

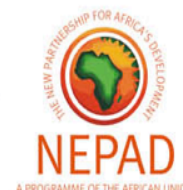
***Advanced Genomics - Bioinformatics Workshop***

**Mark Wamalwa**

*Beca-ILRI Hub, Nairobi, Kenya*

<http://hub.africabiosciences.org/>

[m.wamalwa@cgiar.org](mailto:m.wamalwa@cgiar.org)



7<sup>th</sup> – 18<sup>th</sup> September 2015

**biosciences**

eastern and central **africa**

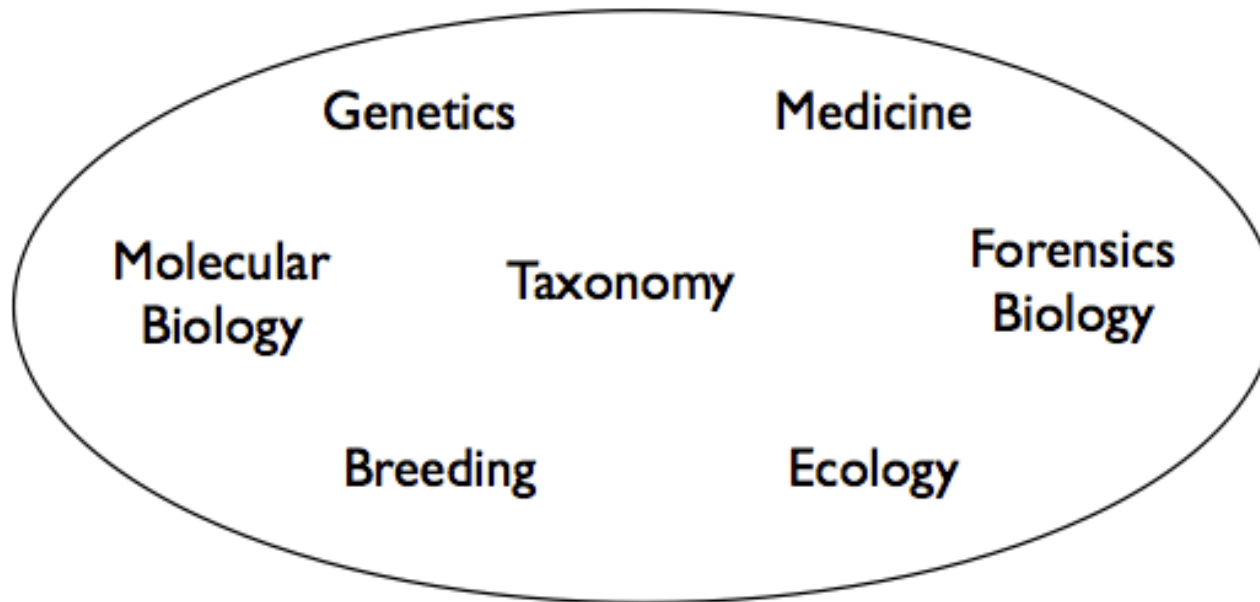
# RNA sequencing, transcriptome and expression quantification

# Lecture Overview

- What is RNA-seq?
- Basic concepts
- Mapping-based transcriptomics (genome - based)
- De novo based transcriptomics (genome-free)
- Expression counts and differential expression
- Transcript annotation

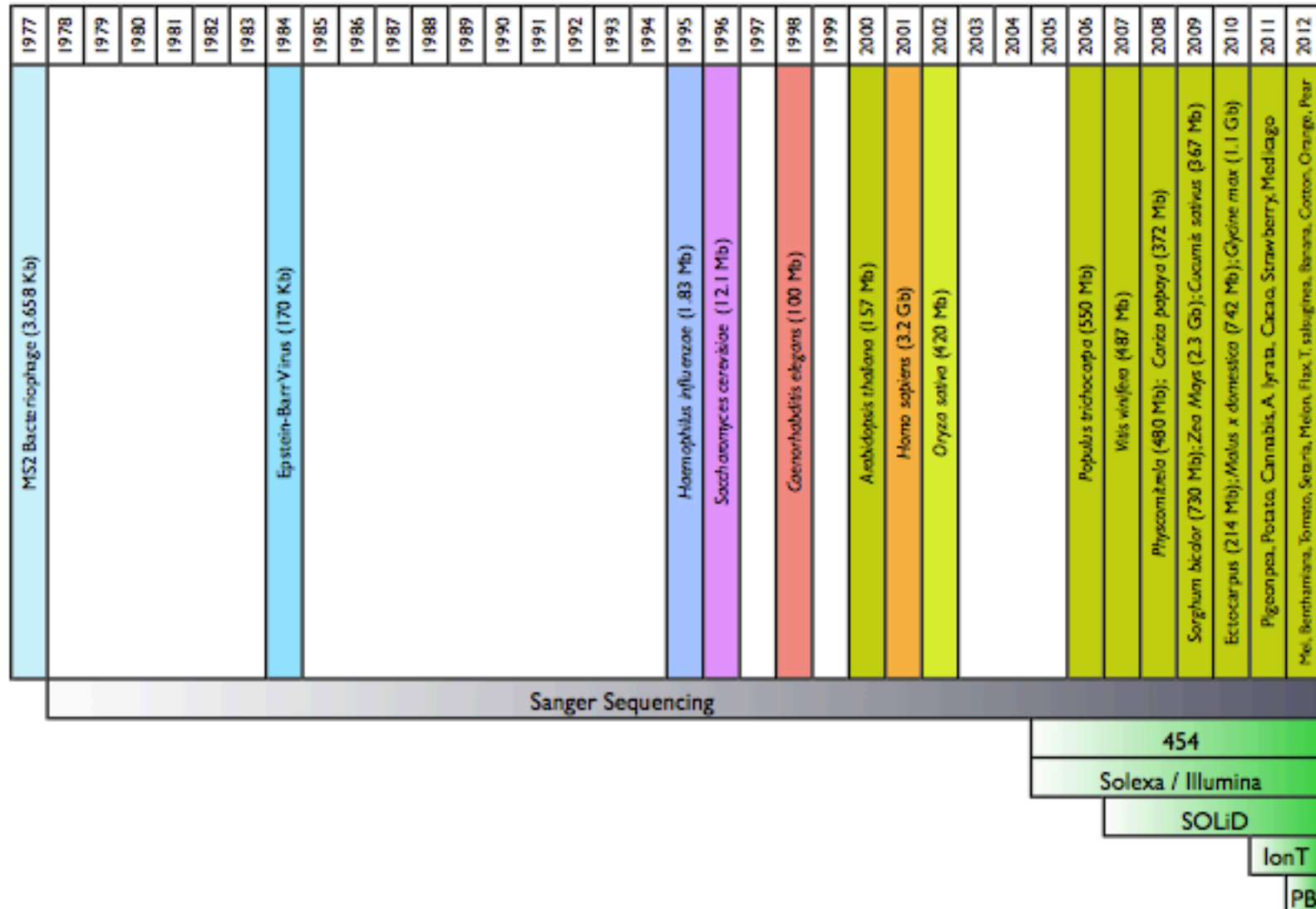
# Recap - Basics of the Next Generation Sequencing (NGS)

- ***DNA Sequencing***: “Process of determining the precise order of nucleotides within a DNA molecule.” -Wikipedia





# Recap - Basics of the Next Generation Sequencing (NGS)



# Features of Next Generation Sequencing.

1. Massive sequence production (from 0.1 to 300 Gb).
2. Wide range of sequence lengths (from 50 to 3,000 bp).
3. Same or bigger error rate than the traditional sequencing (from 87 to 99.9%).
4. Cheap price per base.

# Features of Next Generation Sequencing.

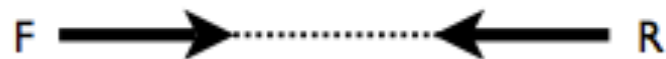
	Inputs	Outputs
454 Pyrosequencing (GS FLX Titanium XL+)	<ul style="list-style-type: none"> <li>- Single Reads Library.</li> <li>- Pair End Library (3 to 20 Kb insert size).</li> <li>- Multiplexed sample.</li> </ul>	<ul style="list-style-type: none"> <li>- sff files</li> <li>- (fasta and fastq files)</li> </ul>
Illumina (HiSeq 2500)	<ul style="list-style-type: none"> <li>- Single Reads Library.</li> <li>- Pair End Library (170-800 bp insert size).</li> <li>- Mate Pair Library (2 to 10 Kb insert Size)</li> <li>- Multiplexed sample.</li> </ul>	<ul style="list-style-type: none"> <li>- fastq files (Phred+64)</li> <li>- fastq files (Phred+33, Illumina 1.8+)</li> </ul>
Illumina (MiSeq)		
SOLID (5500xl system)	<ul style="list-style-type: none"> <li>- Single Reads Library.</li> <li>- Mate Pairs Library (0.6 to 6 Kb insert size).</li> <li>- Multiplexed sample.</li> </ul>	<ul style="list-style-type: none"> <li>- fastq files (Phred+33)</li> </ul>
Ion Torrent (Ion Proton I)	<ul style="list-style-type: none"> <li>- Single Reads Library.</li> <li>- Pair End Library (0.6 to 6 Kb insert size).</li> <li>- Multiplexed sample.</li> </ul>	
PacBio (PacBioRS)	<ul style="list-style-type: none"> <li>- Single Reads Library.</li> </ul>	

# Library types (orientations)

- Single reads

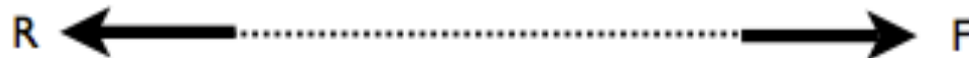


- Pair ends (PE) (150-800 bp insert size)



Illumina

- Mate pairs (MP) (2-20 Kb insert size)



Illumina

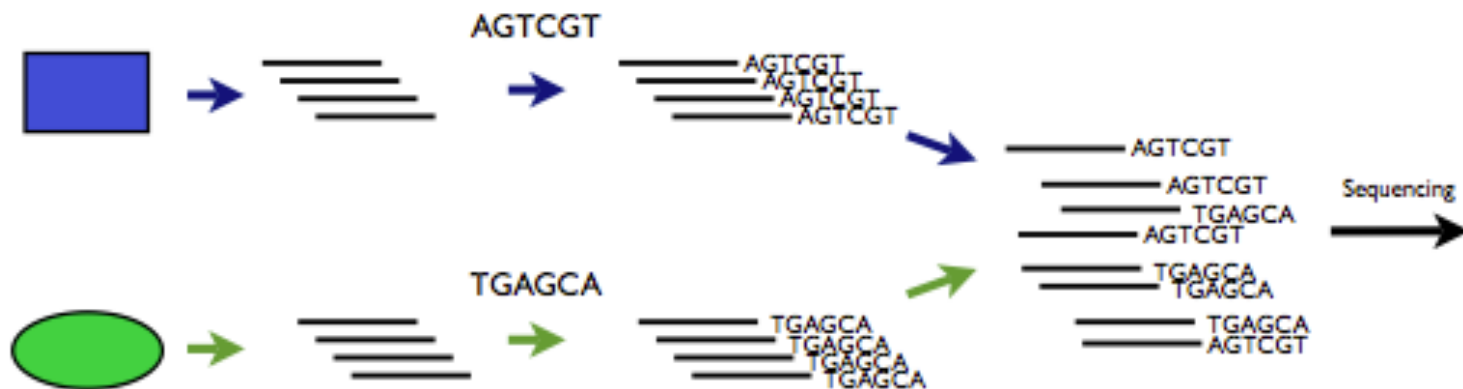


454/Roche

# Library types (orientations)

## ★ Multiplexing:

Use of different tags (4-6 nucleotides) to identify different samples in the same lane/sector.





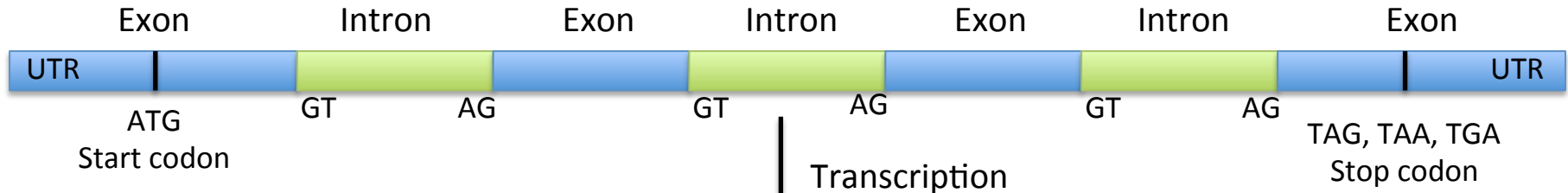
# Lecture Overview

- What is RNA-seq?
- Basic concepts
- Mapping-based transcriptomics (genome - based)
- De novo based transcriptomics (genome-free)
- Expression counts and differential expression
- Transcript annotation



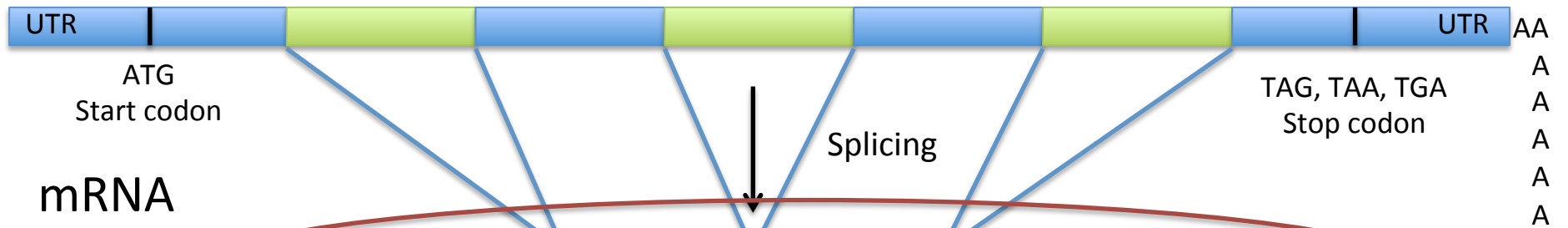
# RNA-seq

DNA



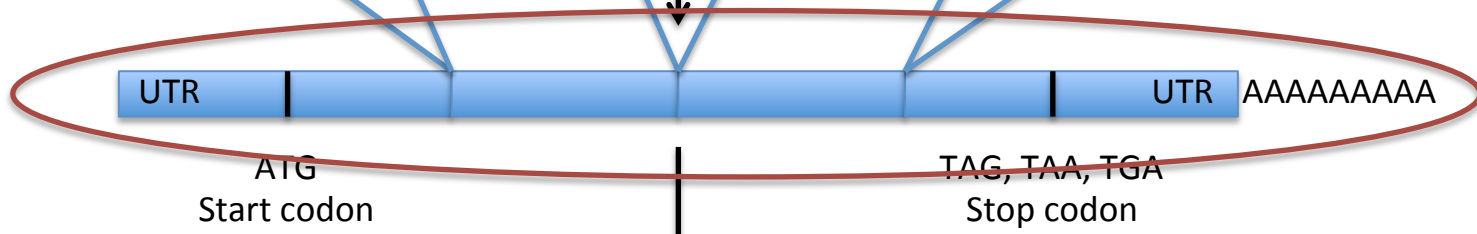
Transcription

Pre-mRNA



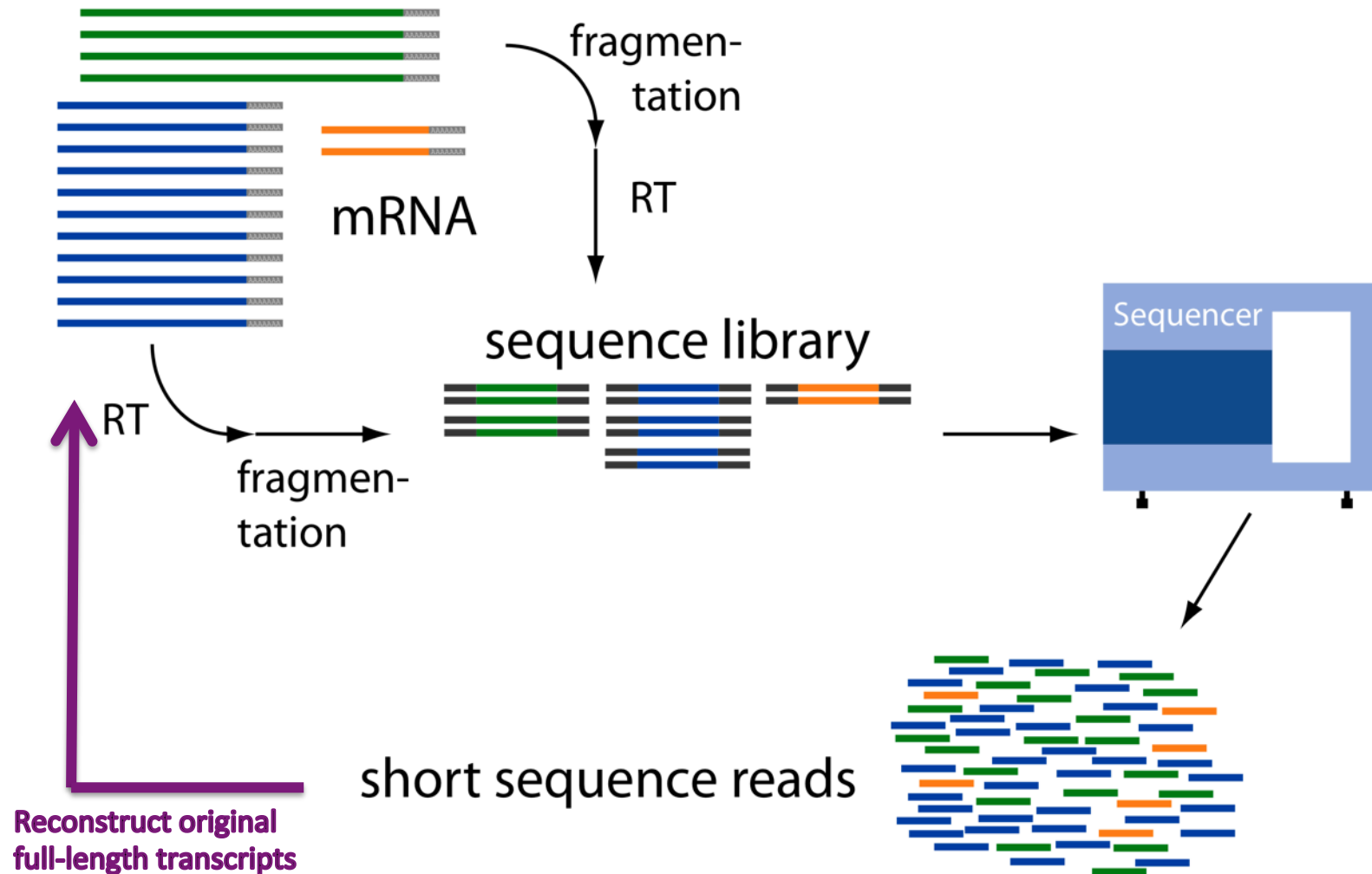
Splicing

mRNA



Translation

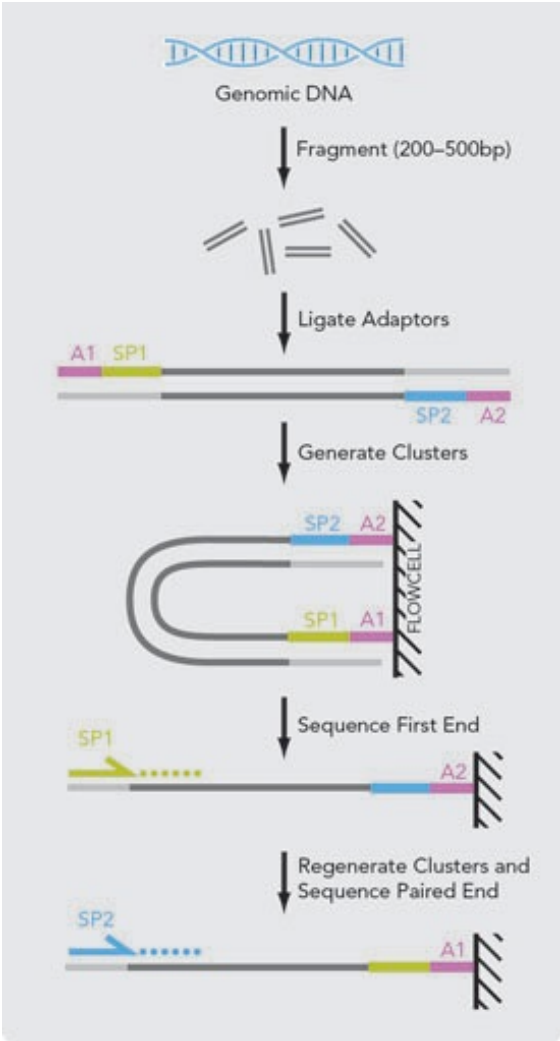
# Overview of RNA-Seq



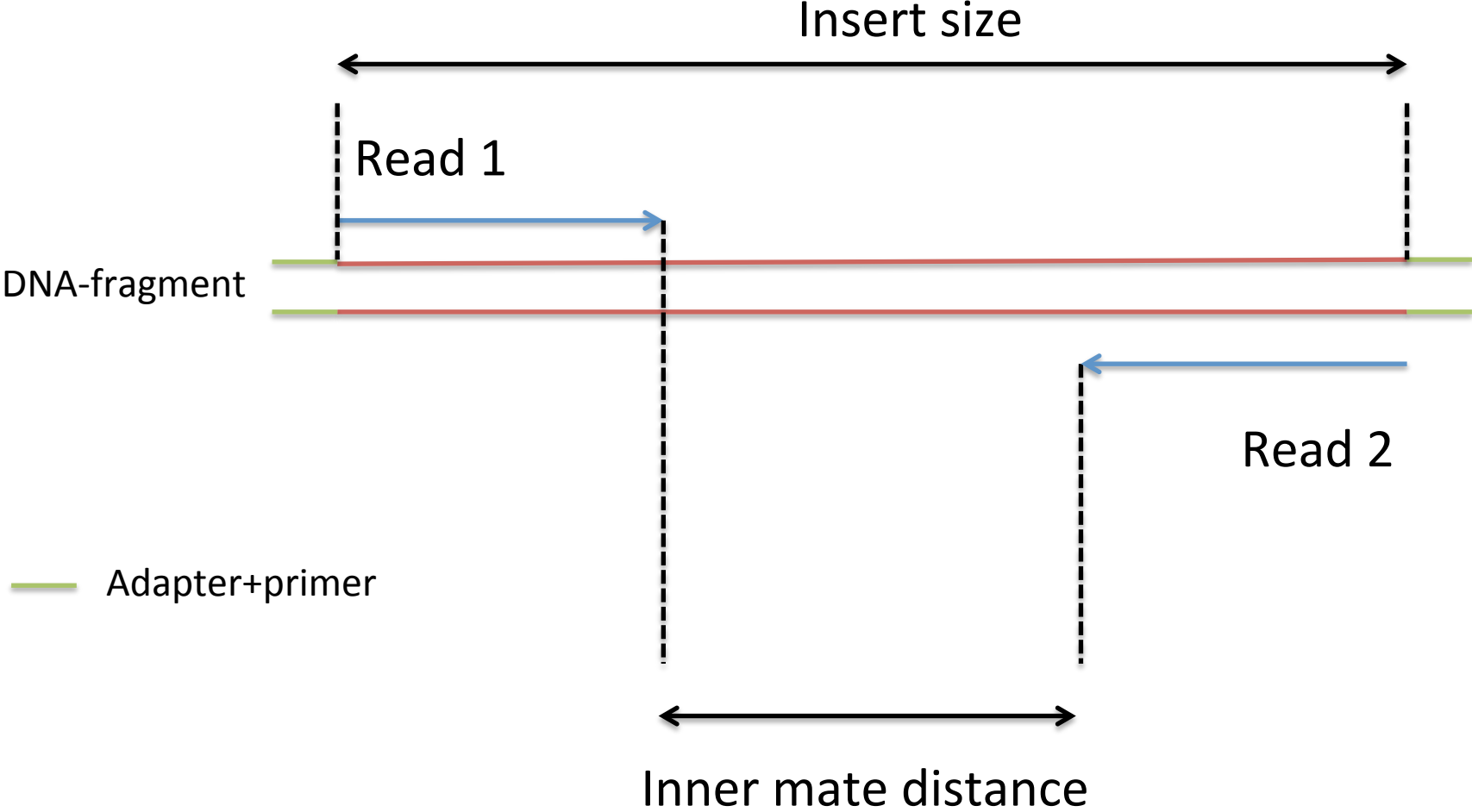
From: <http://www2.fml.tuebingen.mpg.de/raetsch/members/research/transcriptomics.html>



# Paired-End



# Insert size



# Paired-end gives you two files

FASTQ format (old):

```
@61DFRAAXX100204:1:100:10494:3070/1  
AAACAACAGGGCACATTGTCACTCTTGTATTTGAAAAACACTTTCCGGCCAT  
+  
ACCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCBC?CCCCCCCC@@CACCCCCA
```

```
@61DFRAAXX100204:1:100:10494:3070/2  
ATCCAAGTTAAAACAGAGGCCTGTGACAGACTCTTGGCCCATCGTGTTGATA  
+  
_^_a^cccegcgghhgZc`ghhc^egggd^_[d]defcdfd^Z^0XWaQ^ad
```

New: @<instrument>:<run number>:<flowcell ID>:<lane>:<tile>:<x-pos>:<y-pos>  
<read>:<is filtered>:<control number>:<sample number>

Example:

```
@SIM:1:FCX:1:15:6329:1045 1:N:0:2  
TCGCACTCAACGCCCTGCATATGACAAGACAGAATC  
+  
<>;##=><9=AAAAAAAAA9#:<#<;<<????#=#
```

# Transcript Reconstruction from RNA-Seq Reads



## Advancing RNA-Seq analysis

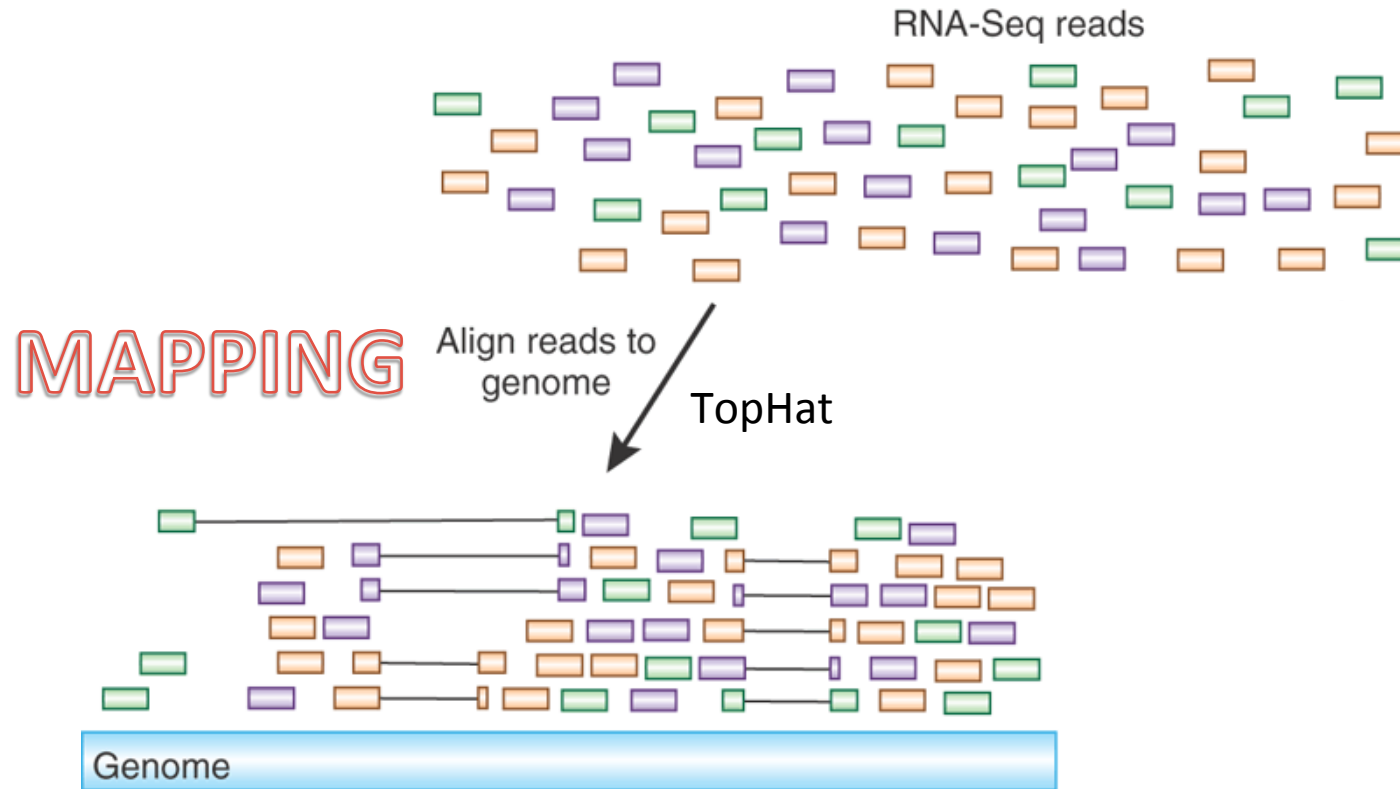
Brian J Haas & Michael C Zody

Nature Biotech, 2010

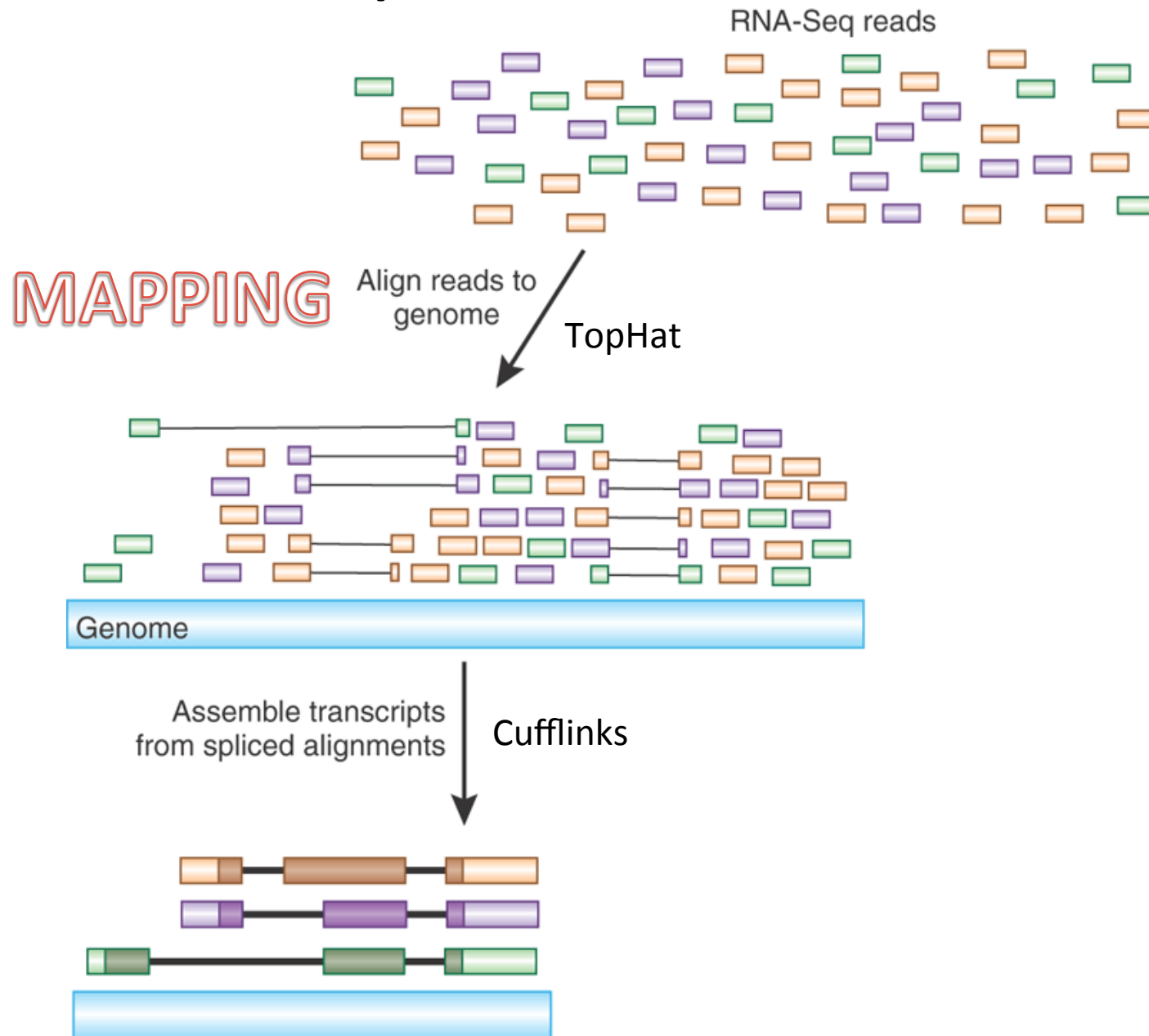
New methods for analyzing RNA-Seq data enable *de novo* reconstruction of the transcriptome.



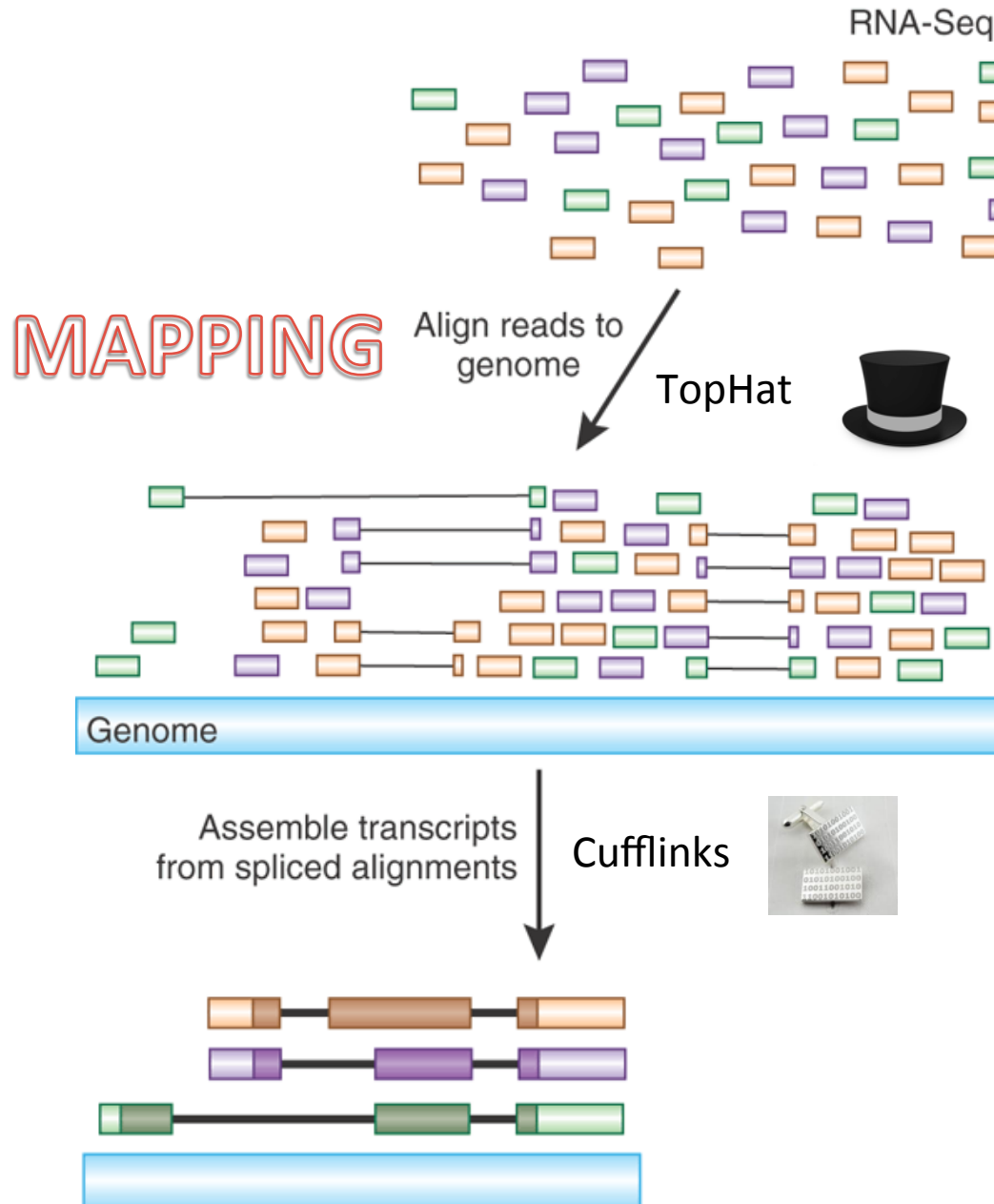
# Transcript Reconstruction from RNA-Seq Reads



# Transcript Reconstruction from RNA-Seq Reads



# Transcript Reconstruction from RNA-Seq Reads



NATURE PROTOCOLS | PROTOCOL

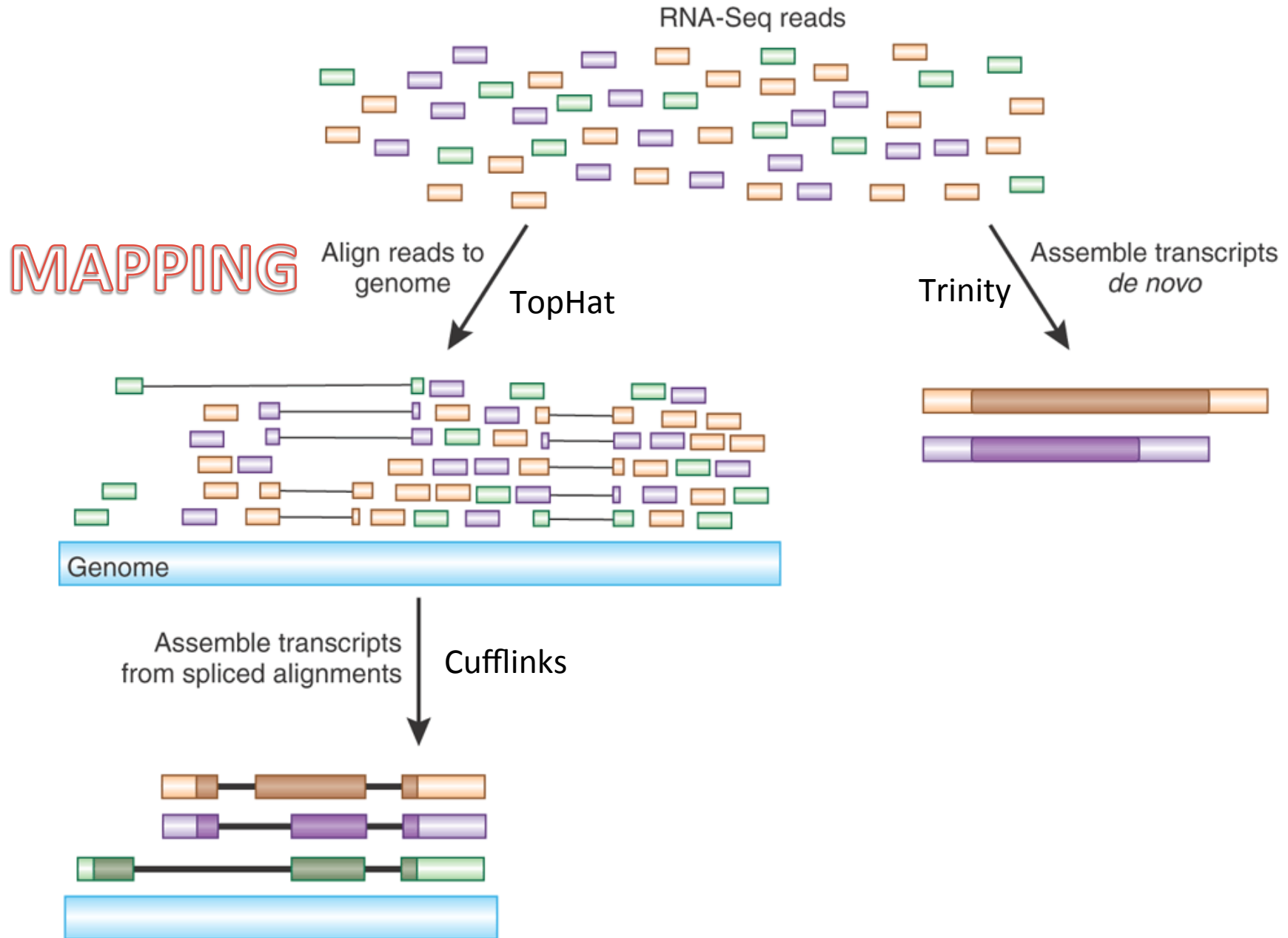
## Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks

Cole Trapnell, Adam Roberts, Loyal Goff, Geo Pertea, Daehwan Kim, David R Kelley, Harold Pimentel, Steven L Salzberg, John L Rinn & Lior Pachter

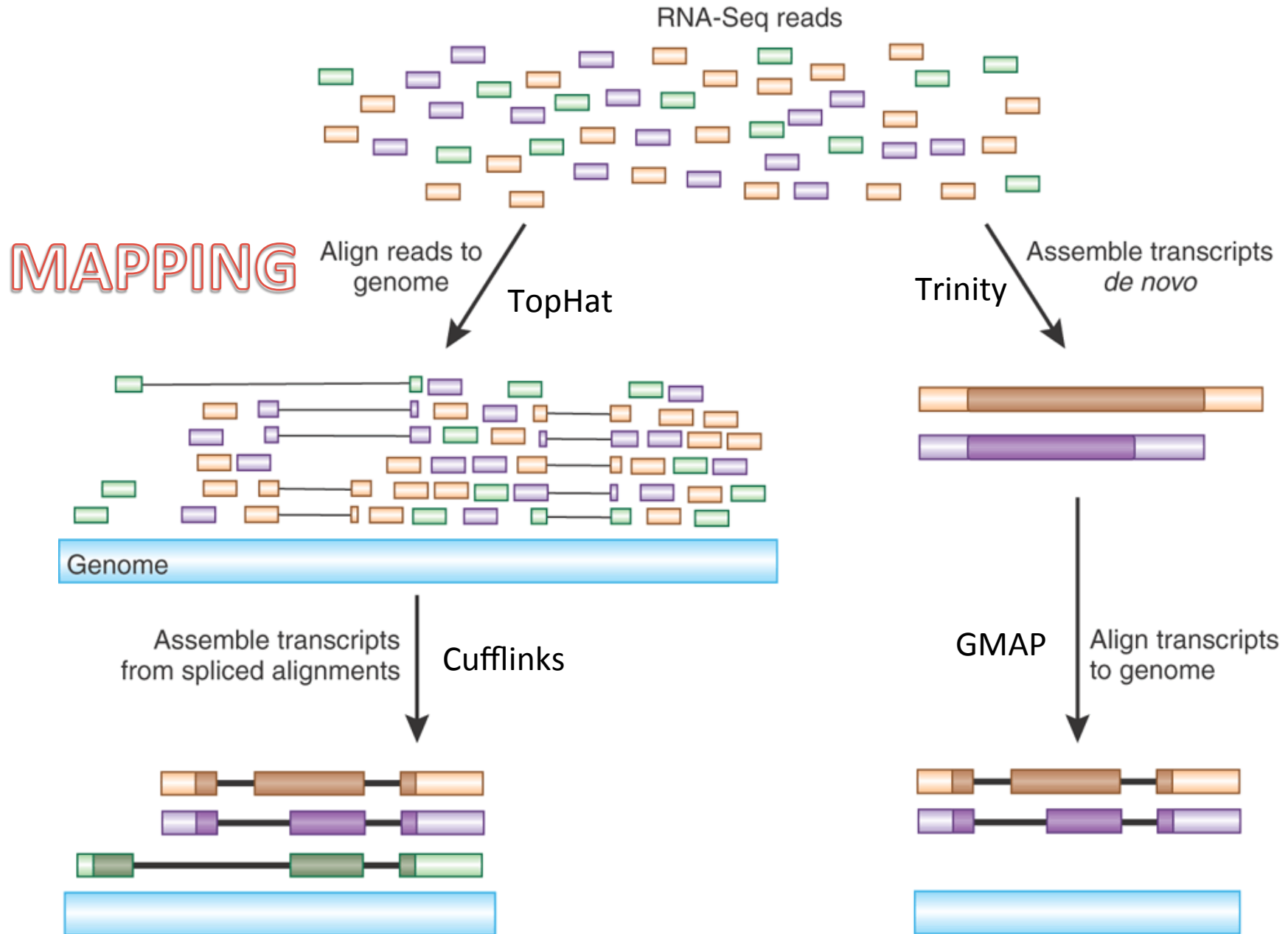
Affiliations | Contributions | Corresponding author

*Nature Protocols* 7, 562–578 (2012) | doi:10.1038/nprot.2012.016  
Published online 01 March 2012

# Transcript Reconstruction from RNA-Seq Reads



# Transcript Reconstruction from RNA-Seq Reads



# Transcript Reconstruction from RNA-Seq Reads



RNA-Seq reads

Assemble transcripts  
*de novo*



Trinity

Align transcripts  
to genome



## End-to-end **Transcriptome**-based RNA-Seq Analysis Software Package

NATURE PROTOCOLS | PROTOCOL

*De novo* transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis

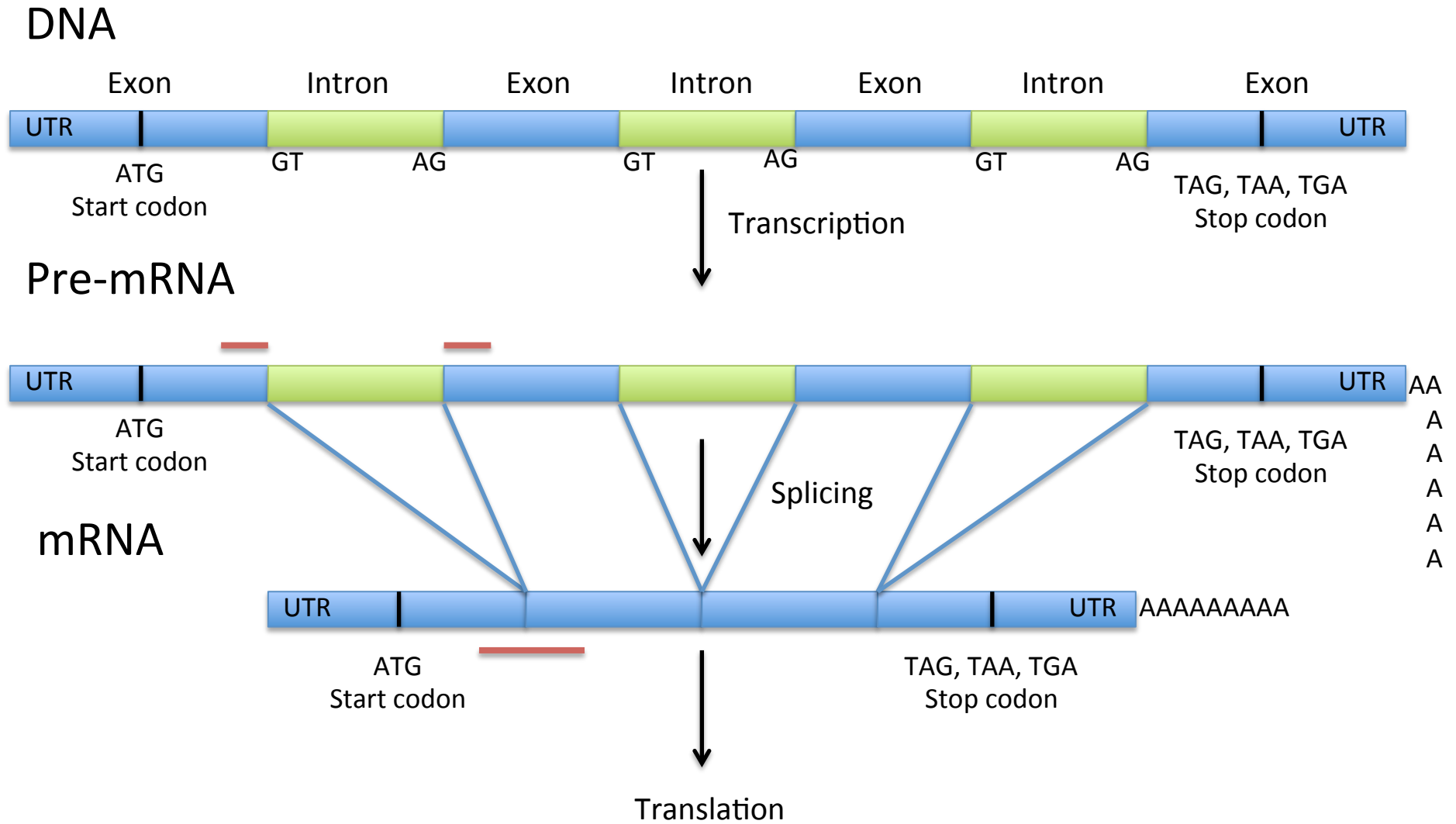
Brian J Haas, Alexie Papanicolaou, Moran Yassour, Manfred Grabherr, Philip D Blood, Joshua Bowden, Matthew Brian Couger, David Eccles, Bo Li, Matthias Lieber, Matthew D MacManes, Michael Ott, Joshua Orvis, Nathalie Pochet, Francesco Strozzi, Nathan Weeks, Rick Westerman, Thomas William, Colin N Dewey, Robert Henschel, Richard D LeDuc, Nir Friedman & Aviv Regev

[Affiliations](#) | [Contributions](#) | [Corresponding authors](#)

Nature Protocols **8**, 1494–1512 (2013) | doi:10.1038/nprot.2013.084

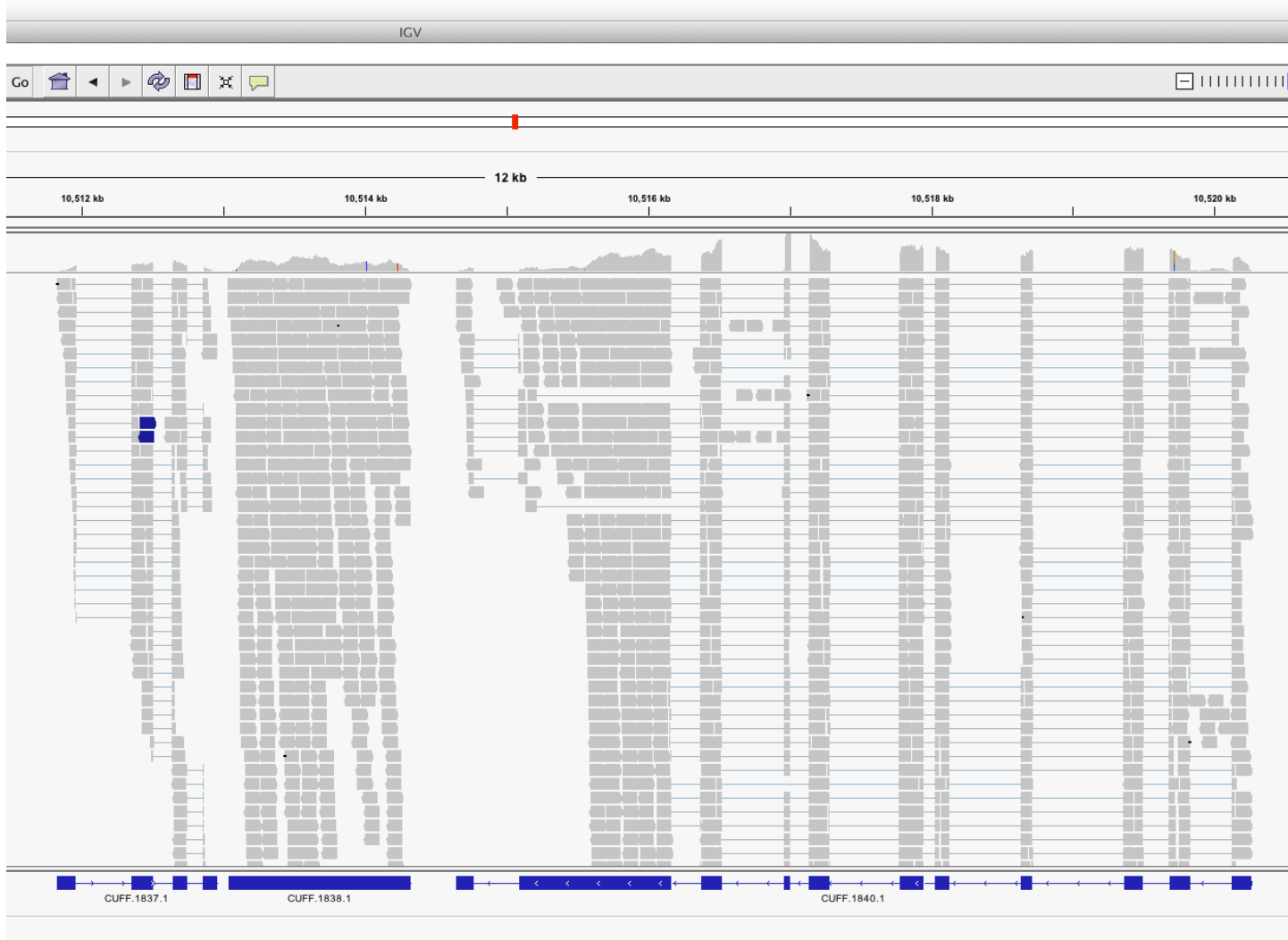
Published online 11 July 2013

# Basic concepts of mapping-based RNA-seq - Spliced reads



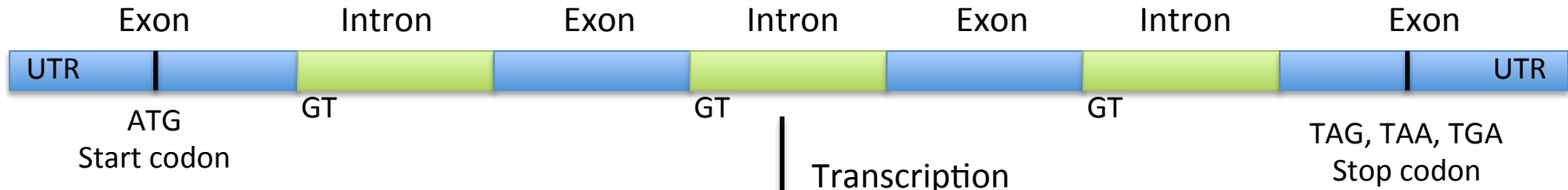


# RNA-seq - Spliced reads



# Pre-mRNA

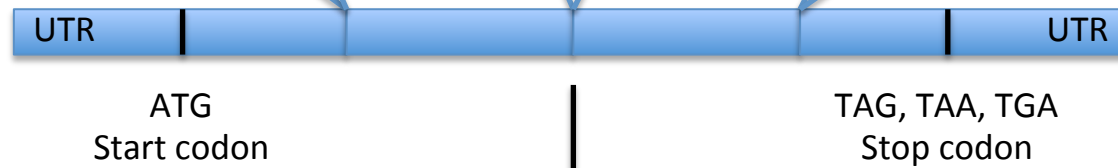
DNA



Pre-mRNA



mRNA

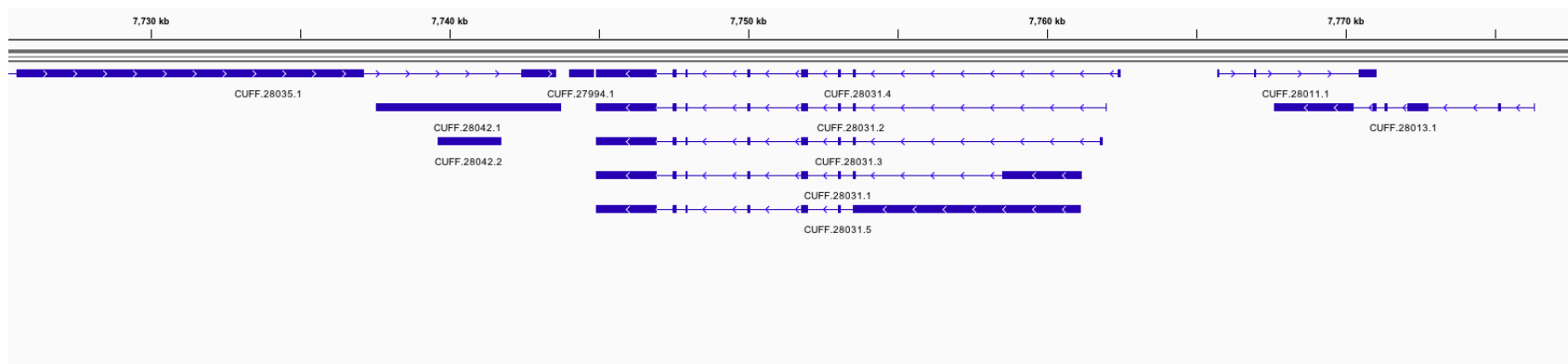


Translation

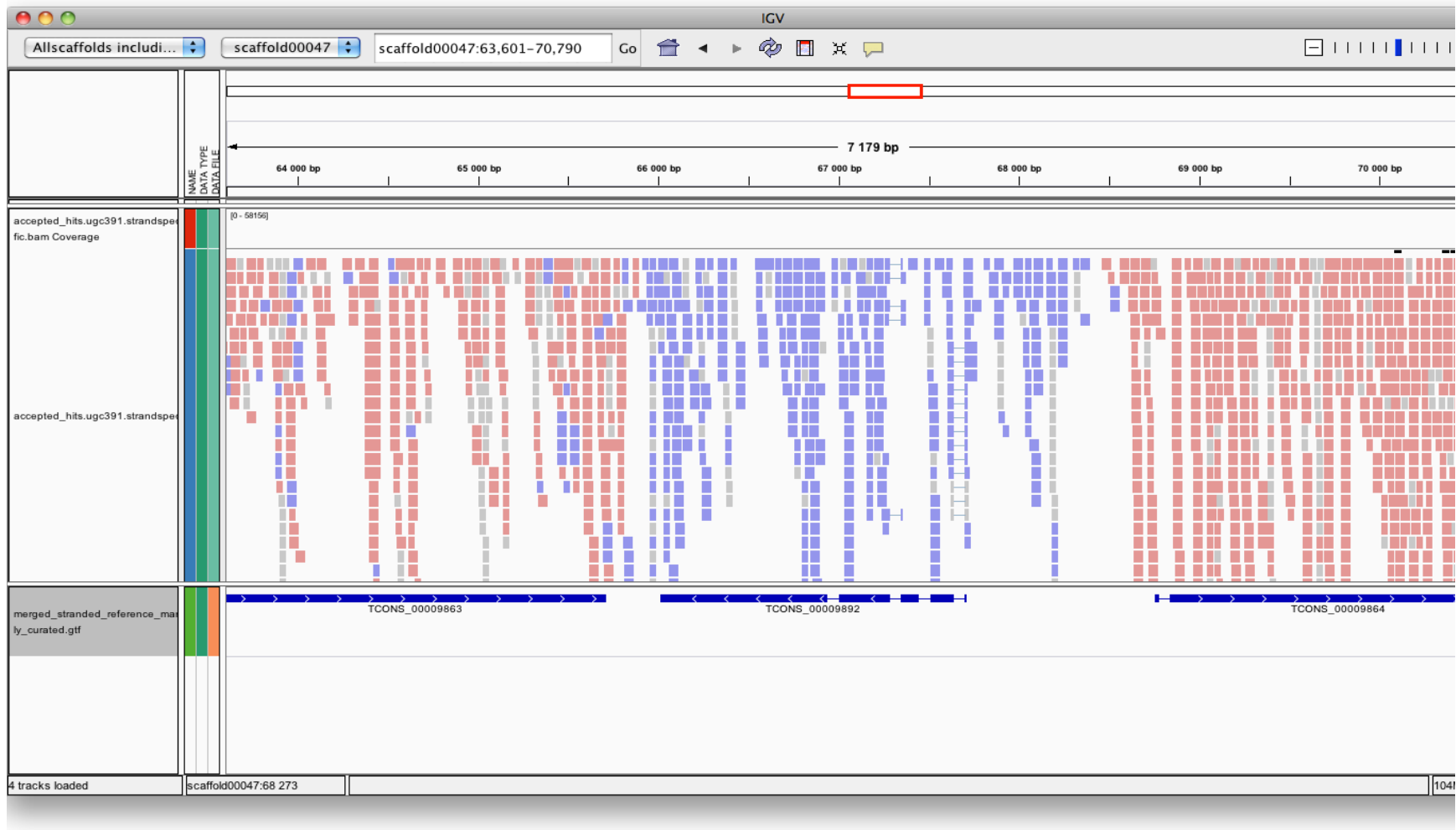
# Pre-mRNA



# Pre-mRNA



# Stranded rna-seq



# Overview of RNASeq

Bowtie (fast short-read alignment)

TopHat (spliced short-read alignment)



Cufflinks (transcript reconstruction from alignments)

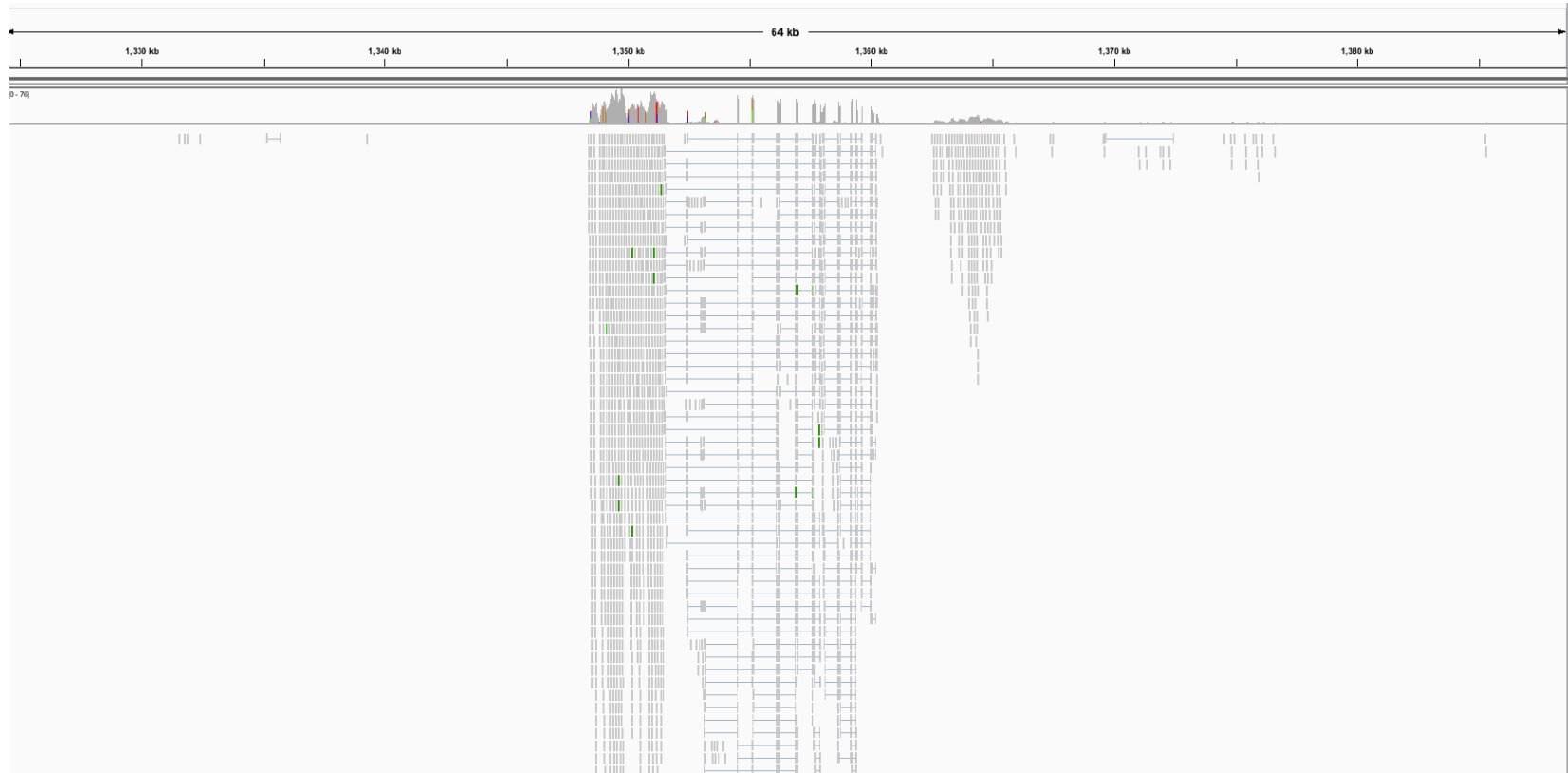


Cuffdiff (differential expression analysis)



CummeRbund (visualization & analysis)

# Tophat-mapped reads





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?=ADDBCBCDEEFFDEFFFDEFFGDBEFGEDGCFGFGGGGG
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
```

SAM format specification: <http://samtools.sourceforge.net/SAM1.pdf>

## 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      CCCAAACAAGCCGAAGCTGATTTGGCTCGTAAAGACCCGGAAA
10     ###CB?=ADDBCBCDEEFFDEFFFDEFFGDBEFGEDEGCFGFGGGGG
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
```

SAM format specification: <http://samtools.sourceforge.net/SAM1.pdf>

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 (metadata)
16      XQ:i:40
17      X2:i:0
```

SAM format specification: <http://samtools.sourceforge.net/SAM1.pdf>

# 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
```

# There is also CRAM...

• <b>CRAM compression rate</b>	<b>File format</b>	<b>File size (GB)</b>
• SAM		7.4
• BAM		1.9
• CRAM lossless		1.4
• CRAM 8 bins		0.8
• CRAM no quality scores		0.26

# Visualizing Alignments of RNA-Seq reads



# IGV

The image shows a screenshot of the Integrative Genomics Viewer (IGV) website homepage. The browser address bar shows the URL [www.broadinstitute.org/igv/](http://www.broadinstitute.org/igv/). The page features a navigation menu on the left, a main banner with the IGV logo and a visualization of genomic data, and several news and citation sections.

**Home**

**Integrative Genomics Viewer**

**What's New**

- July 3, 2012.** Soybean (*Glycine max*) and Rat (*m5*) 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.

**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 Thorvaldsdottir, James T. Robinson, Jill P. Mesirov. [Integrative Genomics Viewer \(IGV\): high-performance genomics data visualization and exploration](#).

**Overview**

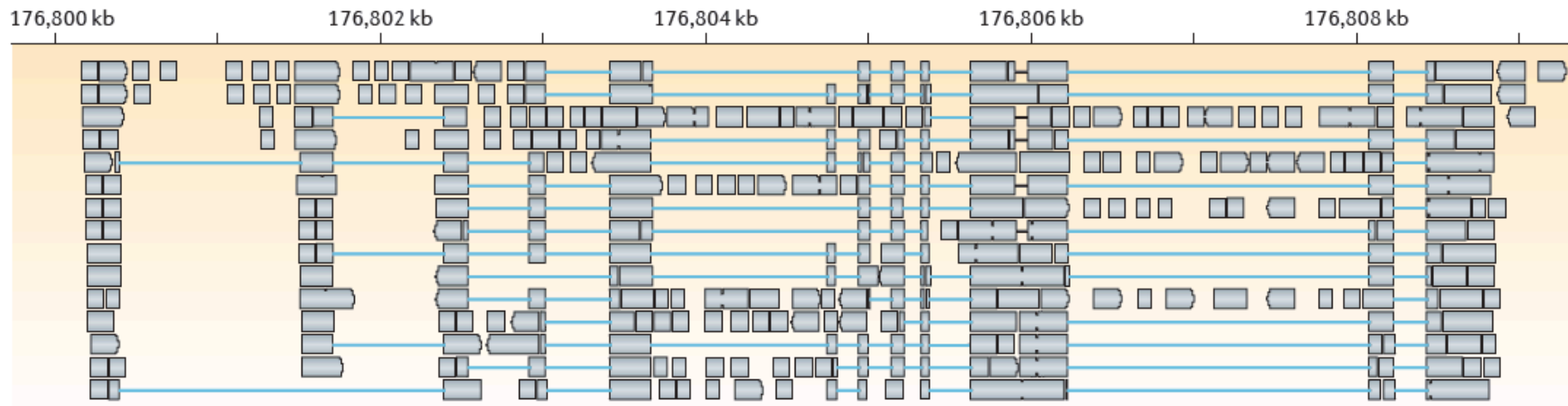


# IGV: Viewing Tophat Alignments



# Transcript Reconstruction Using Cufflinks

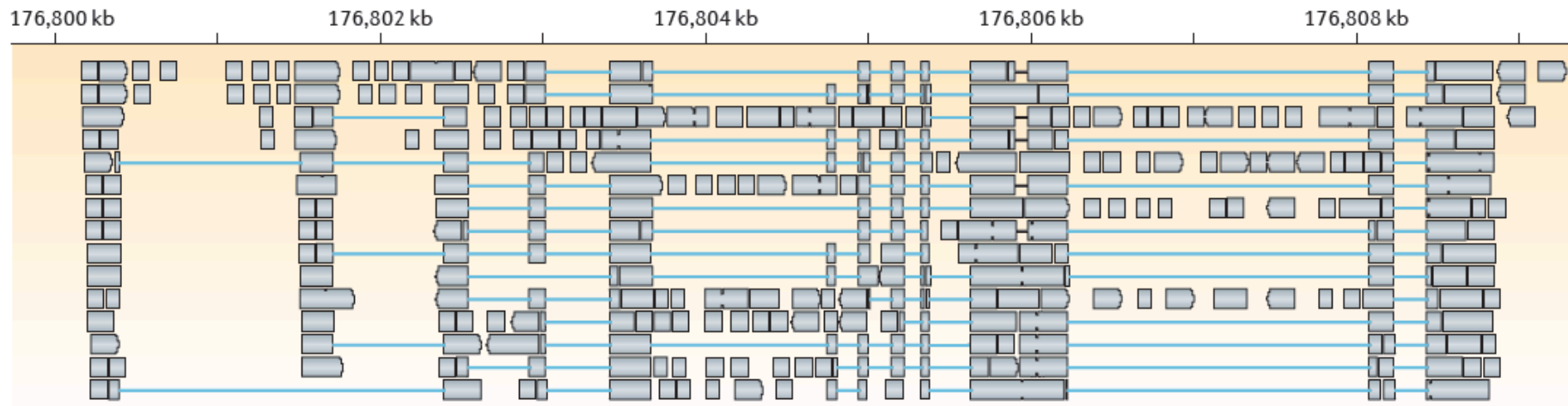
**a** Splice-align reads to the genome



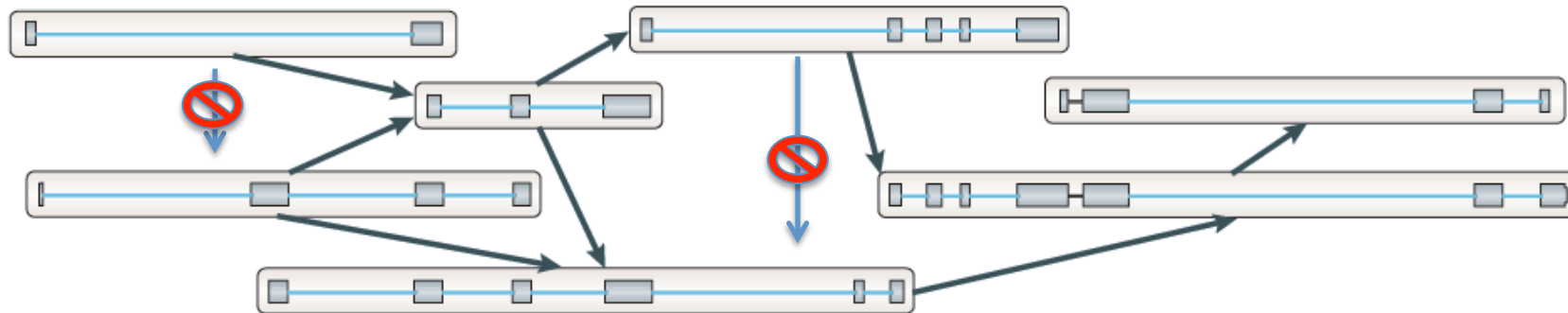
From Martin & Wang. Nature Reviews in Genetics. 2011

# Transcript Reconstruction Using Cufflinks

**a Splice-align reads to the genome**



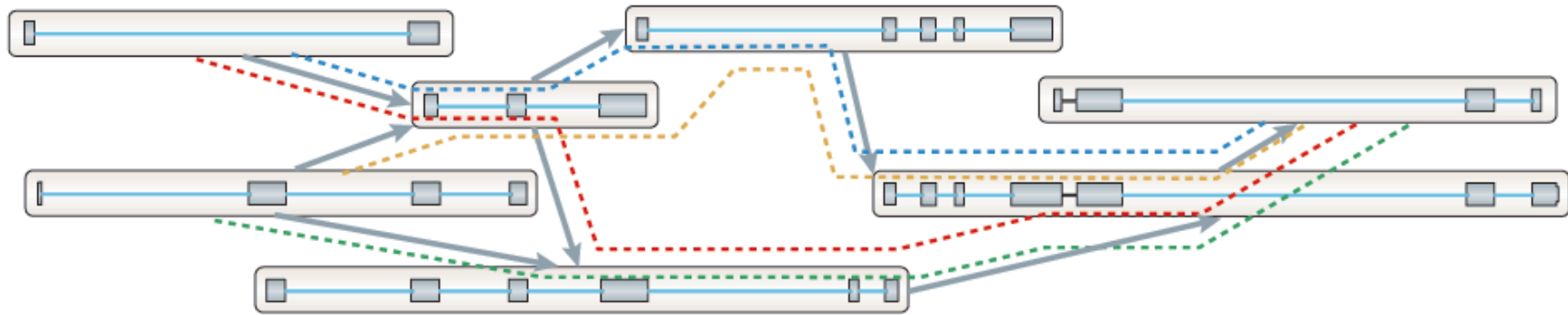
**b Build a graph representing alternative splicing events**



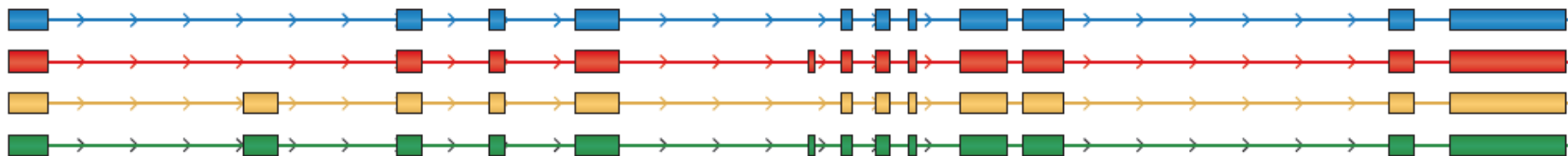
From Martin & Wang. Nature Reviews in Genetics. 2011

# Transcript Reconstruction Using Cufflinks

**c** Traverse the graph to assemble variants



**d** Assembled isoforms



# GFF file format

```
##gff-version 3
scaffold_7 maker gene 133848 144662 . - . ID=C5546228A39E2878A71E2ABDAFB34661;Name=maker-scaffold_7-augustus-gene-0.11
scaffold_7 maker mRNA 133848 144662 . - . ID=A649E923246BADE2184E579FA9124ABD;Parent=C5546228A39E2878A71E2ABDAFB34661;Name=1:cornix-all_reads.72406.1;_AED=1.00;_eAED=1.00;_OI=71|0|0|0|0|0|0|414|347
scaffold_7 maker exon 138974 139077 . - . ID=A649E923246BADE2184E579FA9124ABD;exon:7;Parent=A649E923246BADE2184E579FA9124ABD
scaffold_7 maker exon 135098 135281 . - . ID=A649E923246BADE2184E579FA9124ABD;exon:6;Parent=A649E923246BADE2184E579FA9124ABD
scaffold_7 maker exon 139616 139836 . - . ID=A649E923246BADE2184E579FA9124ABD;exon:5;Parent=A649E923246BADE2184E579FA9124ABD
scaffold_7 maker exon 144511 144662 . - . ID=A649E923246BADE2184E579FA9124ABD;exon:4;Parent=A649E923246BADE2184E579FA9124ABD
scaffold_7 maker exon 136342 136437 . - . ID=A649E923246BADE2184E579FA9124ABD;exon:3;Parent=A649E923246BADE2184E579FA9124ABD
scaffold_7 maker exon 133848 134338 . - . ID=A649E923246BADE2184E579FA9124ABD;exon:2;Parent=A649E923246BADE2184E579FA9124ABD
scaffold_7 maker exon 141262 141383 . - . ID=A649E923246BADE2184E579FA9124ABD;exon:1;Parent=A649E923246BADE2184E579FA9124ABD
scaffold_7 maker exon 144138 144296 . - . ID=A649E923246BADE2184E579FA9124ABD;exon:0;Parent=A649E923246BADE2184E579FA9124ABD
scaffold_7 maker five_prime_UTR 144592 144662 . - . ID=A649E923246BADE2184E579FA9124ABD;five_prime utr;Parent=A649E923246BADE2184E579FA9124ABD
scaffold_7 maker CDS 144511 144591 . - 0 ID=A649E923246BADE2184E579FA9124ABD;cds;Parent=A649E923246BADE2184E579FA9124ABD
scaffold_7 maker CDS 144138 144296 . - 0 ID=A649E923246BADE2184E579FA9124ABD;cds;Parent=A649E923246BADE2184E579FA9124ABD
scaffold_7 maker CDS 141262 141383 . - 0 ID=A649E923246BADE2184E579FA9124ABD;cds;Parent=A649E923246BADE2184E579FA9124ABD
scaffold_7 maker CDS 139616 139836 . - 1 ID=A649E923246BADE2184E579FA9124ABD;cds;Parent=A649E923246BADE2184E579FA9124ABD
scaffold_7 maker CDS 138974 139077 . - 2 ID=A649E923246BADE2184E579FA9124ABD;cds;Parent=A649E923246BADE2184E579FA9124ABD
scaffold_7 maker CDS 136342 136437 . - 0 ID=A649E923246BADE2184E579FA9124ABD;cds;Parent=A649E923246BADE2184E579FA9124ABD
scaffold_7 maker CDS 135098 135281 . - 0 ID=A649E923246BADE2184E579FA9124ABD;cds;Parent=A649E923246BADE2184E579FA9124ABD
scaffold_7 maker CDS 134262 134338 . - 2 ID=A649E923246BADE2184E579FA9124ABD;cds;Parent=A649E923246BADE2184E579FA9124ABD
scaffold_7 maker three_prime_UTR 133848 134261 . - . ID=A649E923246BADE2184E579FA9124ABD;three_prime utr;Parent=A649E923246BADE2184E579FA9124ABD
scaffold_7 maker gene 83101 117593 . + . ID=D3B9A5F27797F56A84A2E890FF6B99;Name=maker-scaffold_7-augustus-gene-0.6
scaffold_7 maker mRNA 83101 117593 . + . ID=CF5DDA190832937C45A002E674C9C26;Parent=D3B9A5F27797F56A84A2E890FF6B99;Name=maker-scaffold_7-augustus-gene-0.6-mRNA-1;_AED=1.00;_eAED=1.00;_OI=0|0|0|0|0|0|21|4|706
scaffold_7 maker exon 95748 95871 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:8;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 99113 99137 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:9;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 90664 90748 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:10;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 110231 110356 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:11;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 113609 113679 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:12;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 94057 94117 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:13;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 84578 84670 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:14;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 115452 115536 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:15;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 111579 111669 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:16;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 102917 103016 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:17;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 96766 96849 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:18;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 86666 86750 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:19;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 99944 100109 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:20;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 109766 109860 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:21;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 93154 93282 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:22;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 114737 114825 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:23;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 83101 83155 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:24;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 108533 108795 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:25;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 117477 117593 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:26;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 106779 106866 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:27;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker exon 105743 105835 . + . ID=CF5DDA190832937C45A002E674C9C26;exon:28;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 83101 83155 . + 0 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 84578 84670 . + 2 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 86666 86750 . + 2 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 90664 90748 . + 1 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 93154 93282 . + 0 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 94057 94117 . + 0 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 95748 95871 . + 2 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 96766 96849 . + 1 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 99113 99137 . + 1 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 99944 100109 . + 0 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 102917 103016 . + 2 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 105743 105835 . + 1 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 106779 106866 . + 1 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 108533 108795 . + 0 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 109766 109860 . + 1 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 110231 110356 . + 2 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 111579 111669 . + 2 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 113609 113679 . + 1 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 114737 114825 . + 2 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 115452 115536 . + 0 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker CDS 117477 117589 . + 2 ID=CF5DDA190832937C45A002E674C9C26;cds;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker three_prime_UTR 117590 117593 . + . ID=CF5DDA190832937C45A002E674C9C26;three_prime utr;Parent=CF5DDA190832937C45A002E674C9C26
scaffold_7 maker gene 22451 37514 . - . ID=E1E94A56071E5D7940C93BE70FD79E56;Name=E1E94A56071E5D7940C93BE70FD79E56
scaffold_7 maker mRNA 22451 37514 . - . ID=2B837928BFCDB81C7A4070E76500C5;Parent=E1E94A56071E5D7940C93BE70FD79E56;Name=scaffold_7.22450-37514;_AED=1.00;_eAED=1.00;_OI=3008|0|0|0|0|0|11|762|1414
scaffold_7 maker mRNA 26810 27085 . - . ID=10629810725FB97175585BA23623FAF6;Parent=E1E94A56071E5D7940C93BE70FD79E56;Name=maker-scaffold_7-augustus-gene-0.10-mRNA-1;_AED=1.00;_eAED=1.00;_OI=3008|0|0|0|0|0|10|0|0|10|0|0|0
scaffold_7 maker exon 27590 28090 . - . ID=2B837928BFCDB81C7A4070E76500C5;exon:39;Parent=2B837928BFCDB81C7A4070E76500C5;_AED=1.00;_eAED=1.00;_OI=3008|0|0|0|0|0|10|0|0|10|0|0|0
scaffold_7 maker exon 33513 33610 . - . ID=2B837928BFCDB81C7A4070E76500C5;exon:38;Parent=2B837928BFCDB81C7A4070E76500C5;_AED=1.00;_eAED=1.00;_OI=3008|0|0|0|0|0|10|0|0|10|0|0|0
maker.gff
```

# GFF3 file format

Seqid	source	type	start	end	score	strand	phase	attributes
Chr1	Snap	gene	234	3657	.	+	.	ID=gene1; Name=Snap1;
Chr1	Snap	mRNA	234	3657	.	+	.	ID=gene1.m1; Parent=gene1;
Chr1	Snap	exon	234	1543	.	+	.	ID=gene1.m1.exon1; Parent=gene1.m1;
Chr1	Snap	CDS	577	1543	.	+	0	ID=gene1.m1.CDS1; Parent=gene1.m1;
Chr1	Snap	exon	1822	2674	.	+	.	ID=gene1.m1.exon2; Parent=gene1.m1;
Chr1	Snap	CDS	1822	2674	.	+	2	ID=gene1.m1.CDS2; Parent=gene1.m1;
		start_ codon						Alias, note, ontology_term ...
		stop_ codon						

# GTF file format

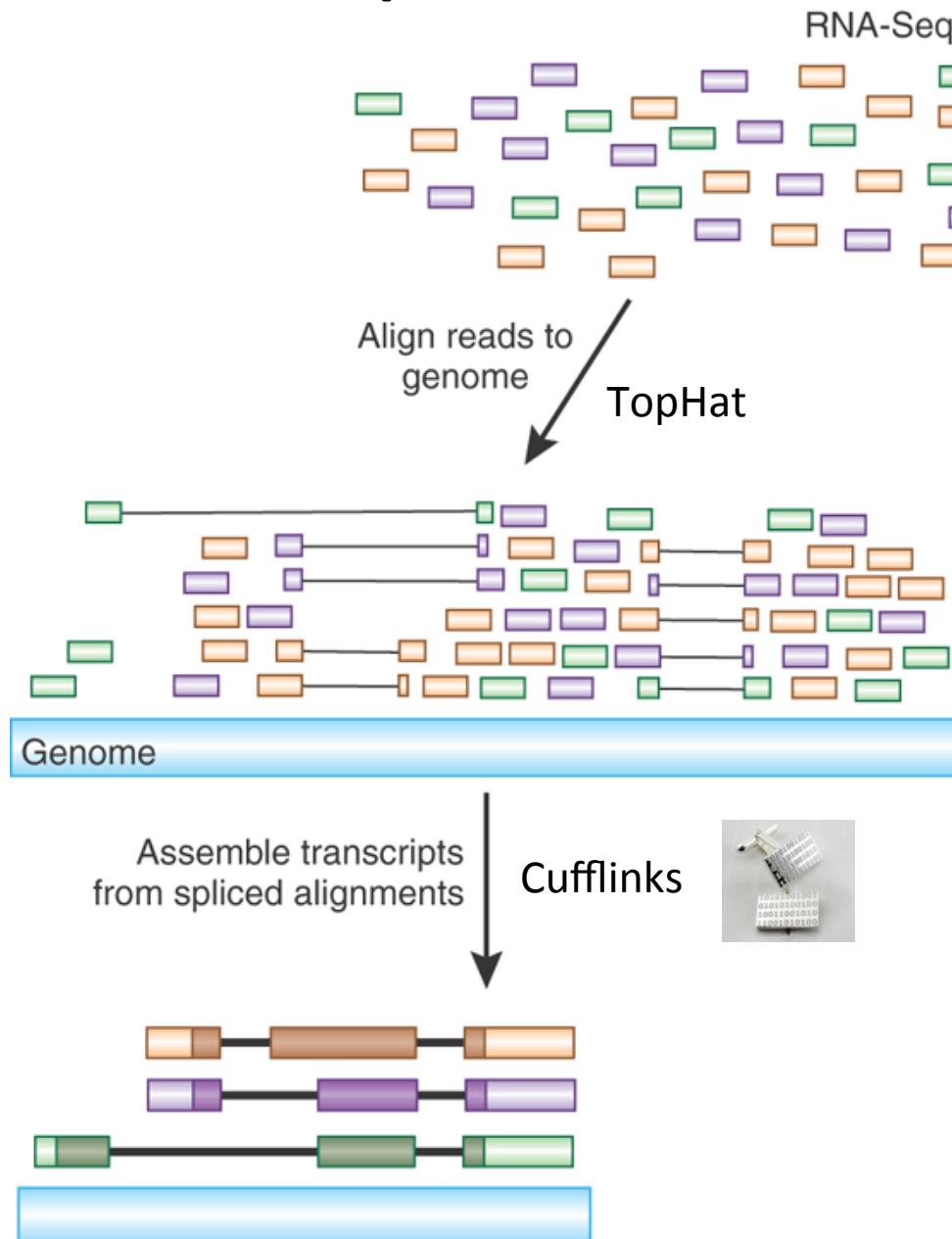
```
Sb_20131119_contig_1 Cufflinks transcript 1522 2095 1000 . . gene_id "CUFF.1"; transcript_id "CUFF.1.1"; FPKM "2.6064385494"; frac "1.000000"; conf_lo "0.948975"; conf_hi "3.440036"; cov "4.817376";
Sb_20131119_contig_1 Cufflinks exon 1522 2095 1000 . . gene_id "CUFF.1"; transcript_id "CUFF.1.1"; FPKM "2.6064385494"; frac "1.000000"; conf_lo "0.948975"; conf_hi "3.440036"; cov "4.817376";
Sb_20131119_contig_1 Cufflinks transcript 2626 4118 1000 . . gene_id "CUFF.2"; transcript_id "CUFF.2.1"; FPKM "3.1548106029"; frac "1.000000"; conf_lo "0.828669"; conf_hi "3.729011"; cov "8.517668";
Sb_20131119_contig_1 Cufflinks exon 2626 4118 1000 . . gene_id "CUFF.2"; transcript_id "CUFF.2.1"; FPKM "3.1548106029"; frac "1.000000"; conf_lo "0.828669"; conf_hi "3.729011"; cov "8.517668";
Sb_20131119_contig_1 Cufflinks transcript 4855 5340 1000 . . gene_id "CUFF.3"; transcript_id "CUFF.3.1"; FPKM "1.000000"; conf_lo "0.828669"; conf_hi "3.729011"; cov "8.517668";
Sb_20131119_contig_1 Cufflinks exon 4855 5340 1000 . . gene_id "CUFF.3"; transcript_id "CUFF.3.1"; FPKM "1.000000"; conf_lo "0.828669"; conf_hi "3.729011"; cov "8.517668";
Sb_20131119_contig_1 Cufflinks transcript 5398 5975 1000 . . gene_id "CUFF.4"; transcript_id "CUFF.4.1"; FPKM "3.1706609980"; frac "1.000000"; conf_lo "1.178010"; conf_hi "6.164435"; cov "16.171080";
Sb_20131119_contig_1 Cufflinks exon 5398 5975 1000 . . gene_id "CUFF.4"; transcript_id "CUFF.4.1"; FPKM "3.1706609980"; frac "1.000000"; conf_lo "1.178010"; conf_hi "6.164435"; cov "16.171080";
Sb_20131119_contig_10 Cufflinks transcript 954 2795 1000 . . gene_id "CUFF.5"; transcript_id "CUFF.5.1"; FPKM "7.0235237898"; frac "1.000000"; conf_lo "2.521814"; conf_hi "6.164435"; cov "16.171080";
Sb_20131119_contig_10 Cufflinks exon 954 2795 1000 . . gene_id "CUFF.5"; transcript_id "CUFF.5.1"; FPKM "7.0235237898"; frac "1.000000"; conf_lo "2.521814"; conf_hi "6.164435"; cov "16.171080";
Sb_20131119_contig_1 Cufflinks transcript 4502 4718 1000 . . gene_id "CUFF.6"; transcript_id "CUFF.6.1"; FPKM "3.1706609980"; frac "1.000000"; conf_lo "1.178010"; conf_hi "6.164435"; cov "16.171080";
Sb_20131119_contig_1 Cufflinks exon 4502 4718 1000 . . gene_id "CUFF.6"; transcript_id "CUFF.6.1"; FPKM "3.1706609980"; frac "1.000000"; conf_lo "1.178010"; conf_hi "6.164435"; cov "16.171080";
Sb_20131119_contig_1 Cufflinks transcript 10522 13208 1000 . . gene_id "CUFF.7"; transcript_id "CUFF.7.1"; FPKM "41.2374406123"; frac "1.000000"; conf_lo "31.982715"; conf_hi "39.274371"; cov "89.788421";
Sb_20131119_contig_1 Cufflinks exon 10522 13208 1000 . . gene_id "CUFF.7"; transcript_id "CUFF.7.1"; FPKM "41.2374406123"; frac "1.000000"; conf_lo "31.982715"; conf_hi "39.274371"; cov "89.788421";
Sb_20131119_contig_1 Cufflinks transcript 14623 14623 1000 . . gene_id "CUFF.8"; transcript_id "CUFF.8.1"; FPKM "41.2374406123"; frac "1.000000"; conf_lo "31.982715"; conf_hi "39.274371"; cov "89.788421";
Sb_20131119_contig_100022 Cufflinks exon 14623 14623 1000 . . gene_id "CUFF.8"; transcript_id "CUFF.8.1"; FPKM "41.2374406123"; frac "1.000000"; conf_lo "31.982715"; conf_hi "39.274371"; cov "89.788421";
Sb_20131119_contig_100022 Cufflinks transcript 991 4547 1000 . . gene_id "CUFF.9"; transcript_id "CUFF.9.1"; FPKM "55.6377793473"; frac "1.000000"; conf_lo "48.931832"; conf_hi "55.241530"; cov "121.429110";
Sb_20131119_contig_100023 Cufflinks exon 991 4547 1000 . . gene_id "CUFF.9"; transcript_id "CUFF.9.1"; FPKM "55.6377793473"; frac "1.000000"; conf_lo "48.931832"; conf_hi "55.241530"; cov "121.429110";
Sb_20131119_contig_100023 Cufflinks transcript 1097 2009 1000 . . gene_id "CUFF.10"; transcript_id "CUFF.10.1"; FPKM "18.4360530685"; frac "1.000000"; conf_lo "2.807793"; conf_hi "8.657363"; cov "24.668034";
Sb_20131119_contig_100023 Cufflinks exon 1097 2009 1000 . . gene_id "CUFF.10"; transcript_id "CUFF.10.1"; FPKM "18.4360530685"; frac "1.000000"; conf_lo "2.807793"; conf_hi "8.657363"; cov "24.668034";
Sb_20131119_contig_100023 Cufflinks transcript 1 123 1000 . . gene_id "CUFF.11"; transcript_id "CUFF.11.1"; FPKM "29.5909083872"; frac "1.000000"; conf_lo "3.216803"; conf_hi "9.918474"; cov "41.309085";
Sb_20131119_contig_100023 Cufflinks exon 1 123 1000 . . gene_id "CUFF.11"; transcript_id "CUFF.11.1"; FPKM "29.5909083872"; frac "1.000000"; conf_lo "3.216803"; conf_hi "9.918474"; cov "41.309085";
Sb_20131119_contig_100040 Cufflinks transcript 1 221 1000 . . gene_id "CUFF.12"; transcript_id "CUFF.12.1"; FPKM "18.4360530685"; frac "1.000000"; conf_lo "2.807793"; conf_hi "8.657363"; cov "24.668034";
Sb_20131119_contig_100040 Cufflinks exon 1 221 1000 . . gene_id "CUFF.12"; transcript_id "CUFF.12.1"; FPKM "18.4360530685"; frac "1.000000"; conf_lo "2.807793"; conf_hi "8.657363"; cov "24.668034";
Sb_20131119_contig_100040 Cufflinks transcript 2 255 1000 . . gene_id "CUFF.13"; transcript_id "CUFF.13.1"; FPKM "25.0500552879"; frac "1.000000"; conf_lo "10.508448"; conf_hi "18.137870"; cov "45.523657";
Sb_20131119_contig_100040 Cufflinks exon 2 255 1000 . . gene_id "CUFF.13"; transcript_id "CUFF.13.1"; FPKM "25.0500552879"; frac "1.000000"; conf_lo "10.508448"; conf_hi "18.137870"; cov "45.523657";
Sb_20131119_contig_100107 Cufflinks transcript 2041 2331 1000 . . gene_id "CUFF.14"; transcript_id "CUFF.14.1"; FPKM "3960.0565774823"; frac "1.000000"; conf_lo "5.622026"; conf_hi "18.740088"; cov "4272.822539";
Sb_20131119_contig_100107 Cufflinks exon 2041 2331 1000 . . gene_id "CUFF.14"; transcript_id "CUFF.14.1"; FPKM "3960.0565774823"; frac "1.000000"; conf_lo "5.622026"; conf_hi "18.740088"; cov "4272.822539";
Sb_20131119_contig_100111 Cufflinks transcript 21 129 1000 . . gene_id "CUFF.15"; transcript_id "CUFF.15.1"; FPKM "418.4602721559"; frac "1.000000"; conf_lo "7.442428"; conf_hi "11.194395"; cov "24.833916";
Sb_20131119_contig_100111 Cufflinks exon 21 129 1000 . . gene_id "CUFF.15"; transcript_id "CUFF.15.1"; FPKM "418.4602721559"; frac "1.000000"; conf_lo "7.442428"; conf_hi "11.194395"; cov "24.833916";
Sb_20131119_contig_100121 Cufflinks transcript 1756 2236 1000 . . gene_id "CUFF.16"; transcript_id "CUFF.16.1"; FPKM "1080.5400540118"; frac "1.000000"; conf_lo "23.150255"; conf_hi "44.484805"; cov "1391.311243";
Sb_20131119_contig_100121 Cufflinks exon 1756 2236 1000 . . gene_id "CUFF.16"; transcript_id "CUFF.16.1"; FPKM "1080.5400540118"; frac "1.000000"; conf_lo "23.150255"; conf_hi "44.484805"; cov "1391.311243";
Sb_20131119_contig_100192 Cufflinks transcript 1840 2212 1000 . . gene_id "CUFF.17"; transcript_id "CUFF.17.1"; FPKM "24.9098132799"; frac "1.000000"; conf_lo "7.484312"; conf_hi "14.786080"; cov "49.916331";
Sb_20131119_contig_100192 Cufflinks exon 1840 2212 1000 . . gene_id "CUFF.17"; transcript_id "CUFF.17.1"; FPKM "24.9098132799"; frac "1.000000"; conf_lo "7.484312"; conf_hi "14.786080"; cov "49.916331";
Sb_20131119_contig_100107 Cufflinks transcript 430 902 1000 . . gene_id "CUFF.18"; transcript_id "CUFF.18.1"; FPKM "69.4702369875"; frac "1.000000"; conf_lo "60.275496"; conf_hi "70.712087"; cov "143.103653";
Sb_20131119_contig_100107 Cufflinks exon 430 902 1000 . . gene_id "CUFF.18"; transcript_id "CUFF.18.1"; FPKM "69.4702369875"; frac "1.000000"; conf_lo "60.275496"; conf_hi "70.712087"; cov "143.103653";
Sb_20131119_contig_100192 Cufflinks transcript 1 616 1000 . . gene_id "CUFF.19"; transcript_id "CUFF.19.1"; FPKM "42.2951435419"; frac "1.000000"; conf_lo "23.433223"; conf_hi "33.823425"; cov "88.806988";
Sb_20131119_contig_100192 Cufflinks exon 1 616 1000 . . gene_id "CUFF.19"; transcript_id "CUFF.19.1"; FPKM "42.2951435419"; frac "1.000000"; conf_lo "23.433223"; conf_hi "33.823425"; cov "88.806988";
Sb_20131119_contig_100111 Cufflinks transcript 219 353 1000 . . gene_id "CUFF.20"; transcript_id "CUFF.20.1"; FPKM "24.9098132799"; frac "1.000000"; conf_lo "7.484312"; conf_hi "14.786080"; cov "49.916331";
Sb_20131119_contig_100111 Cufflinks exon 219 353 1000 . . gene_id "CUFF.20"; transcript_id "CUFF.20.1"; FPKM "24.9098132799"; frac "1.000000"; conf_lo "7.484312"; conf_hi "14.786080"; cov "49.916331";
Sb_20131119_contig_100040 Cufflinks transcript 945 2276 1000 . . gene_id "CUFF.21"; transcript_id "CUFF.21.1"; FPKM "13.6659124939"; frac "1.000000"; conf_lo "11.688997"; conf_hi "14.524101"; cov "32.588772";
Sb_20131119_contig_100040 Cufflinks exon 945 2276 1000 . . gene_id "CUFF.21"; transcript_id "CUFF.21.1"; FPKM "13.6659124939"; frac "1.000000"; conf_lo "11.688997"; conf_hi "14.524101"; cov "32.588772";
Sb_20131119_contig_100040 Cufflinks transcript 2079 2276 1000 . . gene_id "CUFF.22"; transcript_id "CUFF.22.1"; FPKM "1080.5400540118"; frac "1.000000"; conf_lo "23.150255"; conf_hi "44.484805"; cov "1391.311243";
Sb_20131119_contig_100040 Cufflinks exon 2079 2276 1000 . . gene_id "CUFF.22"; transcript_id "CUFF.22.1"; FPKM "1080.5400540118"; frac "1.000000"; conf_lo "23.150255"; conf_hi "44.484805"; cov "1391.311243";
Sb_20131119_contig_100121 Cufflinks transcript 1 150 1000 . . gene_id "CUFF.23"; transcript_id "CUFF.23.1"; FPKM "41.2374406123"; frac "1.000000"; conf_lo "31.982715"; conf_hi "39.274371"; cov "89.788421";
Sb_20131119_contig_100121 Cufflinks exon 1 150 1000 . . gene_id "CUFF.23"; transcript_id "CUFF.23.1"; FPKM "41.2374406123"; frac "1.000000"; conf_lo "31.982715"; conf_hi "39.274371"; cov "89.788421";
Sb_20131119_contig_100022 Cufflinks transcript 1510 2616 1000 . . gene_id "CUFF.24"; transcript_id "CUFF.24.1"; FPKM "13.7863265713"; frac "1.000000"; conf_lo "12.049946"; conf_hi "14.805328"; cov "32.698360";
Sb_20131119_contig_100022 Cufflinks exon 1510 2616 1000 . . gene_id "CUFF.24"; transcript_id "CUFF.24.1"; FPKM "13.7863265713"; frac "1.000000"; conf_lo "12.049946"; conf_hi "14.805328"; cov "32.698360";
Sb_20131119_contig_100112 Cufflinks transcript 933 3862 1000 . . gene_id "CUFF.25"; transcript_id "CUFF.25.1"; FPKM "43.3120795956"; frac "1.000000"; conf_lo "12.379816"; conf_hi "22.244981"; cov "78.396470";
Sb_20131119_contig_100112 Cufflinks exon 933 3862 1000 . . gene_id "CUFF.25"; transcript_id "CUFF.25.1"; FPKM "43.3120795956"; frac "1.000000"; conf_lo "12.379816"; conf_hi "22.244981"; cov "78.396470";
Sb_20131119_contig_10022 Cufflinks transcript 2857 6365 1000 . . gene_id "CUFF.26"; transcript_id "CUFF.26.1"; FPKM "5.9167625988"; frac "1.000000"; conf_lo "4.912969"; conf_hi "6.752149"; cov "12.795227";
Sb_20131119_contig_10022 Cufflinks exon 2857 6365 1000 . . gene_id "CUFF.26"; transcript_id "CUFF.26.1"; FPKM "5.9167625988"; frac "1.000000"; conf_lo "4.912969"; conf_hi "6.752149"; cov "12.795227";
Sb_20131119_contig_100022 Cufflinks transcript 39 1611 1000 . . gene_id "CUFF.27"; transcript_id "CUFF.27.1"; FPKM "49.5933867311"; frac "1.000000"; conf_lo "35.335602"; conf_hi "43.965643"; cov "84.354408";
Sb_20131119_contig_100022 Cufflinks exon 39 1611 1000 . . gene_id "CUFF.27"; transcript_id "CUFF.27.1"; FPKM "49.5933867311"; frac "1.000000"; conf_lo "35.335602"; conf_hi "43.965643"; cov "84.354408";
Sb_20131119_contig_10026 Cufflinks transcript 26 377 1000 . . gene_id "CUFF.28"; transcript_id "CUFF.28.1"; FPKM "44.1667134200"; frac "1.000000"; conf_lo "34.717689"; conf_hi "42.267104"; cov "77.065569";
Sb_20131119_contig_10026 Cufflinks exon 26 377 1000 . . gene_id "CUFF.28"; transcript_id "CUFF.28.1"; FPKM "44.1667134200"; frac "1.000000"; conf_lo "34.717689"; conf_hi "42.267104"; cov "77.065569";
Sb_20131119_contig_10031 Cufflinks transcript 377 1002 1000 . . gene_id "CUFF.29"; transcript_id "CUFF.29.1"; FPKM "386.7080614787"; frac "1.000000"; conf_lo "33.884660"; conf_hi "52.744990"; cov "584.846045";
Sb_20131119_contig_10031 Cufflinks exon 377 1002 1000 . . gene_id "CUFF.29"; transcript_id "CUFF.29.1"; FPKM "386.7080614787"; frac "1.000000"; conf_lo "33.884660"; conf_hi "52.744990"; cov "584.846045";
Sb_20131119_contig_10032 Cufflinks transcript 5174 5386 1000 . . gene_id "CUFF.30"; transcript_id "CUFF.30.1"; FPKM "90.3398746415"; frac "1.000000"; conf_lo "9.098097"; conf_hi "19.663630"; cov "118.213868";
Sb_20131119_contig_10032 Cufflinks exon 5174 5386 1000 . . gene_id "CUFF.30"; transcript_id "CUFF.30.1"; FPKM "90.3398746415"; frac "1.000000"; conf_lo "9.098097"; conf_hi "19.663630"; cov "118.213868";
Sb_20131119_contig_10032 Cufflinks transcript 160 3305 1000 . . gene_id "CUFF.31"; transcript_id "CUFF.31.1"; FPKM "5.2492236693"; frac "1.000000"; conf_lo "3.027818"; conf_hi "5.981787"; cov "12.544434";
Sb_20131119_contig_10032 Cufflinks exon 160 3305 1000 . . gene_id "CUFF.31"; transcript_id "CUFF.31.1"; FPKM "5.2492236693"; frac "1.000000"; conf_lo "3.027818"; conf_hi "5.981787"; cov "12.544434";
Sb_20131119_contig_100403 Cufflinks transcript 822 1053 1000 . . gene_id "CUFF.32"; transcript_id "CUFF.32.1"; FPKM "5.4971611200"; frac "1.000000"; conf_lo "4.259355"; conf_hi "6.262152"; cov "11.356891";
Sb_20131119_contig_100403 Cufflinks exon 822 1053 1000 . . gene_id "CUFF.32"; transcript_id "CUFF.32.1"; FPKM "5.4971611200"; frac "1.000000"; conf_lo "4.259355"; conf_hi "6.262152"; cov "11.356891";
Sb_20131119_contig_100403 Cufflinks transcript 822 1053 1000 . . gene_id "CUFF.33"; transcript_id "CUFF.33.1"; FPKM "55.9644866986"; frac "1.000000"; conf_lo "45.789638"; conf_hi "53.715558"; cov "119.323264";
Sb_20131119_contig_100403 Cufflinks exon 822 1053 1000 . . gene_id "CUFF.33"; transcript_id "CUFF.33.1"; FPKM "55.9644866986"; frac "1.000000"; conf_lo "45.789638"; conf_hi "53.715558"; cov "119.323264";
Sb_20131119_contig_100403 Cufflinks transcript 2827 3748 1000 . . gene_id "CUFF.34"; transcript_id "CUFF.34.1"; FPKM "5.2492236693"; frac "1.000000"; conf_lo "3.027818"; conf_hi "5.981787"; cov "12.544434";
Sb_20131119_contig_100403 Cufflinks exon 2827 3748 1000 . . gene_id "CUFF.34"; transcript_id "CUFF.34.1"; FPKM "5.2492236693"; frac "1.000000"; conf_lo "3.027818"; conf_hi "5.981787"; cov "12.544434";
Sb_20131119_contig_100328 Cufflinks transcript 1505 3535 1000 . . gene_id "CUFF.35"; transcript_id "CUFF.35.1"; FPKM "4.5481695975"; frac "1.000000"; conf_lo "3.453060"; conf_hi "5.296928"; cov "10.333563";
```

# GTF file format

Seqid	source	type	start	end	score	strand	phase	attributes
Chr1	Snap	exon	234	1543	.	+	.	gene_id "gene1"; transcript_id "transcript1";
Chr1	Snap	CDS	577	1543	.	+	0	gene_id "gene1"; transcript_id "transcript1";
Chr1	Snap	exon	1822	2674	.	+	.	gene_id "gene1"; transcript_id "transcript1";
Chr1	Snap	CDS	1822	2674	.	+	2	gene_id "gene1"; transcript_id "transcript1";
		start_codon						
		stop_codon						



# Transcript Reconstruction from RNA-Seq Reads



NATURE PROTOCOLS | PROTOCOL

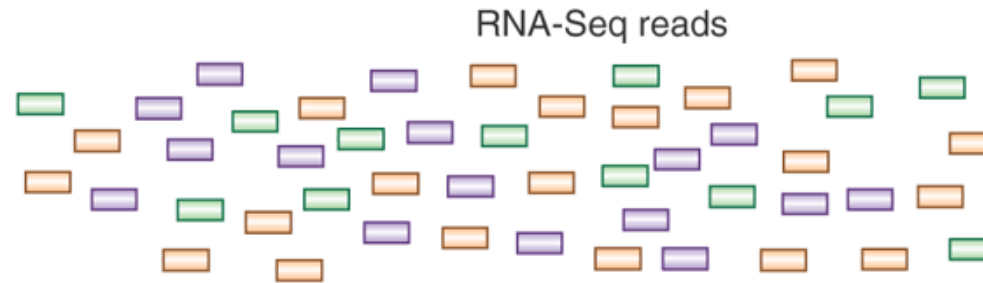
## Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks

Cole Trapnell, Adam Roberts, Loyal Goff, Geo Pertea, Daehwan Kim, David R Kelley, Harold Pimentel, Steven L Salzberg, John L Rinn & Lior Pachter

[Affiliations](#) | [Contributions](#) | [Corresponding author](#)

*Nature Protocols* 7, 562–578 (2012) | doi:10.1038/nprot.2012.016  
Published online 01 March 2012

# Transcript Reconstruction from RNA-Seq Reads

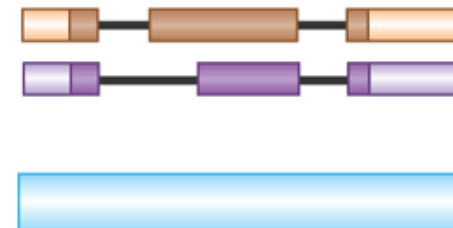


RNA-Seq reads

Assemble transcripts  
*de novo*



Align transcripts  
to genome



## End-to-end **Transcriptome**-based RNA-Seq Analysis Software Package

NATURE PROTOCOLS | PROTOCOL

*De novo* transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis

Brian J Haas, Alexie Papanicolaou, Moran Yassour, Manfred Grabherr, Philip D Blood, Joshua Bowden, Matthew Brian Couger, David Eccles, Bo Li, Matthias Lieber, Matthew D MacManes, Michael Ott, Joshua Orvis, Nathalie Pochet, Francesco Strozzi, Nathan Weeks, Rick Westerman, Thomas William, Colin N Dewey, Robert Henschel, Richard D LeDuc, Nir Friedman & Aviv Regev

[Affiliations](#) | [Contributions](#) | [Corresponding authors](#)

Nature Protocols **8**, 1494–1512 (2013) | doi:10.1038/nprot.2013.084

Published online 11 July 2013

# *De novo* transcriptome assembly

No genome required

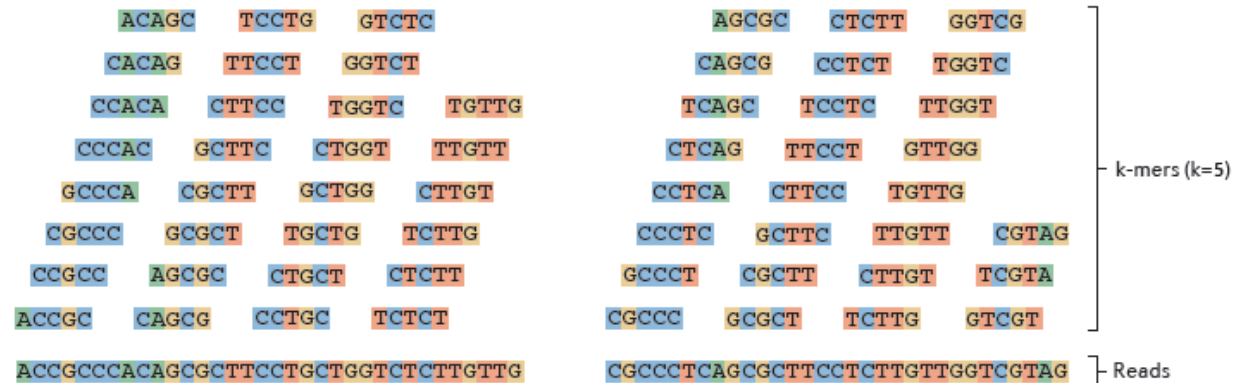
Empower studies of non-model organisms

- expressed gene content
- transcript abundance
- differential expression

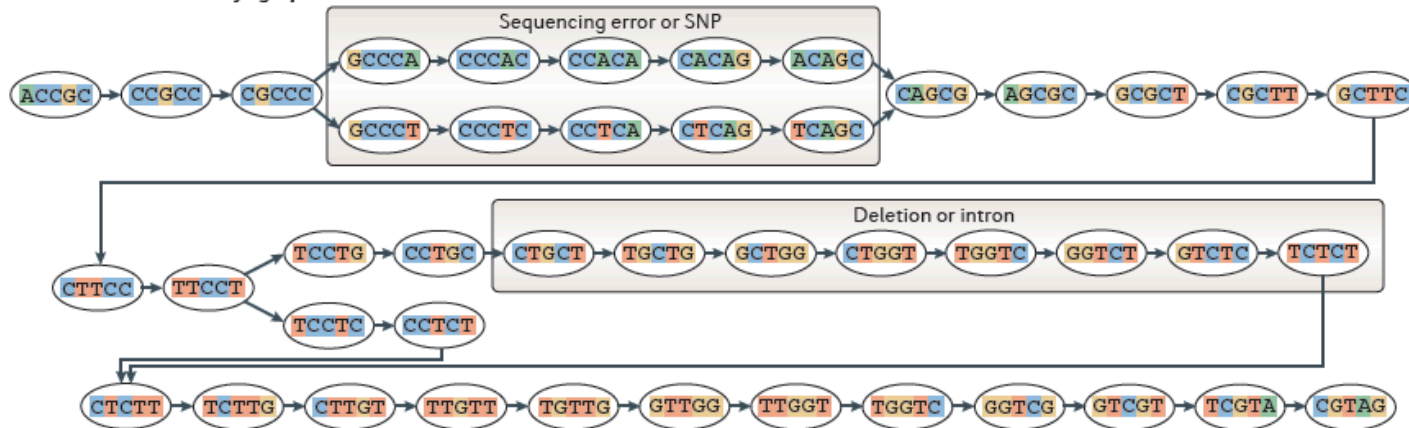
The General Approach to  
*De novo* RNA-Seq Assembly  
Using De Bruijn Graphs

# Sequence Assembly via De Bruijn Graphs

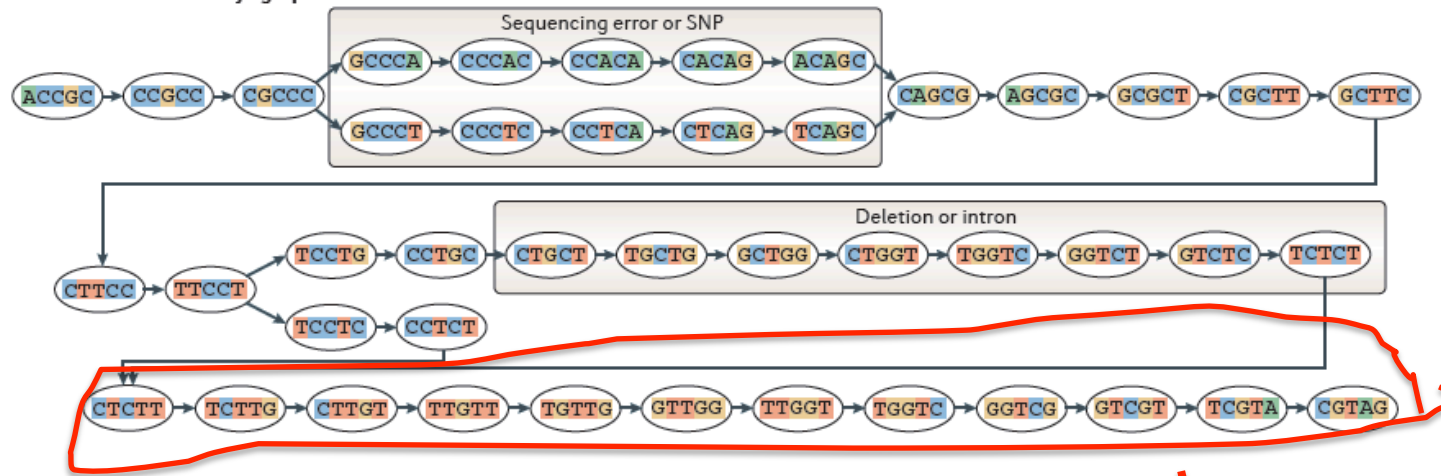
**a** Generate all substrings of length  $k$  from the reads



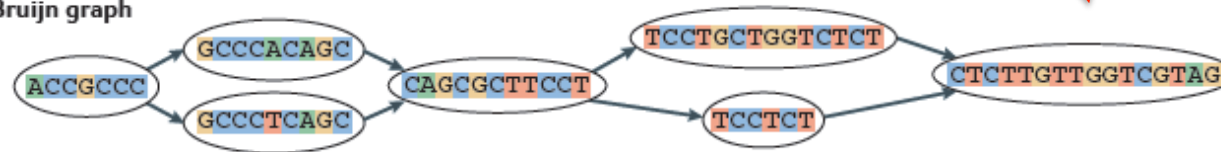
**b** Generate the De Bruijn graph



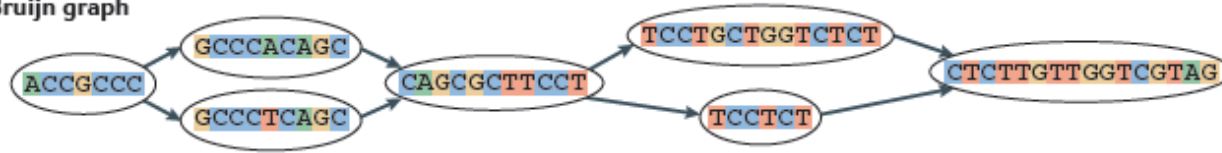
**b Generate the De Bruijn graph**



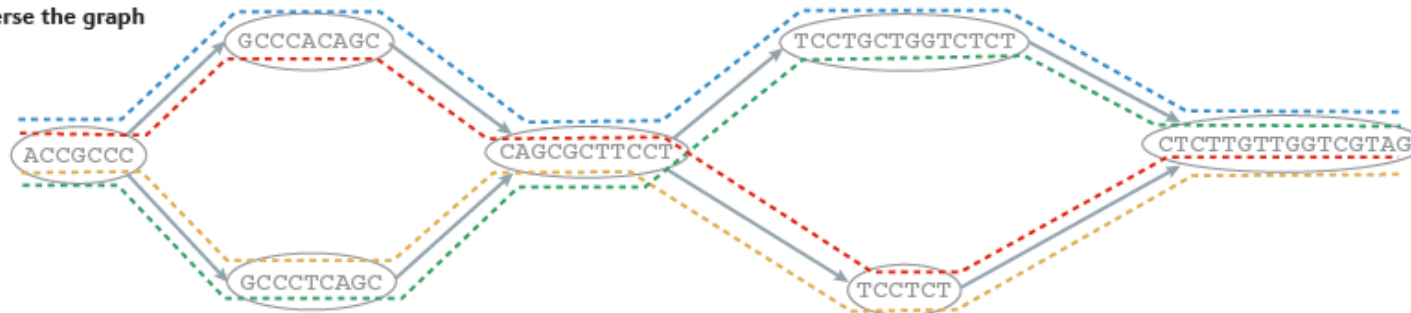
**c Collapse the De Bruijn graph**



**c Collapse the De Bruijn graph**



**d Traverse the graph**



**e Assembled isoforms**

- - - - - ACCGCCACAGCGCTTCCTGCTGGTCTCTTGGTGGTCGTAG  
 - . . . . - ACCGCCACAGCGCTTCCT - - - - - CTTGGTGGTCGTAG  
 - . . . . - ACCGCCCTCAGCGCTTCCT - - - - - CTTGGTGGTCGTAG  
 - . . . . - ACCGCCCTCAGCGCTTCCTGCTGGTCTCTTGGTGGTCGTAG

# Contrasting Genome and Transcriptome Assembly

## Genome Assembly

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

## Transcriptome Assembly

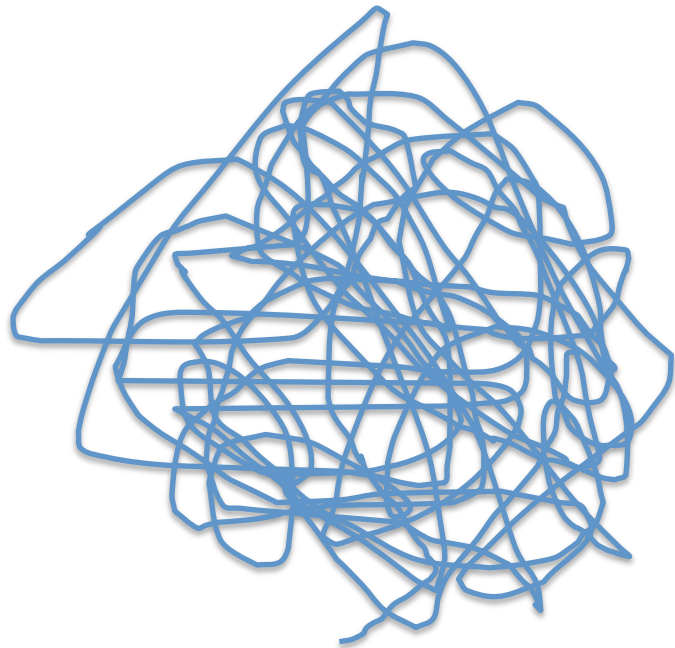
- Exponentially distributed coverage levels
- Multiple contigs per locus (alt splicing)
- Strand-specific





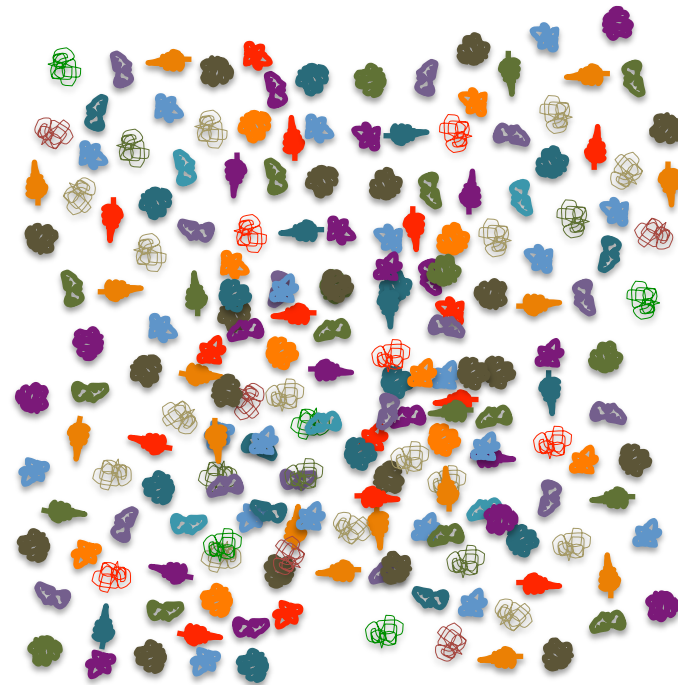
# Trinity Aggregates Isolated Transcript Graphs

**Genome Assembly**  
Single Massive Graph



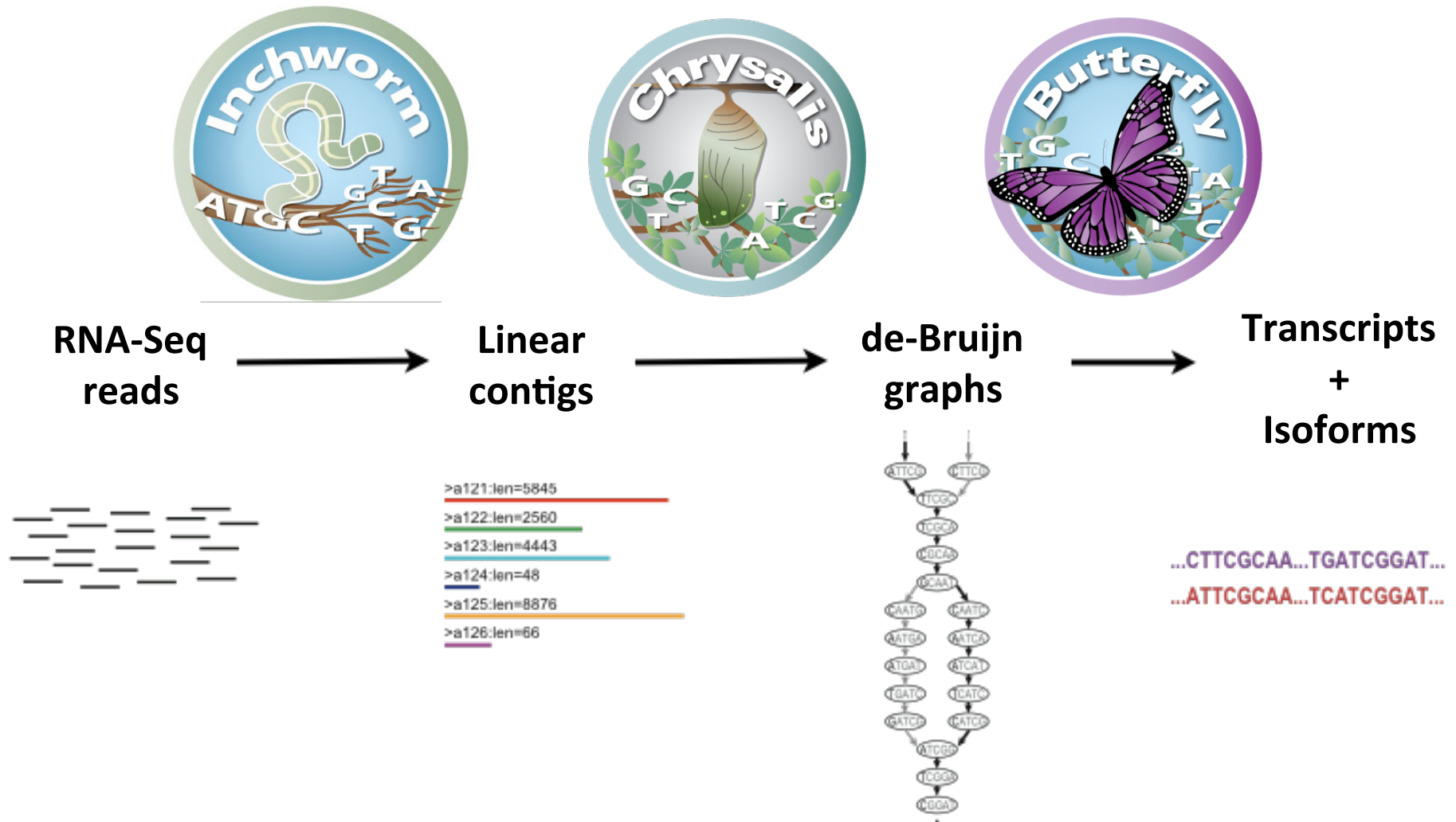
Entire chromosomes represented.

**Trinity Transcriptome Assembly**  
Many Thousands of Small Graphs



Ideally, one graph per expressed gene.

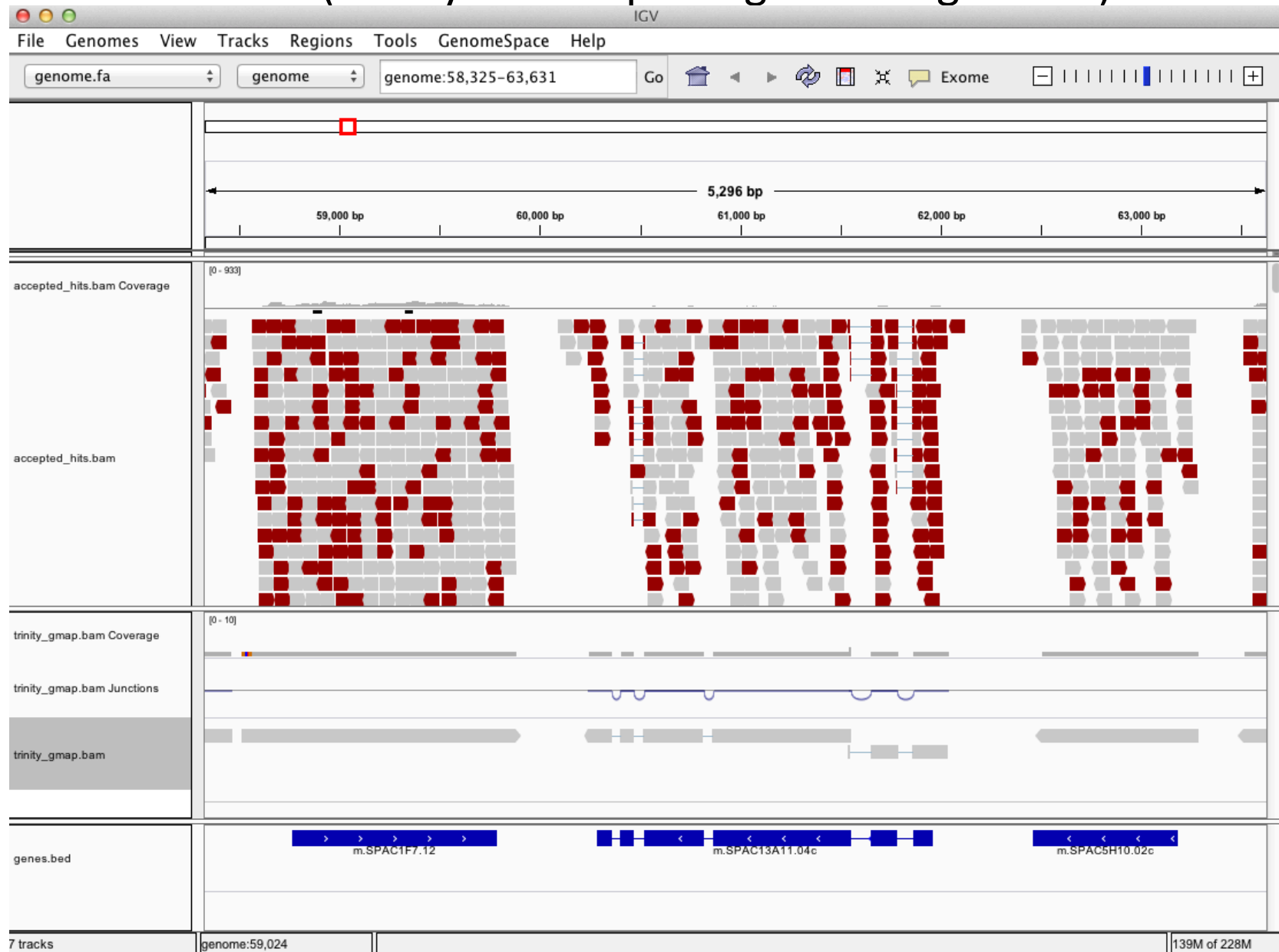
# Trinity – How it works:



Thousands of disjoint graphs

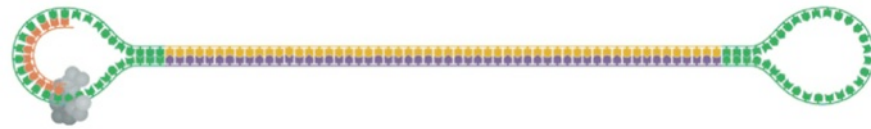
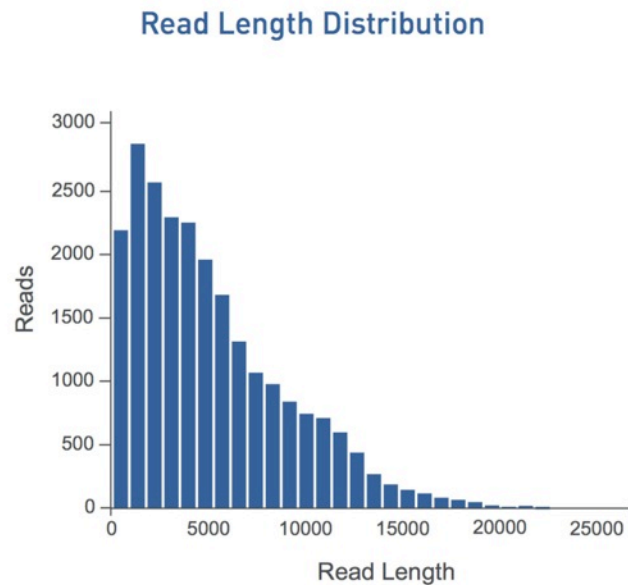


# Can align Trinity transcripts to genome scaffolds to examine intron/exon structures (Trinity transcripts aligned using GMAP)



## An alternative: Pacific Biosciences (PacBio)

- Pros: Long reads (average 4.5 kbp), can give you full length transcripts in one read
- Cons: High error rate on longer fragments (15%), expensive



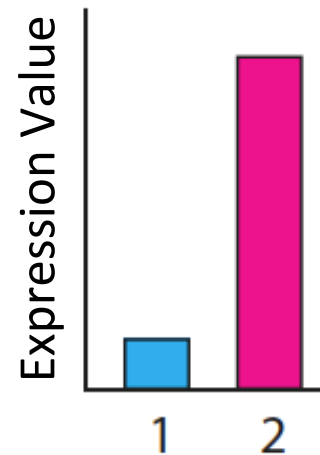
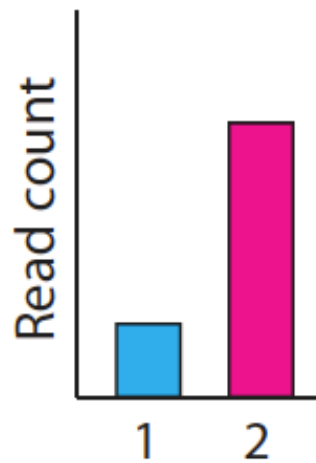
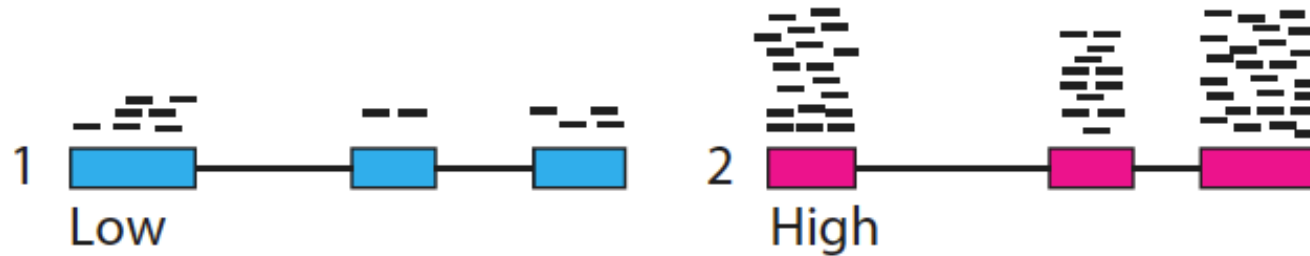


# Abundance Estimation

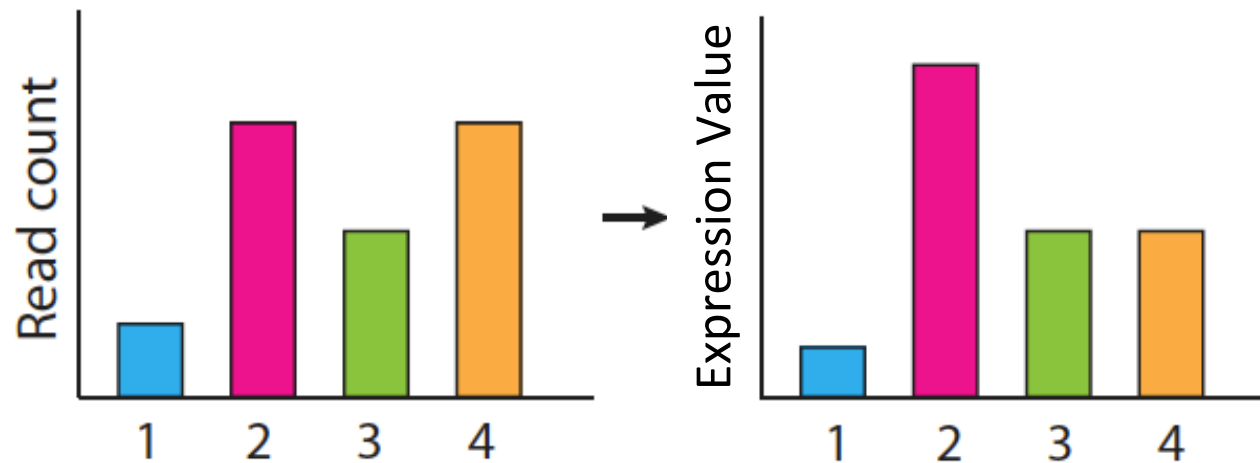
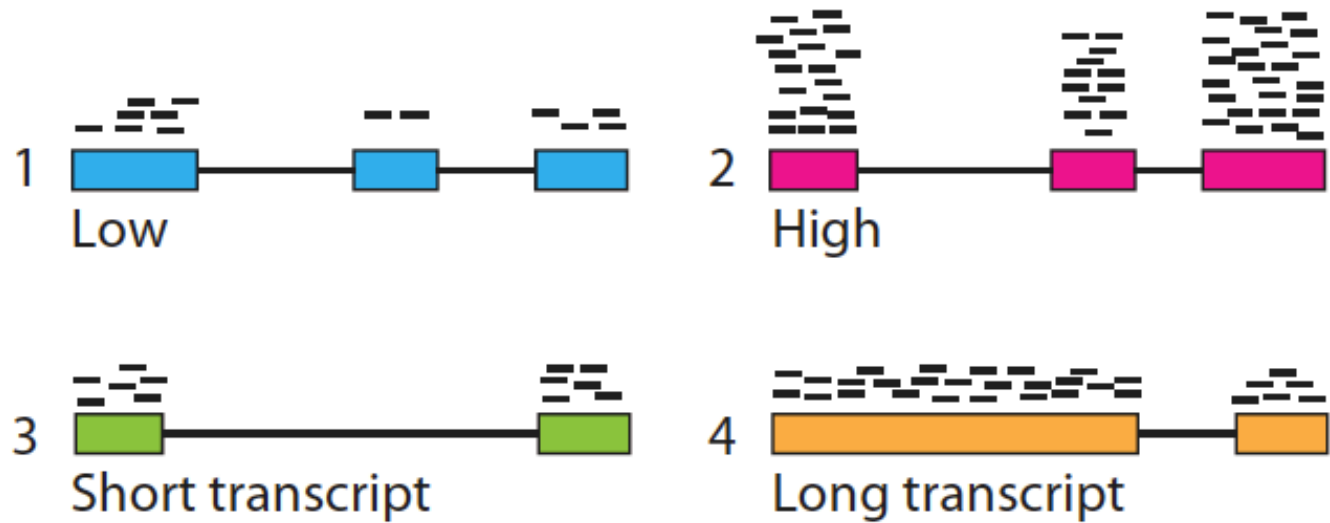
(Computing Expression Values)

# Calculating expression of genes and transcripts

---



# Calculating expression of genes and transcripts





# Normalized Expression Values

Gene expression for RNAseq analysis is based in how many reads map to an specific gene. For comparison purposes the counts needs to be normalized. There are different methodologies.

- **RPKM** (Mortazavi et al. 2008): Reads per Kilobase of Exon perMillion of Mapped reads.
- **Upper-quartile** (Bullard et al. 2010): Counts are divided per upper quartile of counts with at least one read.
- **TMM** (Robinson and Oshlack, 2010): Trimmed Means of M values (EdgeR).
- **FPKM** (Trapnell et al. 2010): Fragment per Kilobase of exon per Million of Mapped fragments (Cufflinks).

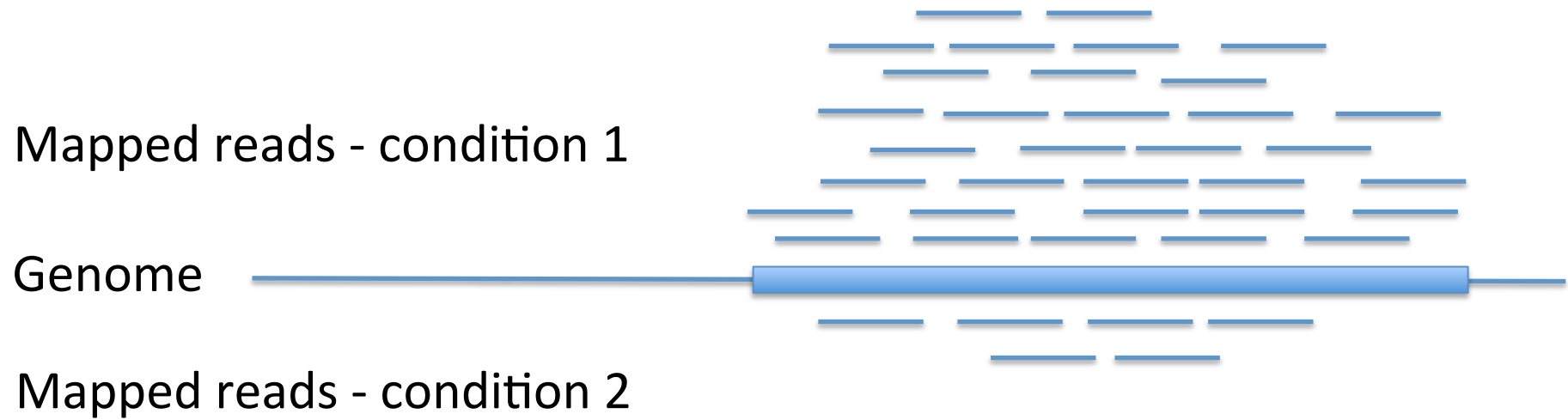
# Normalized Expression Values

- 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  
**P**er **K**ilobase of transcript  
per total **M**illion fragments mapped

**FPKM**

# Differential Expression Analysis Using RNA-Seq

# Differential expression



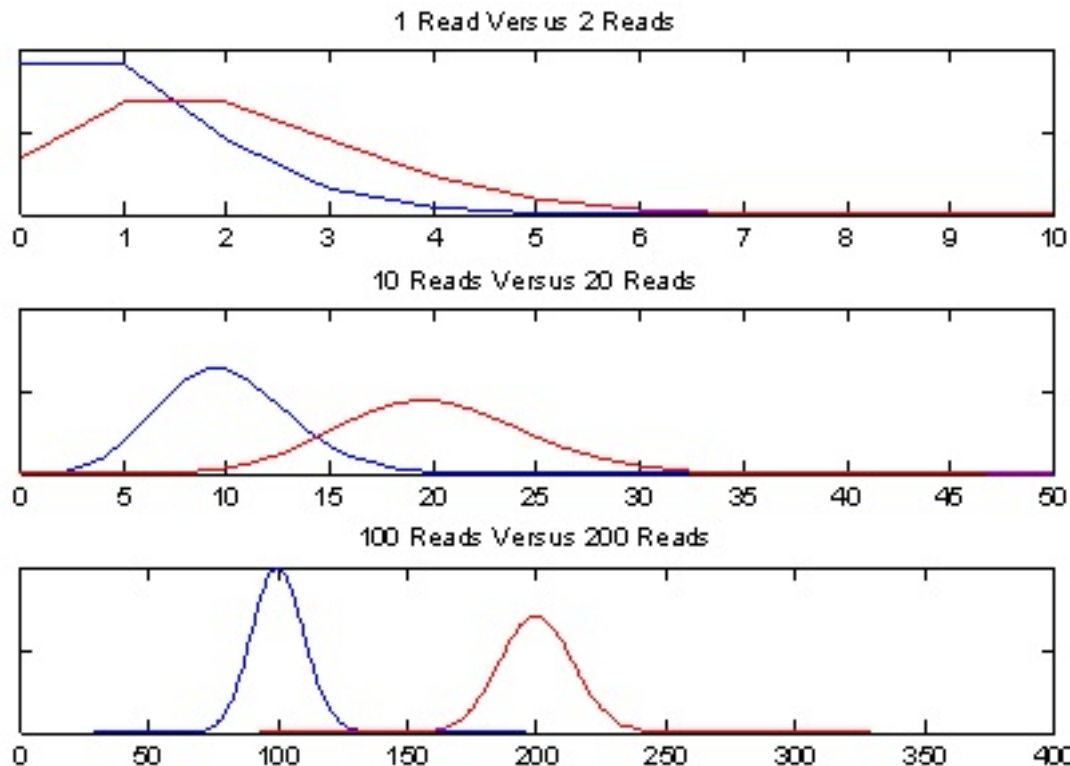
# Diff. Expression Analysis Involves

- 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.

# More Counts = More Statistical Power

Example: 5000 total reads per sample.

Observed 2-fold differences in read counts.

	<b>SampleA</b>	<b>Sample B</b>	<b>Fisher's Exact Test (P-value)</b>
<b>geneA</b>	1	2	1.00
<b>geneB</b>	10	20	0.098
<b>geneC</b>	100	200	<b>&lt; 0.001</b>

# Tools for DE analysis with RNA-Seq



ShrinkSeq  
NoiSeq  
baySeq  
Vsf  
Voom  
SAMseq  
TSPM  
DESeq  
EBSeq  
NBPSeq  
edgeR

+ other (not-R)  
including CuffDiff

See: <http://www.biomedcentral.com/1471-2105/14/91>



# Tools for DE analysis with RNA-Seq

Software	Normalization	Notes	URL
ERANGE	RPKM	Python	<a href="http://woldlab.caltech.edu/wiki/RNASeq">http://woldlab.caltech.edu/wiki/RNASeq</a>
Scripture	RPKM	Java	<a href="http://www.broadinstitute.org/software/scripture">http://www.broadinstitute.org/software/scripture</a>
BitSeq*	RPKM	R/Bioconductor, Calculate DE	<a href="http://www.bioconductor.org/packages/2.12/bioc/html/BitSeq.html">http://www.bioconductor.org/packages/2.12/bioc/html/BitSeq.html</a>
EdgeR	TMM	R/Bioconductor, Calculate DE	<a href="http://www.bioconductor.org/packages/2.11/bioc/html/edgeR.html">http://www.bioconductor.org/packages/2.11/bioc/html/edgeR.html</a>
Cufflinks*	FPKM	Isoforms, Calculate DE	<a href="http://cufflinks.cbc.umd.edu/">http://cufflinks.cbc.umd.edu/</a>
MMSEQ*	FPKM	Isoforms, Haplotypes	<a href="http://bgx.org.uk/software/mmseq.html">http://bgx.org.uk/software/mmseq.html</a>
RSEM*	FPKM	Calculate DE (EBSeq)	<a href="http://deweylab.biostat.wisc.edu/rsem/README.html">http://deweylab.biostat.wisc.edu/rsem/README.html</a>

Glaus P. *et al* (2012) *Bioinformatics* 28:1721-1728 doi:10.1093/bioinformatics/bts260

# Differential Gene Expression

Statistical test to evaluate if one gene has an differential expression between two or more conditions. These test can be based in different methodologies.

- **Negative binomial distribution** (DESeq, CuffLinks).
- **Bayesian methods for the negative binomial distribution** (EdgeR, BaySeq, BitSeq).
- **Non-parametric:** models the noise distribution of count changes by contrasting fold-change differences (M) and absolute expression differences (D) (NOISeq).

# Tools for DE analysis with RNA-Seq

Software	Normalization	Need Replicas	Input	URL
EdgeR	Library Size / TMM	Yes	Raw Counts	<a href="http://www.bioconductor.org/packages/2.11/bioc/html/edgeR.html">http://www.bioconductor.org/packages/2.11/bioc/html/edgeR.html</a>
DESeq	Library Size	No	Raw Counts	<a href="http://bioconductor.org/packages/release/bioc/html/DESeq.html">http://bioconductor.org/packages/release/bioc/html/DESeq.html</a>
baySeq	Library Size	Yes	Raw Counts	<a href="http://www.bioconductor.org/packages/2.11/bioc/html/baySeq.html">http://www.bioconductor.org/packages/2.11/bioc/html/baySeq.html</a>
NOISeq	Library Size / RPKM / UpperQ	No	Raw or Normalized Counts	<a href="http://bioinfo.cipf.es/noiseq/doku.php?id=start">http://bioinfo.cipf.es/noiseq/doku.php?id=start</a>

# Explorative Data Mining Methods

For **gene expression** there are some common tasks and associated methods for the **data mining**:

- Clustering of the expression values and principal component analysis to reduce the variables.
- Classification using Gene Ontology terms and metabolic annotations
- Summarization visualizing the expression data through heat maps.

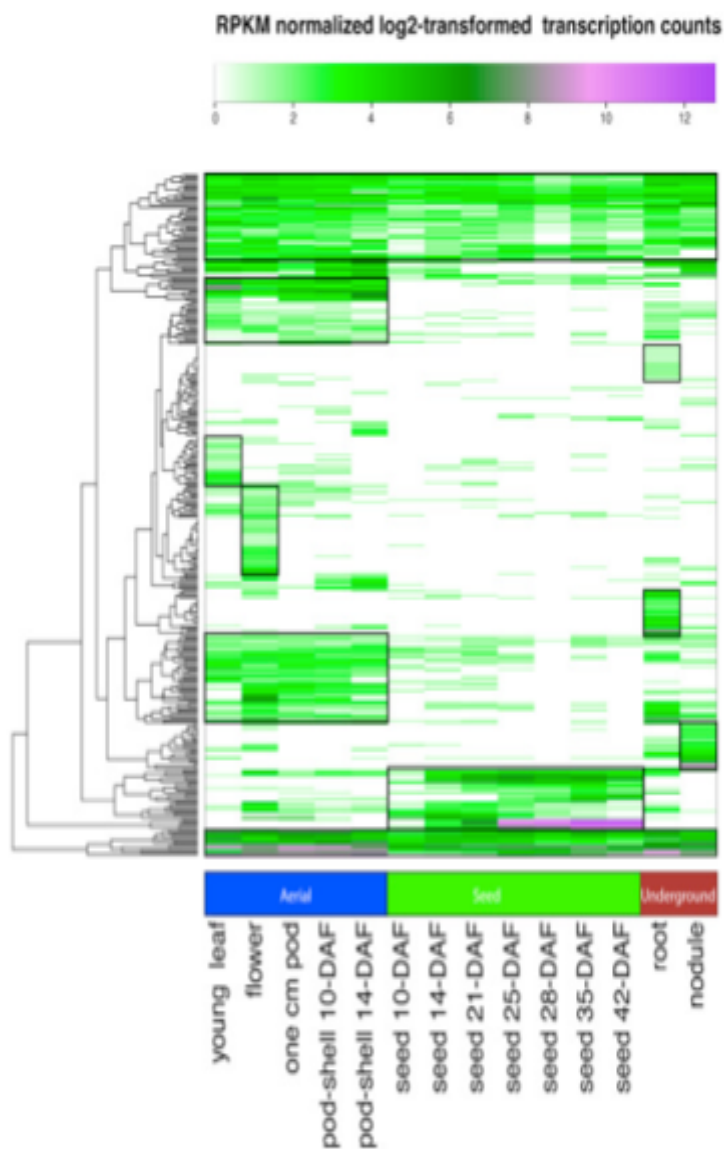
# Cluster Analysis and Visualization

Cluster analysis or clustering is the task of **assigning a set of objects into groups** (called clusters) so that the objects in the same cluster are more similar (in some sense or another) to each other than to those in other clusters. Clustering is a main task of **explorative data mining**. The most used **clustering algorithm** for gene expression are:

- **Hierarchical clustering (HCL)**, where the distance between elements is used to build the clusters.
- **K-means clustering (KMC)**, where clusters are represented by a vector. The number of clusters is fixed and the elements are assigned based in its distance to the vector.

# Cluster Analysis and Visualization

Severin AJ et al., 2010 BMC Plant Biology, 10:160



One of the most common classification data mining method is the use of gene annotations such as GO terms or metabolic annotations. These methodologies compare two groups between them to find if there are term more represented in one group than in other. Some examples are:

- **Gene Set Enrichment Analysis (GSEA)**, computational method that determines whether an a priori defined set of genes shows statistically significant.
- **Profile comparisons**, each group defines a profile based in the annotation groups (generally GO terms). Profiles are compared to find if they are significantly different.

# Use of transcripts

- Transcripts can be assembled de novo or from mapped reads and then used in gene expression/differential expression studies
- Can be functionally annotated

# Functional annotation

- Take transcripts from Cufflinks or Trinity
- Annotate the sequences functionally in Blast2GO





# Blast2GO



/Users/hobbe/Documents/Artemis\_files\_current/blast2go\_20101001\_0816.dat - Blast2GO V.2.4.4

File Blast Mapping Annotation Analysis Statistics Select Tools View Info

GO:0007067,GO:0016021 transport;binding;apoptos SPO\_2518,DDX18\_HUMAN

nr	sequence name	seq description	length	#...	min. eValue	sim mean	#G...	GO IDs	Enzyme	InterPro	
<input type="checkbox"/>	3884	gene_3884 GeneMar...	c6 transcription	977	20	1.0E-171	59.85%	7	F:transcription factor activity; F:zinc ion binding; P:regulation of transcription, DNA-dependent; C:transcription factor complex; F:transporter activity; C:membrane; P:transmembrane transport		IPR005829; IPR007219
<input type="checkbox"/>	3885	gene_3885 GeneMar...	hypothetical protein NFIA_039100 [Neosartorya fischeri NRRL 181]	312	20	1.0E-39	63.15%	1	C:viral capsid		no IPS match
<input type="checkbox"/>	3886	gene_3886 GeneMar...	sin3 complex subunit	870	20	0.0	73.2%	0			
<input type="checkbox"/>	3887	gene_3887 GeneMar...	mitochondrial intermembrane space translocase subunit	87	20	1.0E-40	88.55%	5	F:metal ion binding; P:protein import into mitochondrial inner membrane; C:mitochondrial inner membrane; C:mitochondrial intermembrane space protein transporter complex; P:transmembrane transport		IPR004217; PTHR11038 (PANTHER); PTHR11038:SF8 (PANTHER)
<input type="checkbox"/>	3888	gene_3888 GeneMar...	lysyl-tRNA synthetase	592	20	0.0	73.55%	7	C:cytoplasm; P:auxin biosynthetic process; F:nucleic acid binding; F:lysine-tRNA ligase activity; P:lysyl-tRNA aminoacylation; F:ATP binding; P:lysine biosynthetic process	EC:6.1.1.6	IPR004364; IPR004365; IPR006195; IPR012340; IPR016027; IPR018149; IPR018150; G3DSA:3.30.930.10 (GENE3D), SSF5568 (SUPERFAMILY)
<input type="checkbox"/>	3889	gene_3889 GeneMar...	transcription factor conserved	1569	20	0.0	70.9%	0			
<input type="checkbox"/>	3890	gene_3890 GeneMar...	hypothetical protein [Aspergillus clavatus NRRL 1]	240	20	1.0E-51	56.25%	0			
			udp-glc gal endoplasmic reticulum nucleotide						C:integral to membrane; C:endoplasmic reticulum membrane; P:transmembrane transport; P:carbohydrate transport		IPR013657; PTHR10778 (PANTHER)

GO Graphs Application Messages Blast/IPS Results Statistics Kegg Maps

```

17:59 InterProScan for gene_8871|GeneMark.hmm|286_aa done.
17:59 -----
17:59 InterProScan Result:
17:59 InterProId: IPR001715
17:59 InterProName: Calponin-like actin-binding
17:59 InterProType: Domain
17:59 DB-Name: GENE3D - G3DSA:1.10.418.10
17:59 InterProId: IPR016146
17:59 InterProName: Calponin-homology
17:59 InterProType: Domain
17:59 DB-Name: SUPERFAMILY - SSF47576
17:59 InterProId: noIPR
17:59 InterProName: unintegrated
17:59 InterProType: unintegrated
17:59 DB-Name: PANTHER - PTHR19961
17:59 DB-Name: PANTHER - PTHR19961:SF9
  
```

Annotation already running

# KEGG-mapping

/Users/hobbe/Documents/Artemis\_files\_current/blast2go\_20101001\_0816.dat - Blast2GO V.2.4.4

File Blast Mapping Annotation Analysis Statistics Select Tools View Info

GO:0007067,GO:0016021 transport:binding:apoptosis SPO\_2518.DDX18\_HUMAN

nr	sequence name	seq description	length	#...	min. eValue	sim mean	#C...	GO IDs	Enzyme	InterPro
		succinyl- synthetase subunit						F:ATP binding; F:succinate-CoA ligase (GDP-forming)		IPR003781;
								activity: P:tricarboxylic acid cycle; C:succinate-CoA ligase		IPR005810;

GO Graphs Application Messages Blast/IPS Results Statistics **Kegg Maps**

**GLYCEROLIPID METABOLISM**

Key enzymes and EC numbers shown in the map:

- 1.1.1.2 (red)
- 1.1.1.21 (yellow)
- 1.1.1.72 (pink)
- 2.3.1.15 (orange)
- 2.3.1.13 (blue)
- 2.3.1.20 (green)
- 2.3.1.18 (yellow)
- 2.7.1.107 (light-red)
- 2.7.1.29
- 2.7.1.31
- 1.2.1.3
- 4.2.1.30
- 3.1.3.21
- 2.7.1.30
- 3.1.1.23
- 2.7.1.94
- 3.1.1.3
- 3.1.1.34
- 3.1.1.3
- 3.1.1.34
- 2.3.1.20
- 2.7.8.20
- 2.4.1.
- 241157
- 3.1.3.4
- 2.7.1.107

Pathways

- Pentose phosphate pathway
- Fructose and mannose metabolism
- Butanoate metabolism
- Carbon fixation in photosynthetic organisms
- Lysine degradation
- Tyrosine metabolism
- Methane metabolism
- Glyoxylate and dicarboxylate metabolism
- Glycerolipid metabolism**
- Glutathione metabolism
- Selenoamino acid metabolism
- Phenylalanine metabolism
- Benzoate degradation via CoA ligation
- Valine, leucine and isoleucine biosynthesis
- Reductive carboxylate cycle (CO2 fixation)
- Galactose metabolism
- Phenylalanine, tyrosine and tryptophan biosynthesis
- N-Glycan biosynthesis
- Photosynthesis
- Drug metabolism - other enzymes
- Sulfur metabolism
- Fatty acid biosynthesis
- Inositol phosphate metabolism
- beta-Alanine metabolism
- Drug metabolism - cytochrome P450
- Pantothenate and CoA biosynthesis
- Biosynthesis of unsaturated fatty acids
- Cyanoamino acid metabolism
- Terpenoid backbone biosynthesis
- Histidine metabolism
- T cell receptor signaling pathway
- Tropine, piperidine and pyridine alkaloid biosynthesis
- One carbon pool by folate
- Pentose and glucuronate interconversions
- Phosphatidylinositol signaling system

Color	Enzyme	Sequences
red	ec:1.1.1.2 - alcohol dehydrogenase (NADP+)	gene_674 GeneMark.hmm 333_aa, gene_5801 GeneMark.hmm 312_aa
yellow	ec:2.3.1.158 - phospholipid:diacylglycerol acyltransferase	gene_2604 GeneMark.hmm 188_aa, gene_6532 GeneMark.hmm 505_aa
orange	ec:2.3.1.51 - 1-acylglycerol-3-phosphate O-acyltransferase	gene_176 GeneMark.hmm 429_aa, gene_6693 GeneMark.hmm 292_aa
green	ec:2.3.1.20 - diacylglycerol O-acyltransferase	gene_176 GeneMark.hmm 429_aa, gene_7213 GeneMark.hmm 521_aa, gene_8170 GeneMark.hmm 470_aa
blue	ec:2.3.1.15 - glycerol-3-phosphate O-acyltransferase	gene_886 GeneMark.hmm 748_aa, gene_2640 GeneMark.hmm 823_aa
pink	ec:1.1.1.72 - glycerol dehydrogenase (NADP+)	gene_3376 GeneMark.hmm 325_aa, gene_4577 GeneMark.hmm 326_aa
violet	ec:1.2.1.3 - aldehyde dehydrogenase (NAD+)	gene_2201 GeneMark.hmm 497_aa, gene_5247 GeneMark.hmm 502_aa, gene_5611 GeneMark.hmm 471_aa
light-red	ec:2.7.1.107 - diacylglycerol kinase	gene_5292 GeneMark.hmm 409_aa

Annotation already running

# Acknowledgement

Henrik, SciLifeLab