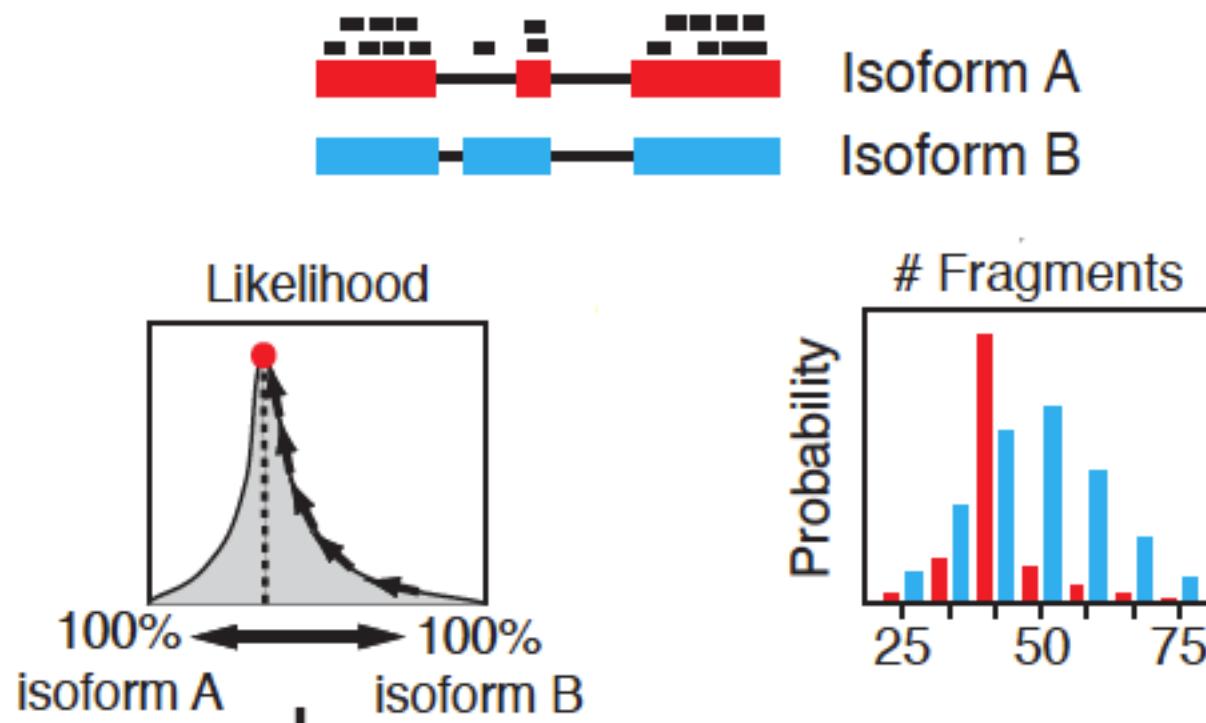
**FPKM**

Note, RPKM : Reads per ... instead of Fragments is often used with single-end reads.

Sophisticated computations are required to estimate isoform expression where there is read mapping ambiguity.



Sophisticated computations are required to estimate isoform expression where there is read mapping ambiguity.



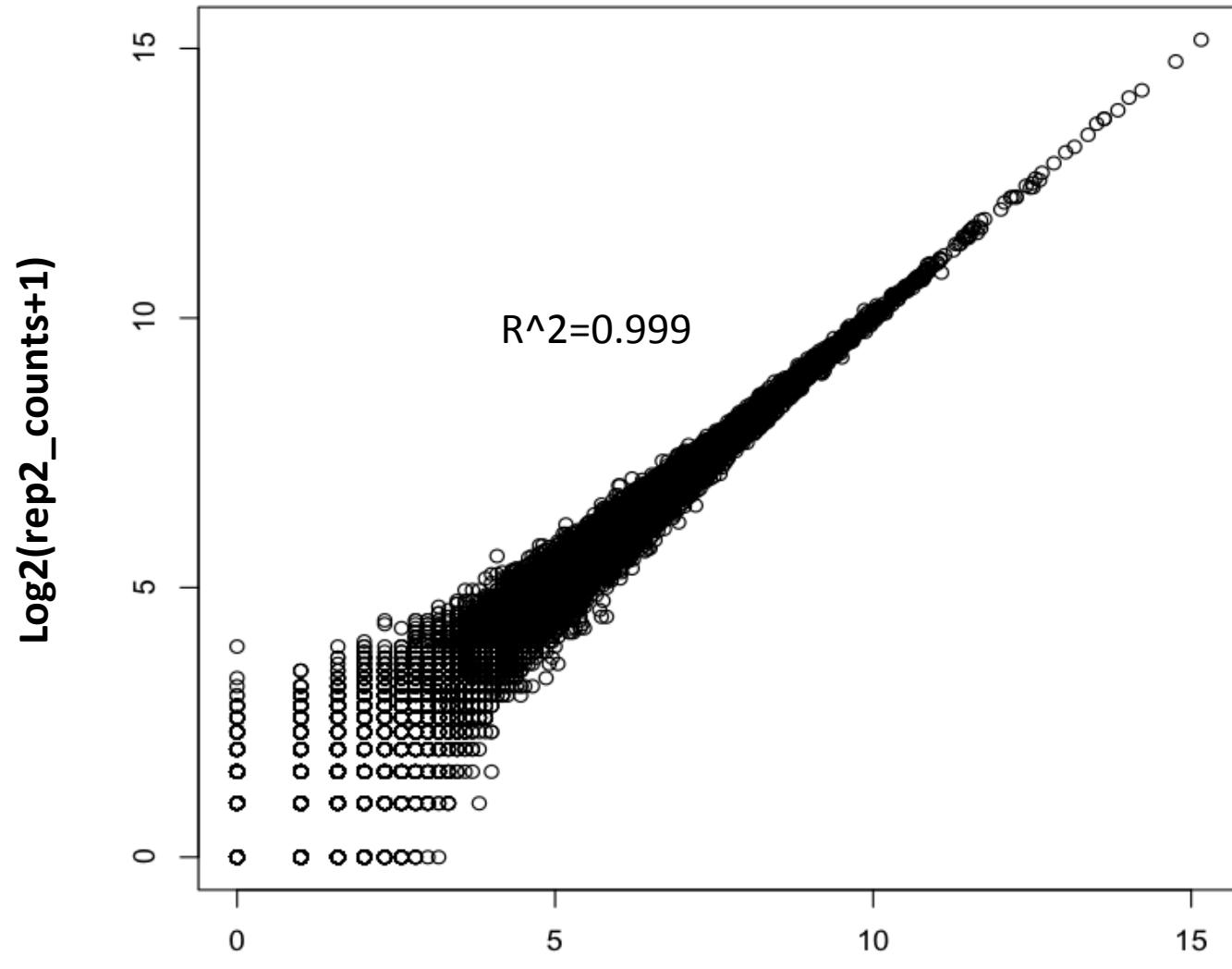
Model considers unique and ambiguously mapping reads and the length of transcripts.

Illustrations courtesy of Cole Trapnell

# Tools that perform abundance estimation

Cuffdiff			RSEM		
0	tracking_id	XLOC_000001	0	transcript_id	comp100_c0_seq1
1	class_code	-	1	gene_id	comp100_c0
2	nearest_ref_id	-	2	length	727
3	gene_id	XLOC_000001	3	effective_length	534.74
4	gene_short_name	-	4	expected_count	14.00
5	tss_id	TSS1	5	TPM	328.11
6	locus	Chr1:180422-180902	6	FPKM	532.77
7	length	-	7	IsoPct	100.00
8	coverage	-			
9	condA_FPKM	10042.1			
10	condA_conf_lo	0			
11	condA_conf_hi	20319.6			
12	condA_status	OK			

# Technical Replicates



Kidney rna-seq samples

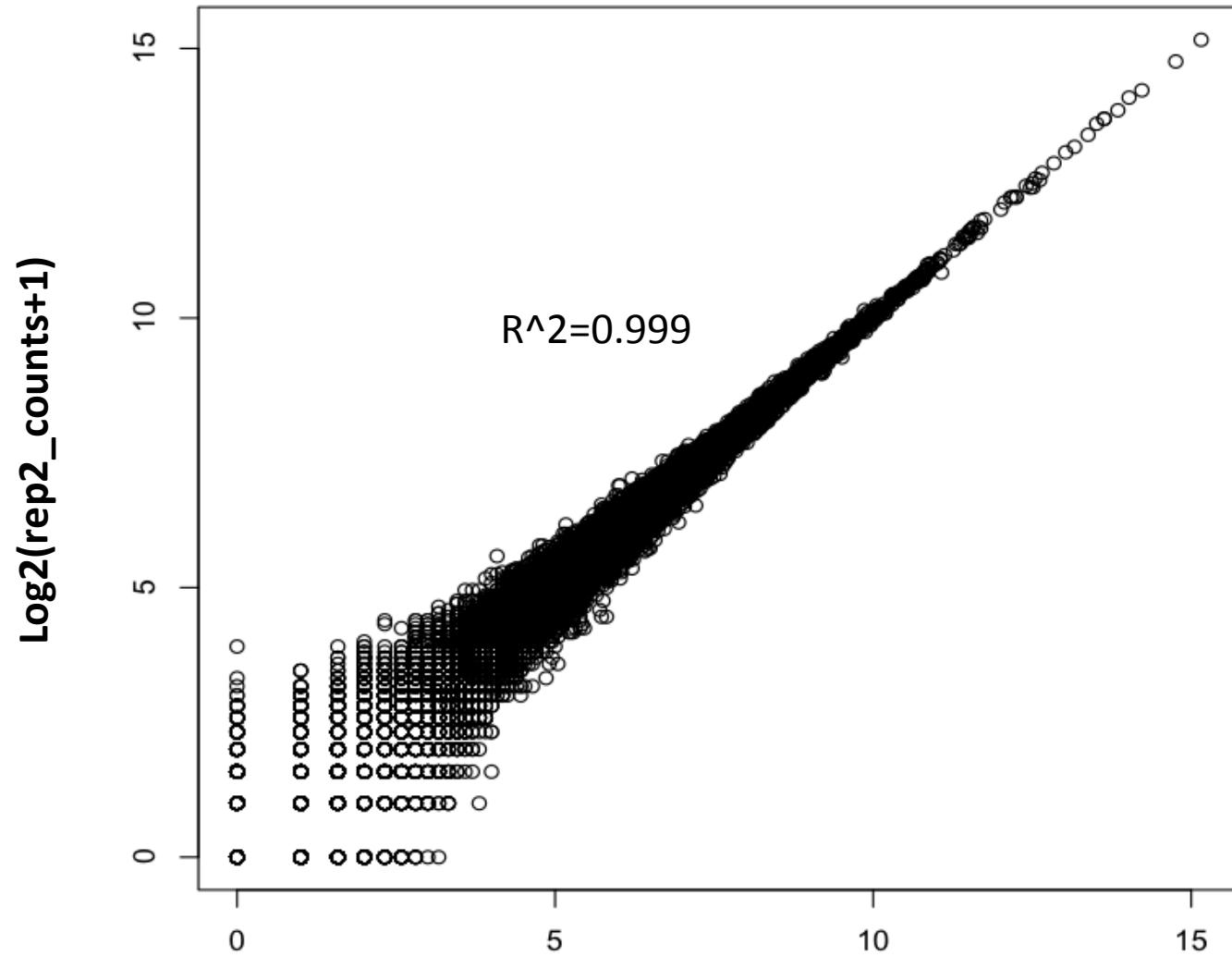
Data from Marioni et al. [2008]

$\text{Log2}(\text{rep1\_counts}+1)$



See above the toasters in  
Blackford Hall, CSHL

# Technical Replicates



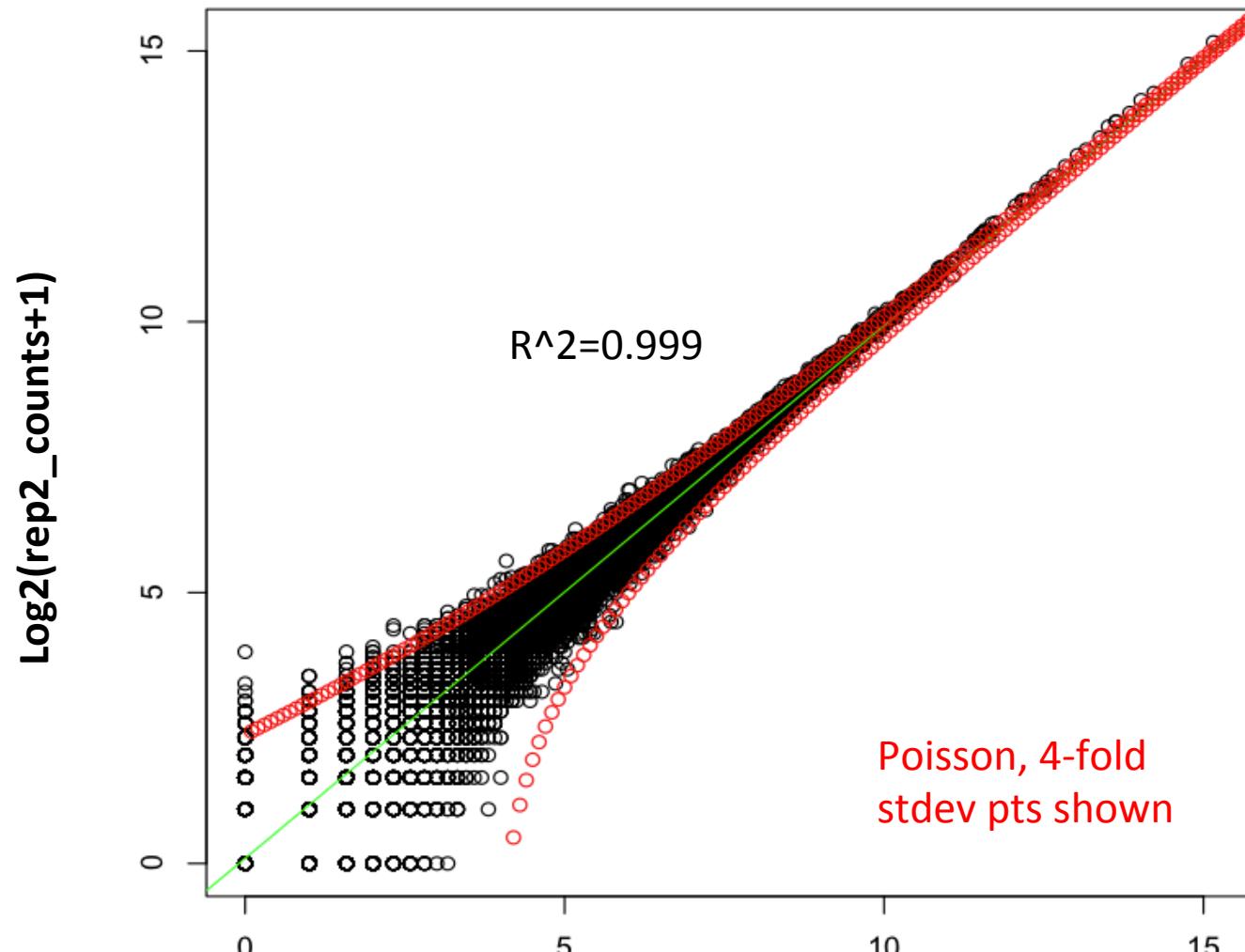
Kidney rna-seq samples

Data from Marioni et al. [2008]

$\text{Log2}(\text{rep1\_counts}+1)$

# Technical Replicates

Variation observed matches expectations due to random sampling (Poisson distribution)



Kidney rna-seq samples

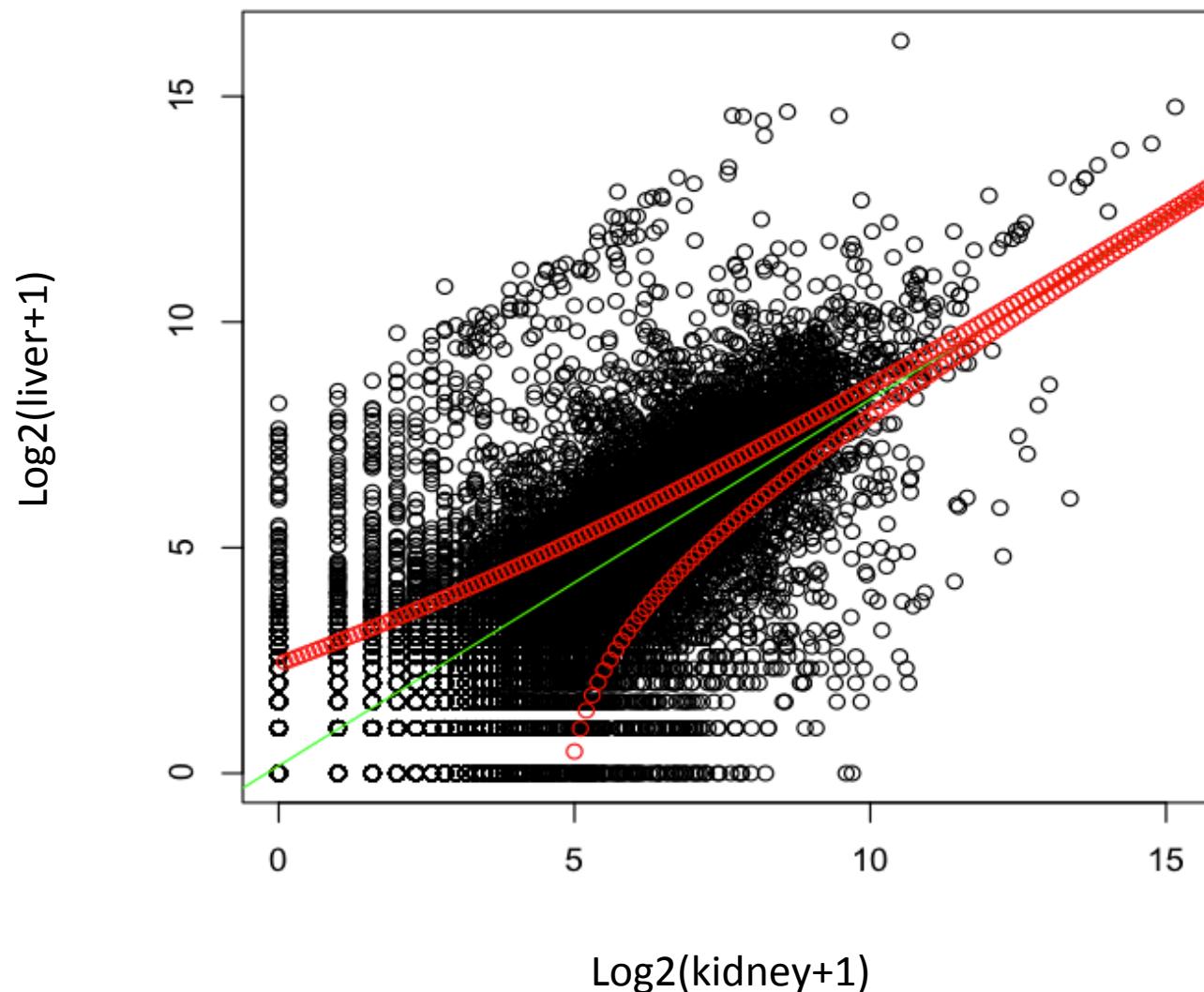
Data from Marioni et al. [2008]

Log2(rep1\_counts+1)

- Poisson well-describes variation observed in technical replicates.
- Negative binomial (overly dispersed poisson) better models biological replicates.

# Comparing Samples and Identifying Differentially Expressed Transcripts

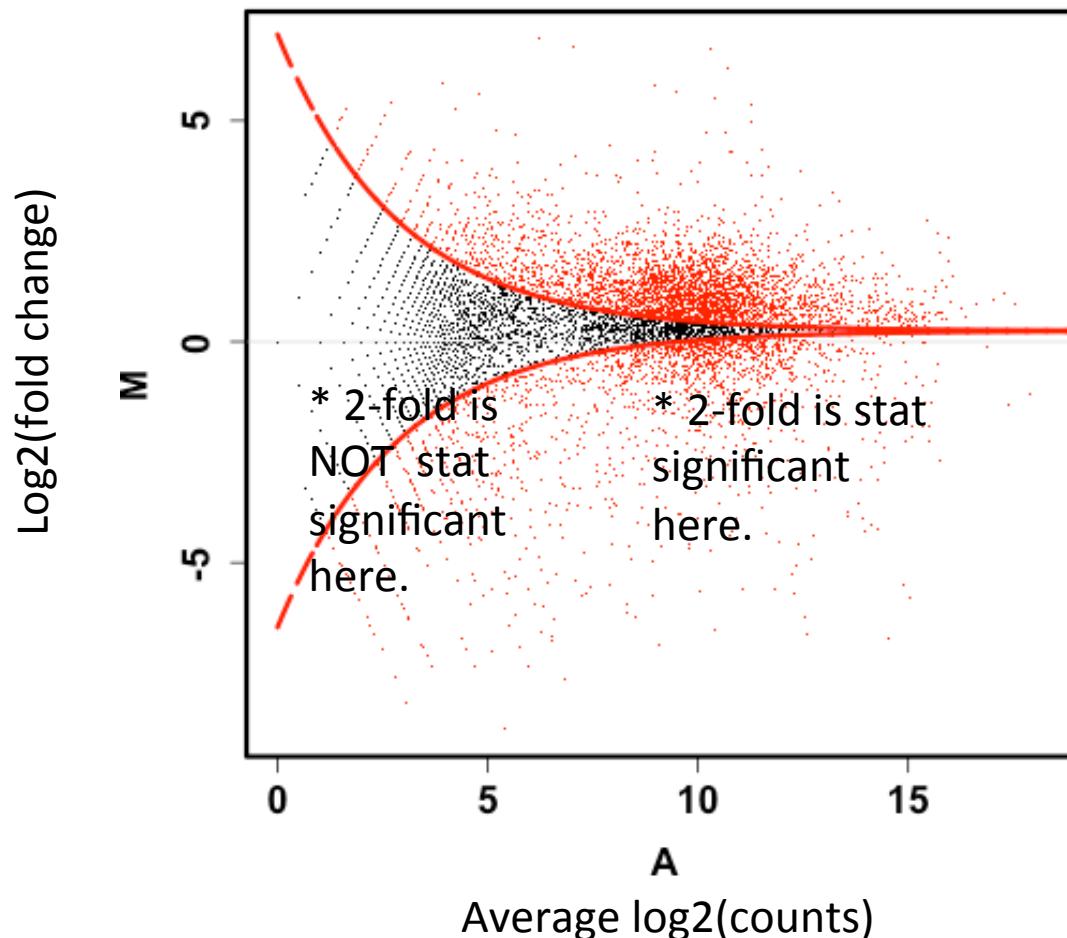
# Kidney vs. Liver



# Increased Power for Identifying Differentially Expressed Transcripts With Deeper Sequencing

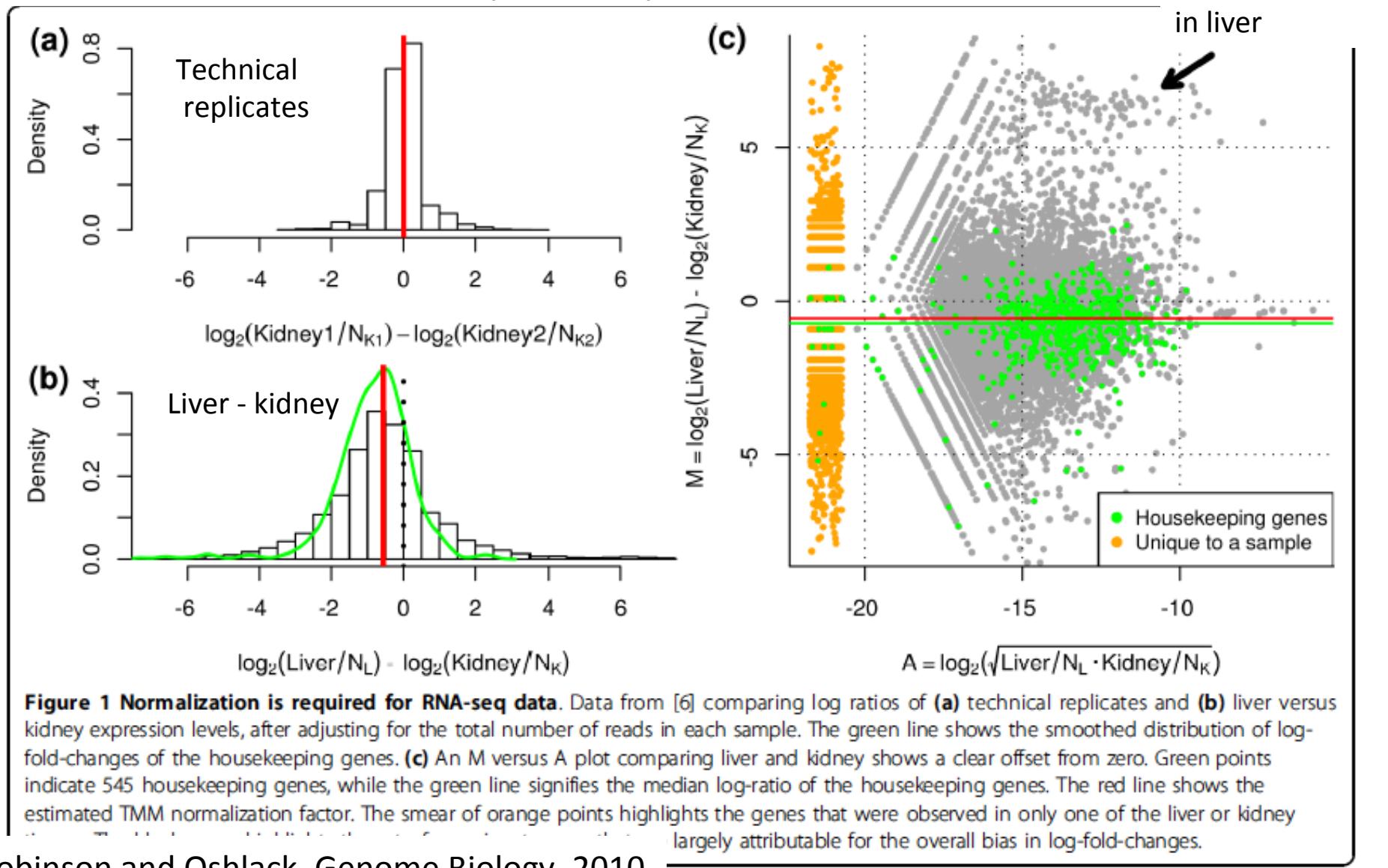
MA plot:  $\log(\text{Counts})$  vs.  $\log(\text{Fold change})$

Log Phase VS Heat Shock



# Normalization Required

Otherwise, housekeeping genes look diff expressed  
due to sample composition differences



# Identifying Differentially Expressed Transcripts

- Statistical tests performed on fragment counts (not FPKM values).
- Given observed read counts for a transcript in each of two samples, what's the probability they were derived from the same distribution (null hypothesis)? (ex. Fishers exact test)  
If ( $P \leq 0.05$ ), significantly different
- Don't forget to adjust P-values due to false discovery rate (FDR) resulting from running many (thousands of) statistical tests. (ex. use Q-values)

# Experimental Design

- Forego technical replicates
- Ideally, have at least 3 biological replicates
- Without biological replicates, can still model variation based on parametric distributions (eg. Negative binomial), but expect lower accuracy.

# Statistical Analysis Software for Identifying Differentially Expressed Transcripts

- Bioconductor
  - EdgeR
  - DEGseq
  - DESeq
  - And others...
- Tuxedo suite
  - Cuffdiff
  - (analysis enabled with CummeRbund/Bioconductor)

# Examples of Results

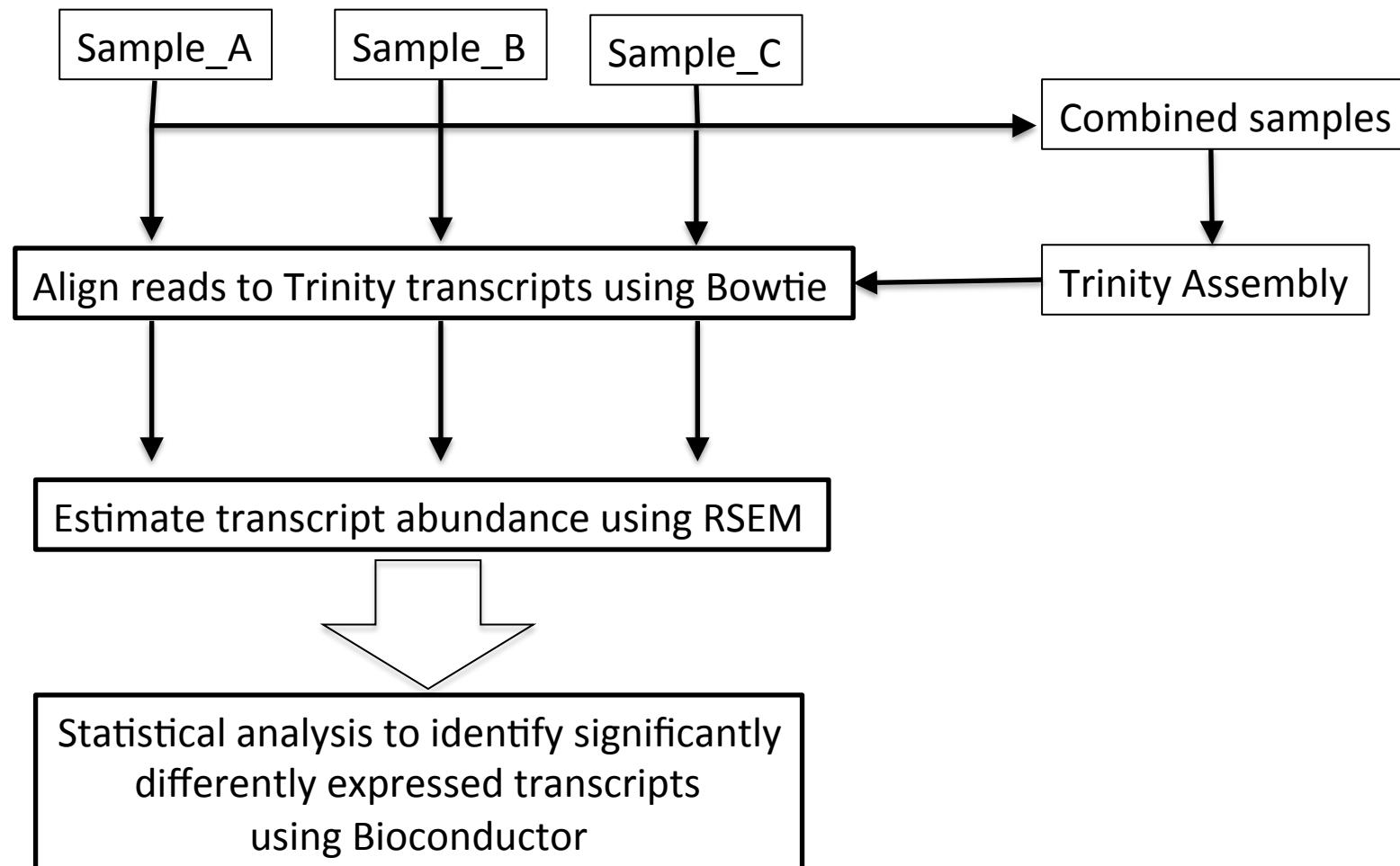
(Cuffdiff)

0	test_id	XLOC_000024
1	gene_id	XLOC_000024
2	gene	-
3	locus	7000000090838467:1335927-1338056
4	sample_1	condA
5	sample_2	condB
6	status	OK
7	value_1	680.167
8	value_2	68932
9	log2(fold_change)	6.66314
10	test_stat	-2.91993
11	p_value	0.00350111
12	q_value	0.0424377
13	significant	yes

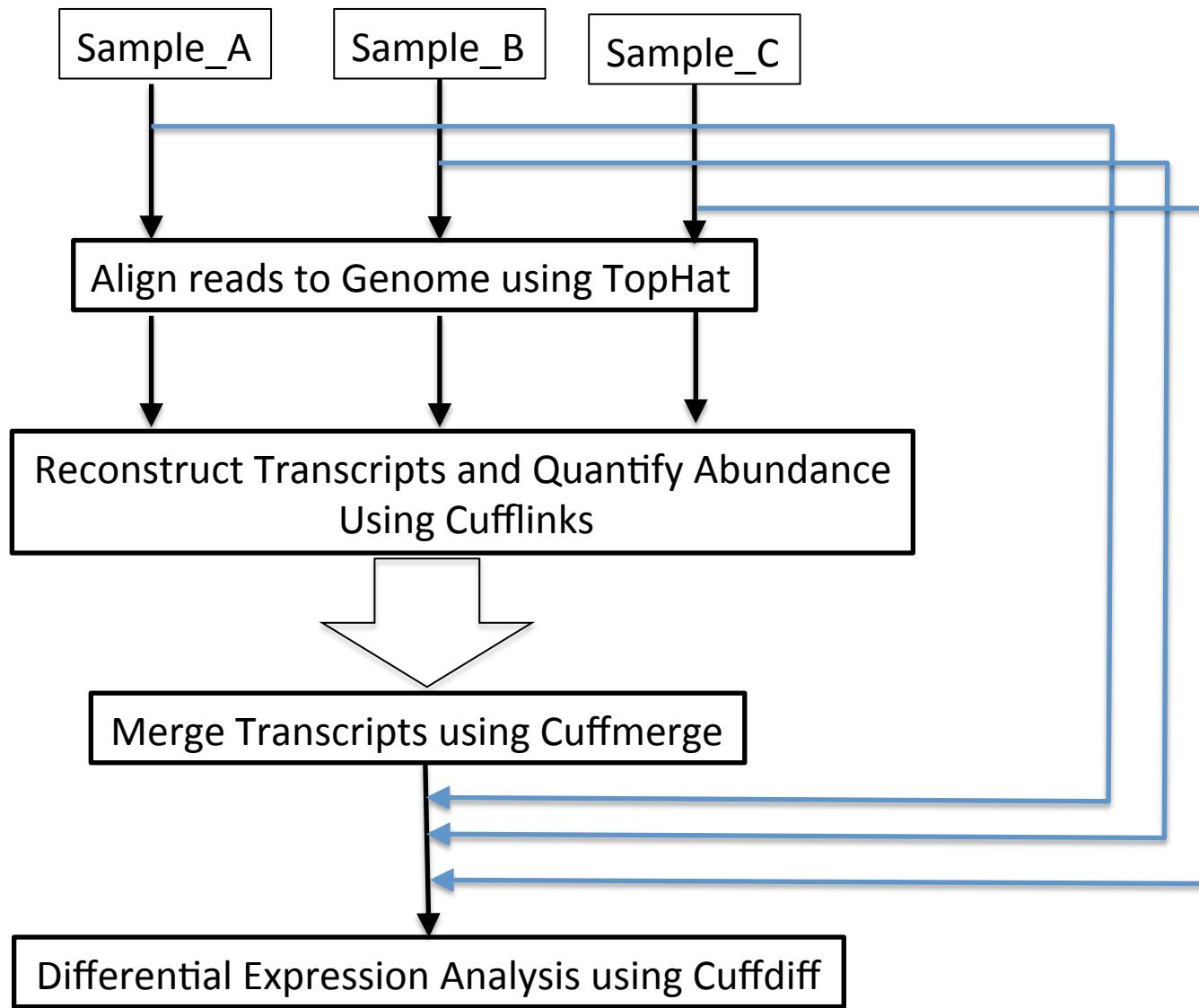
# Examples of Results (example edgeR)

		comp217_c0_seq1
0	logFC	6.69056684523186
1	logCPM	16.1146897543805
2	PValue	2.06844466442231e-15
3	FDR	9.01969581996253e-13

# Trinity Differential Expression Workflow

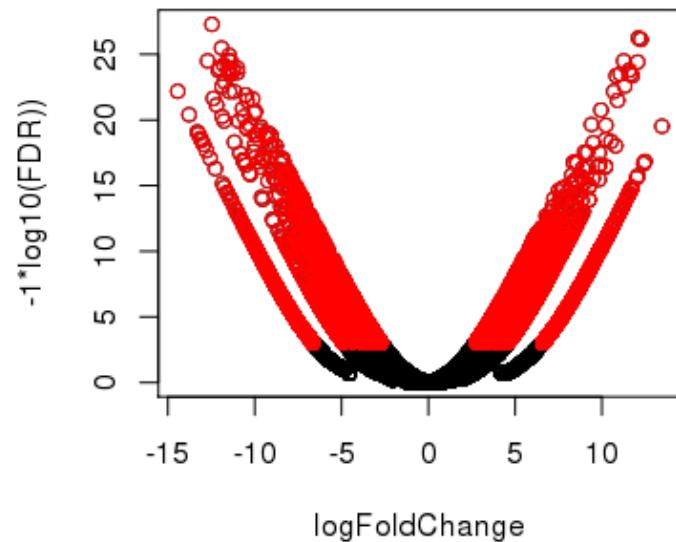


# Tuxedo Differential Expression Workflow

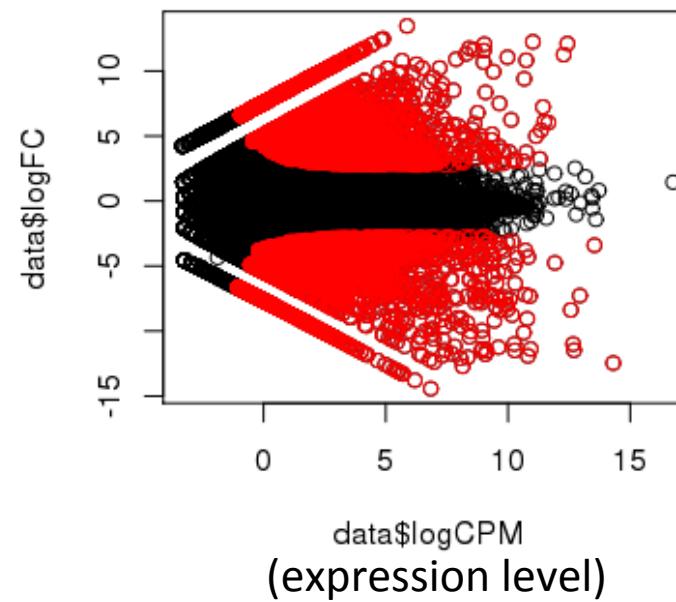


# Plotting Pairwise Differential Expression Data

Volcano plot  
( fold change vs. significance)



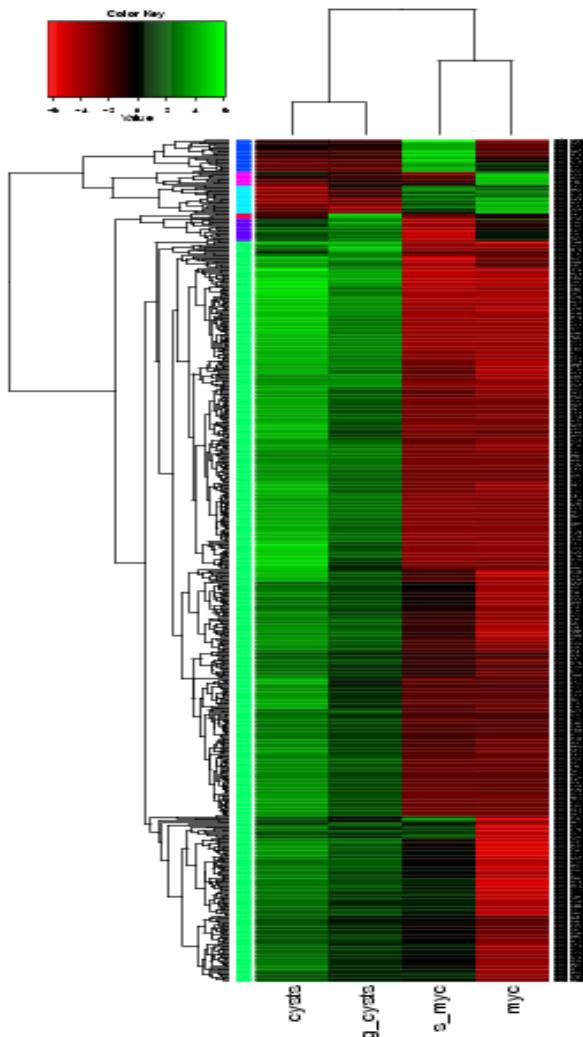
MA plot  
(abundance vs. fold change)



Significantly differently expressed transcripts have FDR  $\leq 0.001$   
(shown in red)

No replicates available, so modeled by edgeR using the  
Negative Binomial with dispersion manually set to 0.1

# Comparing Multiple Samples

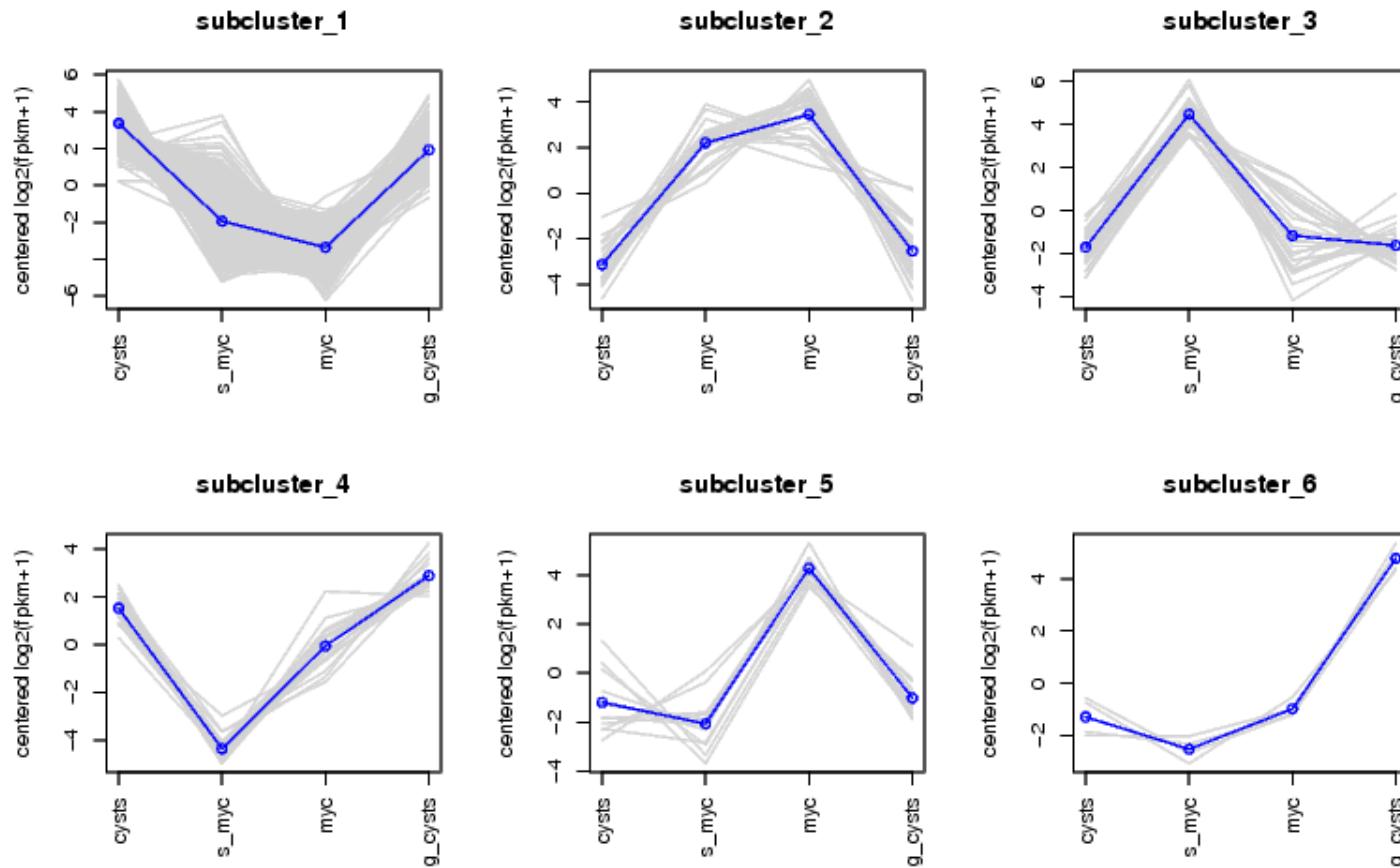


**Heatmaps** provide an effective tool for navigating differential expression across multiple samples.

**Clustering** can be performed across both axes:  
-cluster transcripts with similar expression patterns.  
-cluster samples according to similar expression values among transcripts.

# Examining Patterns of Expression Across Samples

Can extract clusters of transcripts and examine them separately.



# Hands-on Tutorials

- Tuxedo
  - Tophat alignment
  - Cufflinks transcript reconstructions
  - GenomeView for navigating the alignments
  - Cuffdiff for differential expression analysis
  - cummeRbund for exploring diff. express. results.
- Trinity
  - De novo assembly using Trinity
  - Bowtie and RSEM for abundance estimation
  - edgeR for differential expression analysis