

# Introduction to RNA-Seq

Dhivya Arasappan

(With some slides borrowed from Scott Hunicke-Smith and Jeff Barrick)

# Logistics

- Commands that I will run today also on Biolteam wiki:
  - <https://wikis.utexas.edu/display/bioiteam/Introduction+to+RNA+Seq+Short+Course+Commands>
- If you want to run them too, open terminal and ssh into lonestar.

# Resources

- Biolteam Wiki- Bookmark it!

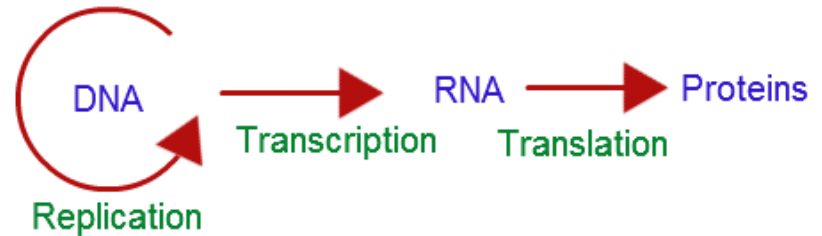
<https://wikis.utexas.edu/display/bioiteam>

- Summer School course materials:

[https://wikis.utexas.edu/display/bioiteam/  
Introduction+to+RNA+Seq+Course+2014](https://wikis.utexas.edu/display/bioiteam/Introduction+to+RNA+Seq+Course+2014)

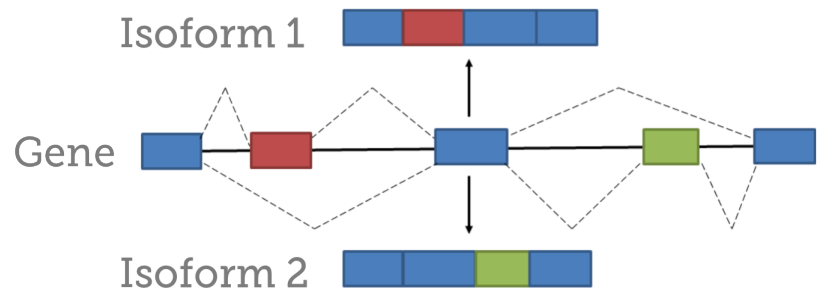
# The Purpose of RNA-Seq

- Examine the state of the transcriptome.



- Genes expression patterns vary in:

- Tissue types
- Cell types
- Development stages
- Disease conditions
- Time points



- RNA-Seq measures these expression variations using high-throughput sequencing technologies.
- Additionally, RNA-Seq allows detection of novel isoforms of genes.



# Advantages of RNA-Seq

Technology	Tiling microarray	RNA-Seq
<b>Technology specifications</b>		
Principle	Hybridization	High-throughput sequencing
Resolution	From several to 100 bp	Single base
Throughput	High	High
Reliance on genomic sequence	Yes	In some cases
Background noise	High	Low
<b>Application</b>		
Simultaneously map transcribed regions and gene expression	Yes	Yes
Dynamic range to quantify gene expression level	Up to a few-hundredfold	>8,000-fold
Ability to distinguish different isoforms	Limited	Yes
Ability to distinguish allelic expression	Limited	Yes
<b>Practical issues</b>		
Required amount of RNA	High	Low
Cost for mapping transcriptomes of large genomes	High	Relatively low

## RNA-Seq: a revolutionary tool for transcriptomics

Zhong Wang, Mark Gerstein, and Michael Snyder

*Nat Rev Genet.* 2009 January ; 10(1): 57–63. doi:10.1038/nrg2484.

# Decisions, Decisions...

- What sort of library?
- What sort of preprocessing?
- What sort of analysis pipeline?
- What sort of downstream analysis?

# What are your biological questions ?

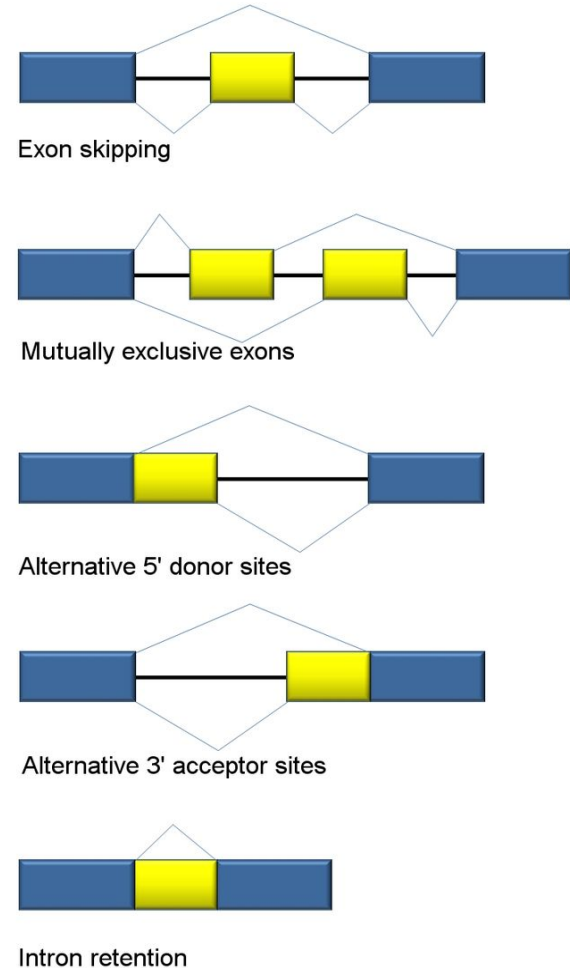
This determines how you set up your RNA Seq study and how you analyze the data.

What are you looking for?

- Annotated transcripts or Novel transcripts?
- Differential Gene expression?
- Differential exon level counts?
- Differential regulation?
- Differential splicing?

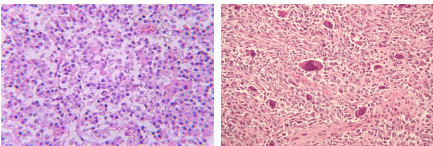
# Why is RNA-Seq Difficult?

- Biases may mean what we are seeing is not reflective of true state of the transcriptome.
- Ugh, splicing!
- Gene level, exon level?
- Multimapping, partial mapping, not mapping.
- Normalization issues
  - some datasets are larger than others, some genes are larger than others

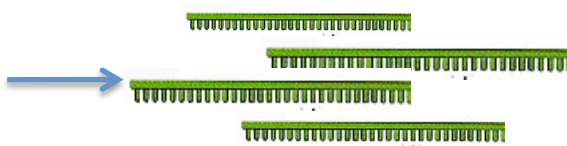


From Wikipedia- alternative splicing

# RNA-Seq... at it's Most Basic Form



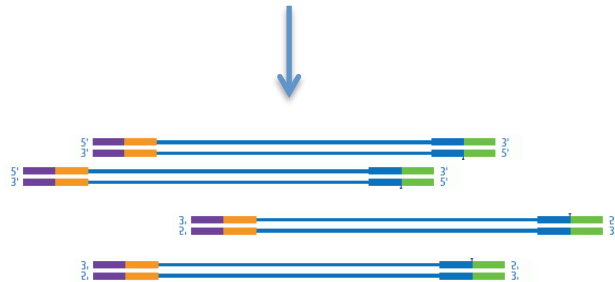
Samples from two conditions



Isolate RNA



Generate cDNA



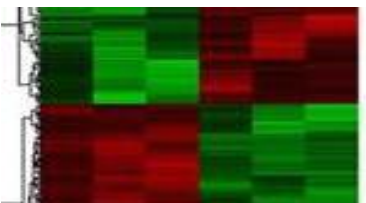
Create sequencing library by fragmenting, size selection and adding adaptors



Run sequencer



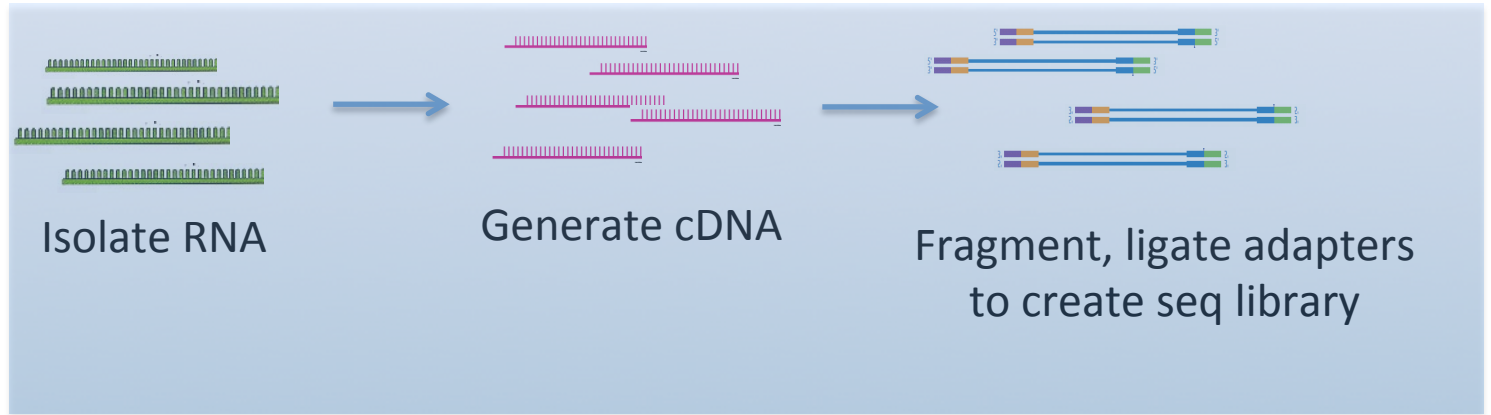
Generate short reads



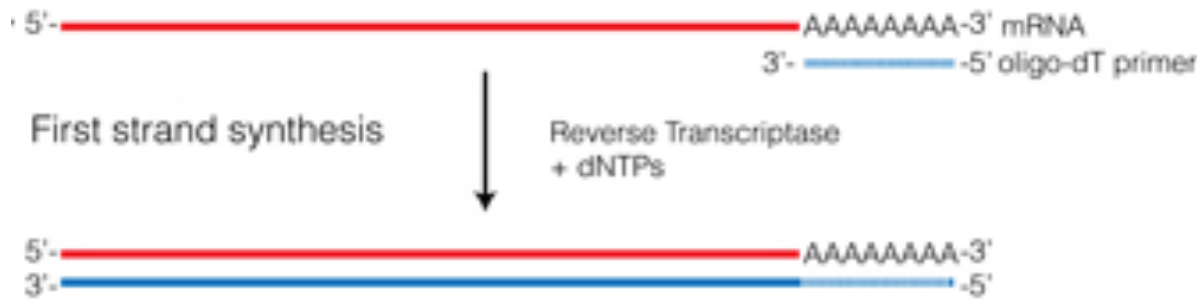
- defense response
- negative regulation of programmed cell death
- negative regulation of cell death
- chromosome organization
- regulation of cell proliferation
- response to DNA damage stimulus
- cell cycle process
- programmed cell death
- response to organic substance
- cellular response to stress
- regulation of cell death

Identify differentially expressed genes

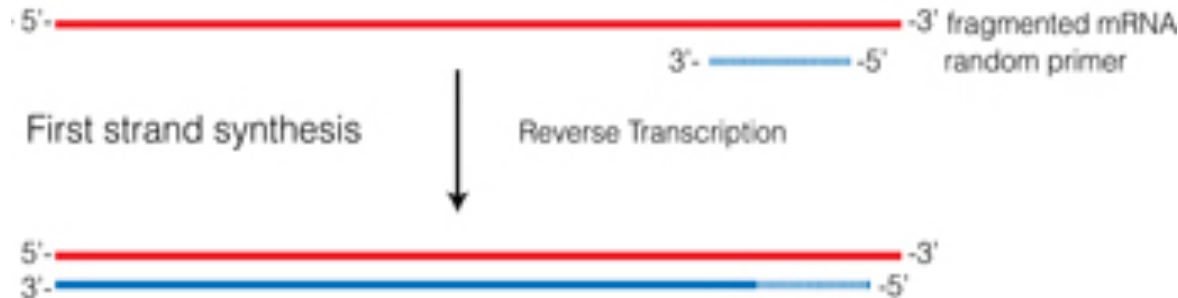
# RNA-Seq Libraries... with More Details



## A. Poly A Priming



## B. Random Priming

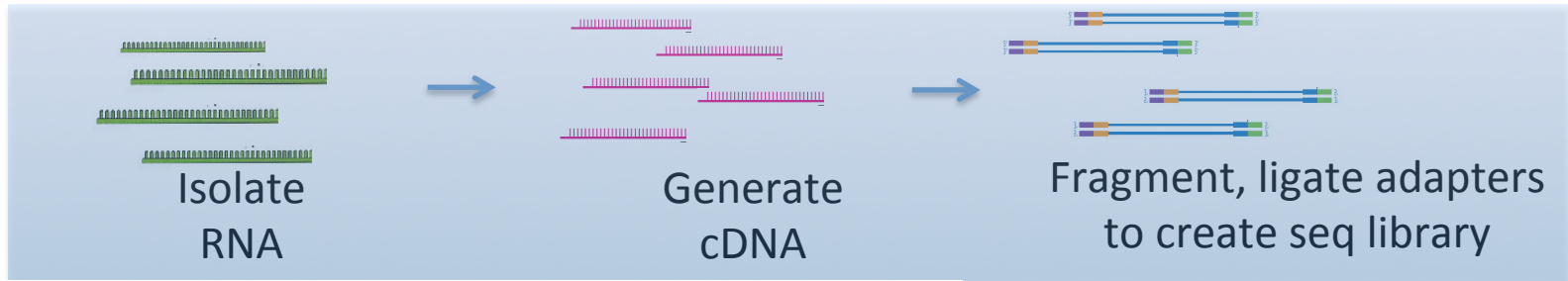


## C. rRNA Depletion



Ribominus kit

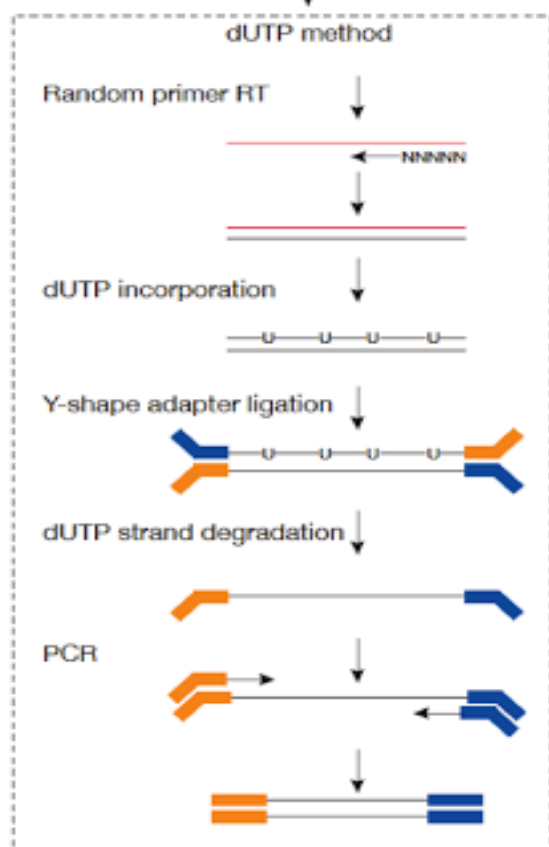
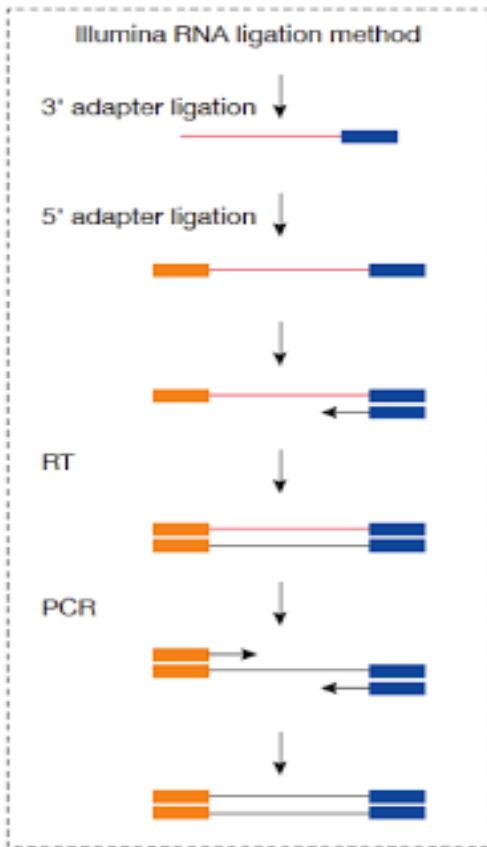
# RNA-Seq Libraries... with More Details



RNA after rRNA depletion

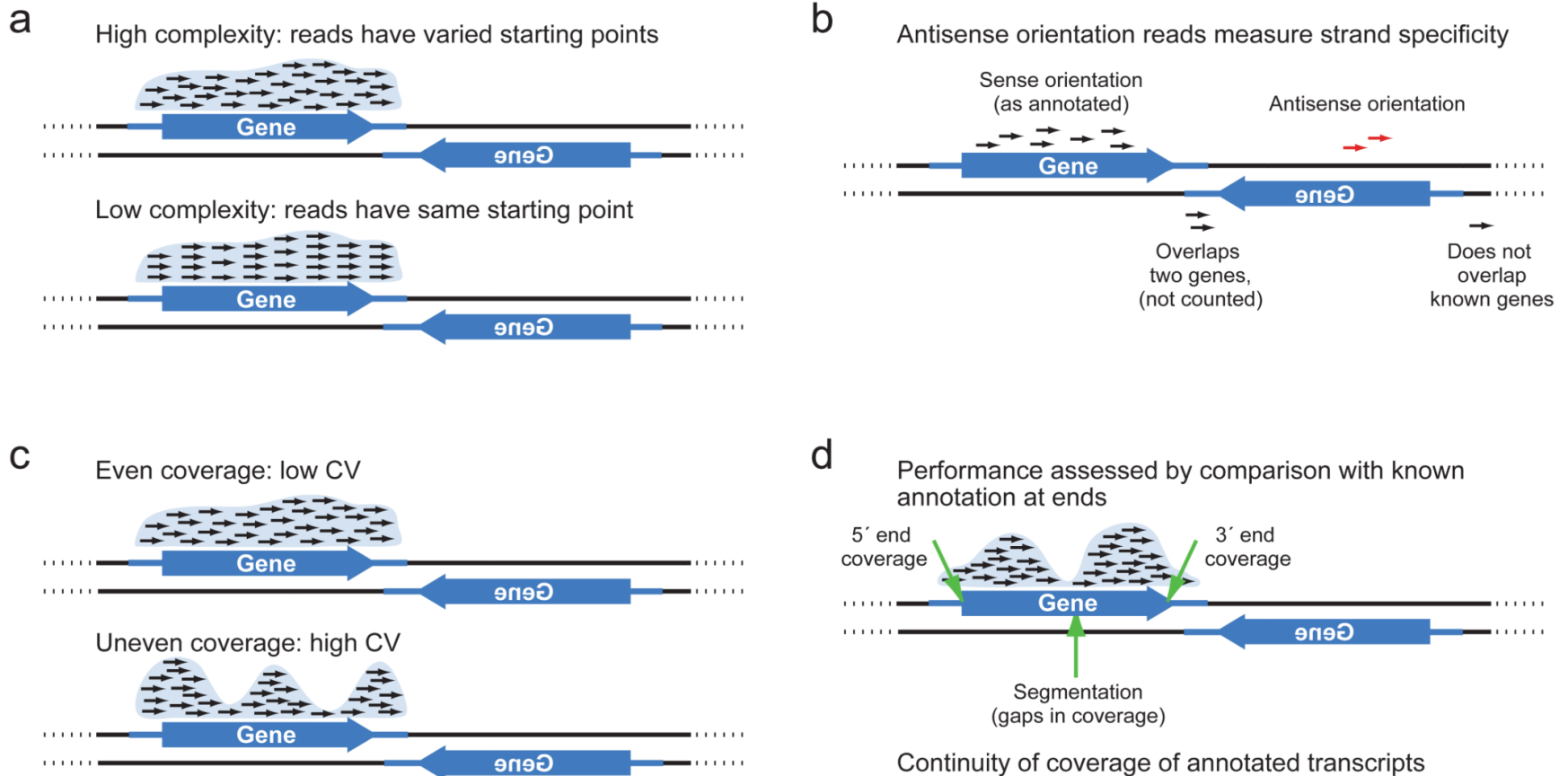
RNA fragmentation

**Second Strand Synthesis-  
Many Strand Specific  
Methods.**



Strand-specific libraries for high throughput RNA sequencing prepared without poly(A) selection, Zhang et al.

# Comparing Stranded RNA-Seq Library Protocols



**Figure 2. Key criteria for evaluation of strand-specific RNAseq libraries**

Four categories of quality assessment. Double stranded genome (black parallel lines), with Gene ORF orientation (thick blue arrow) and UTRs (thin blue line), along with mapped reads (short black arrows – reads mapped to sense strand; red – reads mapped to antisense strand). (a) Complexity. (b) Strand Specificity. (c) Evenness of coverage. (d) Comparison to known transcript structure..



# Types of Illumina Fragment Libraries

single-end



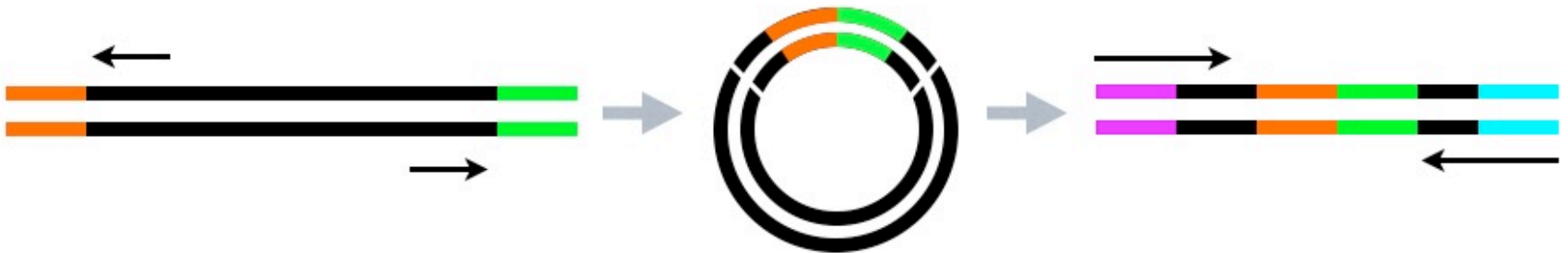
independent reads

paired-end



two inwardly oriented reads separated by ~200 nt

mate-paired



two outwardly oriented reads separated by ~3000 nt

# How do we analyze RNA-Seq data?

- **STEP 1:** EVALUATE AND MANIPULATE RAW DATA
- **STEP 2:** MAP TO REFERENCE, ASSESS RESULTS
- **STEP 3:** ASSEMBLE TRANSCRIPTS
- **STEP 4:** QUANTIFY TRANSCRIPTS
- **STEP 5:** TEST FOR DIFFERENTIAL EXPRESSION
- **STEP 6:** VISUALIZE AND PERFORM OTHER DOWNSTREAM ANALYSIS

# STEP 1 - Evaluate Raw Data

## FASTQ FORMAT

```
@HWI-EAS216_91209:1:2:454:192#0/1  
GTTGATGAATTTCTCCAGCGCGAATTTGTGGGCT  
+HWI-EAS216_91209:1:2:454:192#0/1  
B@BBBBBB@BBBBAAAA>@AABA?BBBAAB??>A?
```

**Line 1:** @read name

**Line 2:** called base sequence

**Line 3:** +read name (optional after +)

**Line 4:** base quality scores

# STEP 1 - Evaluate Raw Data

## Illumina Base Quality Scores

<http://www.asciitable.com/>

Quality character	!"#\$%&'()*+,-./0123456789:;<=>?@ABCDEFGHI
ASCII Value	33 43 53 63 73
Base Quality (Q)	0 10 20 30 40

$$\text{Probability of Error} = 10^{-Q/10}$$

(This is a **Phred** score, also used for other types of qualities.)

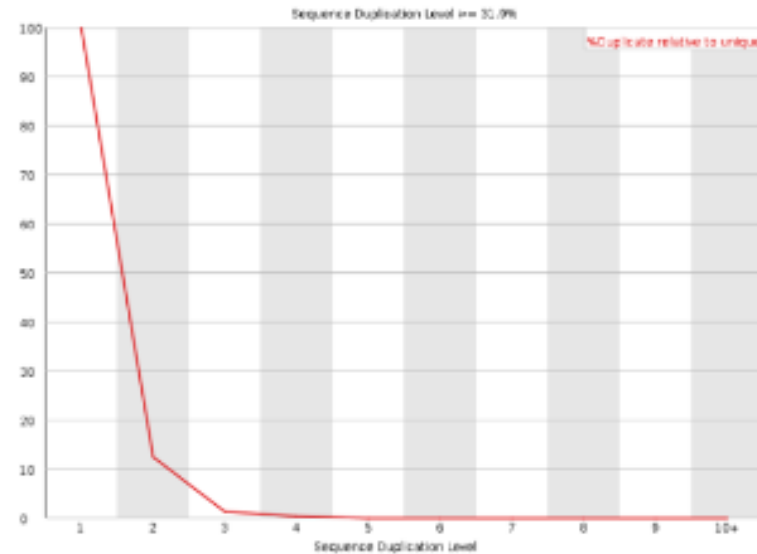
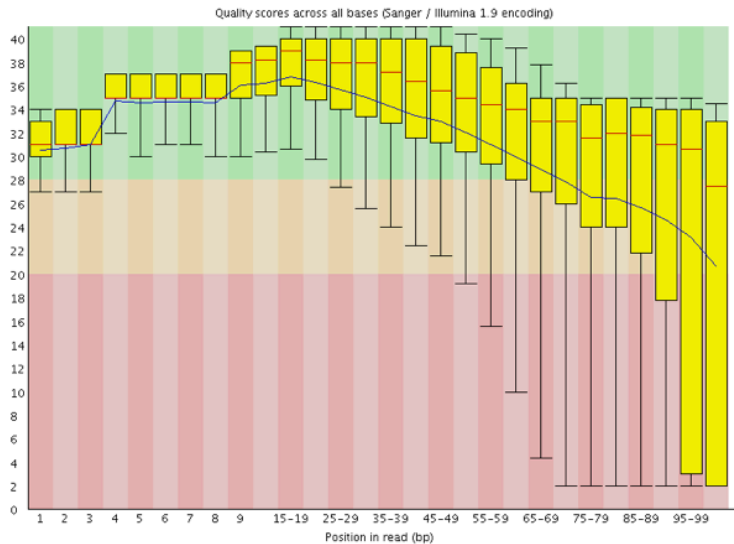
Phred Quality Score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%

Quality scores are ASCII encoded in fastq files. Different platforms/older sequencing data can have different encoding! Illumina HiSeq 2500 produces Sanger encoded data.

**Phred +33 =ASCII**

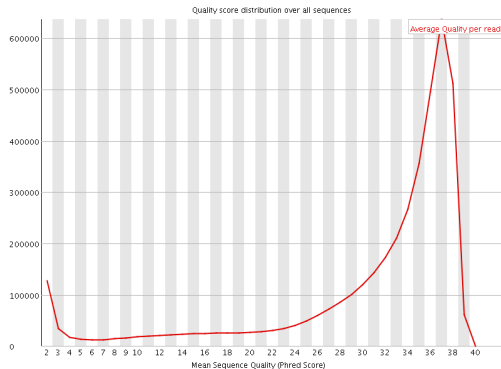
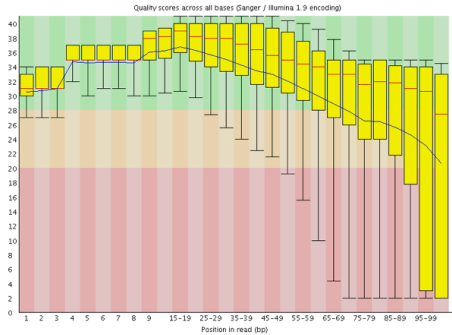
# STEP 1 - Evaluate Raw Data

- Count your reads!
- Assess quality using FastQC reports



Sequence	Count	Percentage	Possible Source
AGATCGGAAGAGCACACGTCTGAACTCCAGTCACCTCAGAATCTCGTATG	60030	5.01369306977828	TruSeq Adapter, Index 1 (97% over 37bp)
GATCGGAAGAGCACACGTCTGAACTCCAGTCACCTCAGAATCTCGTATGC	42955	3.5875926338884896	TruSeq Adapter, Index 1 (97% over 37bp)
CACACGTCTGAACTCCAGTCACCTCAGAATCTCGTATGCCGTCTTCTGCT	3574	0.29849973398946483	RNA PCR Primer, Index 40 (100% over 41bp)
CAGATCGGAAGAGCACACGTCTGAACTCCAGTCACCTCAGAATCTCGTAT	2519	0.2103863542024236	TruSeq Adapter, Index 1 (97% over 37bp)
GAGATCGGAAGAGCACACGTCTGAACTCCAGTCACCTCAGAATCTCGTAT	1251	0.10448325887543942	TruSeq Adapter, Index 1 (97% over 37bp)

# STEP 1 – Manipulate Raw Data



- Trim low quality bases
  - Fastx toolkit- **fastx\_trimmer**
    - Trim X number of low quality bases from each read.
- Filter out low quality reads
  - Fastx toolkit- **fastq\_quality\_filter**
    - Filter out reads with more than X percent of low quality bases.

- Trim Adaptor

- Fastx toolkit- **fastx\_clipper**

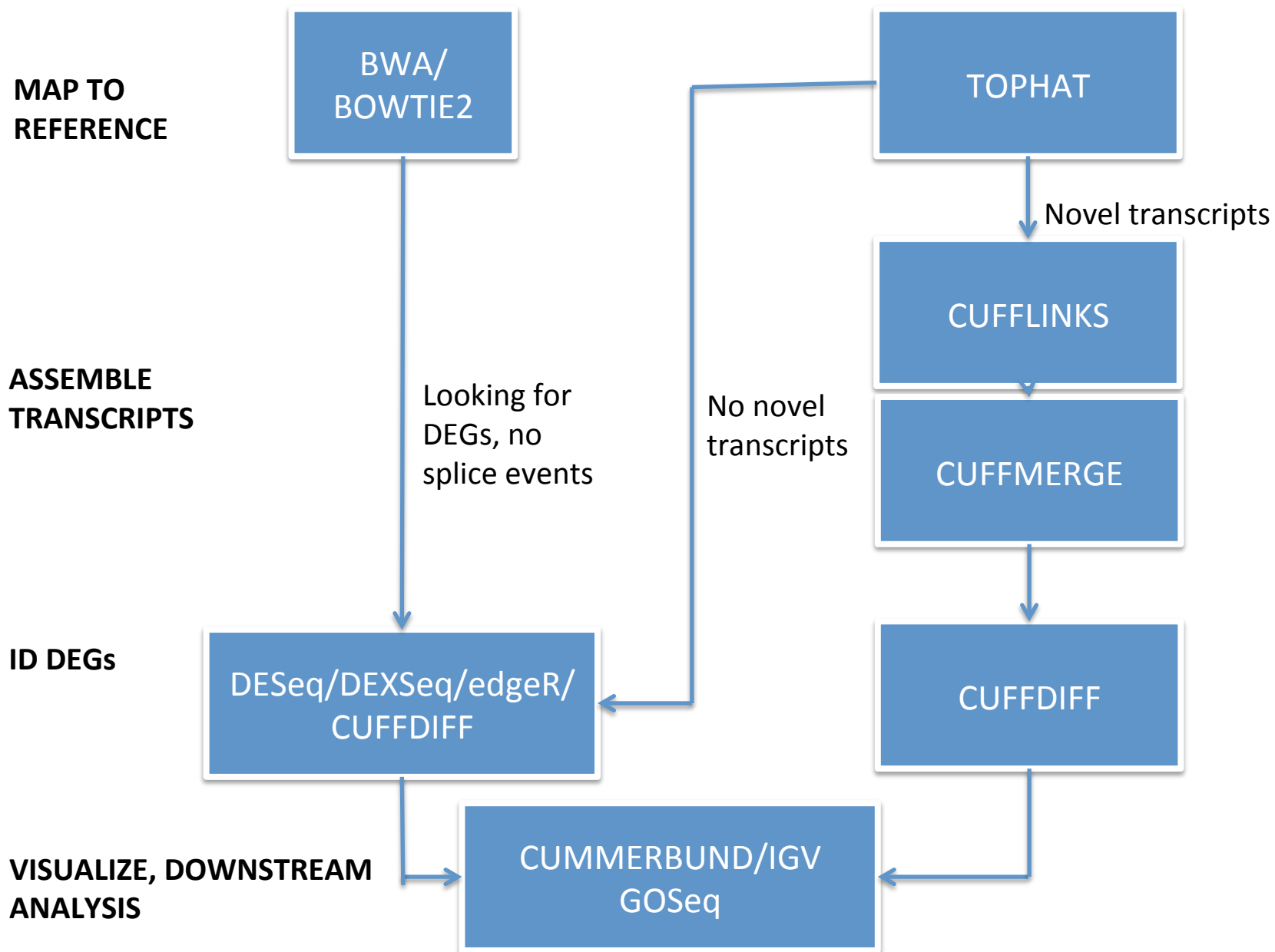
- Look for and clip a given sequence from the end of reads

- **Cutadapt**

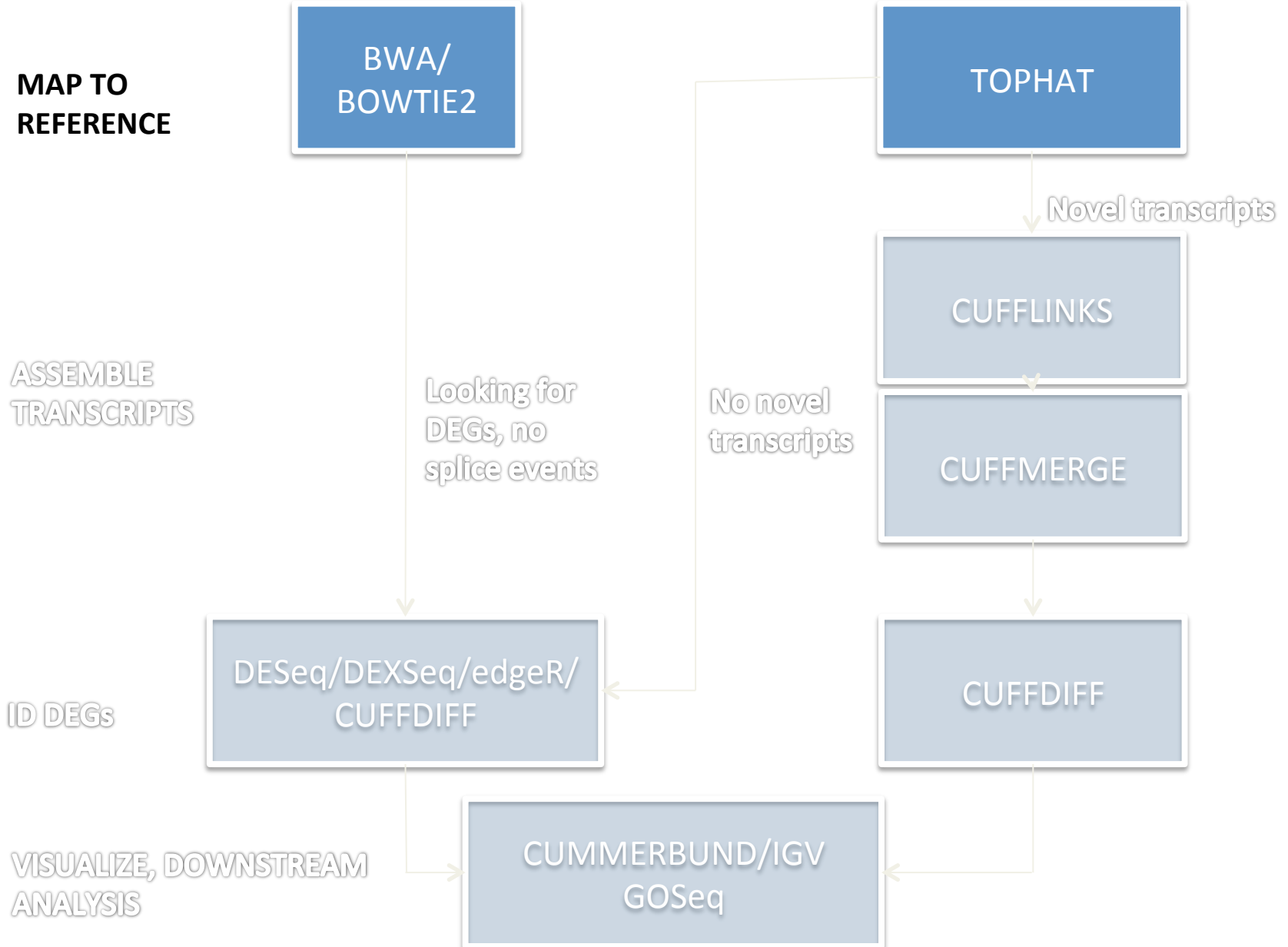
- Allows for mismatches
- Paired -end support

Sequence	Count	Percentage	Possible Source
AGATCGGAAGAGCACACGTCTGAACTCCAGTCACCTCAGAATCTCGTATG	60030	5.01369306977828	TruSeq Adapter, Index 1 (97% over 37bp)
GATCGGAAGAGCACACGTCTGAACTCCAGTCACCTCAGAATCTCGTATG	42955	3.5875926338884896	TruSeq Adapter, Index 1 (97% over 37bp)
CACACGTCTGAACTCCAGTCACCTCAGAATCTCGTATGCCGCTTCTCTGT	3574	0.2984997339894683	BNA PCR Primer, Index 40 (100% over 41bp)
CAGATCGGAAGAGCACACGTCTGAACTCCAGTCACCTCAGAATCTCGTAT	2519	0.2103863542024236	TruSeq Adapter, Index 1 (97% over 37bp)
GAGATCGGAAGAGCACACGTCTGAACTCCAGTCACCTCAGAATCTCGTAT	1251	0.10448325887543942	TruSeq Adapter, Index 1 (97% over 37bp)

# RNA-Seq Analysis Pipelines

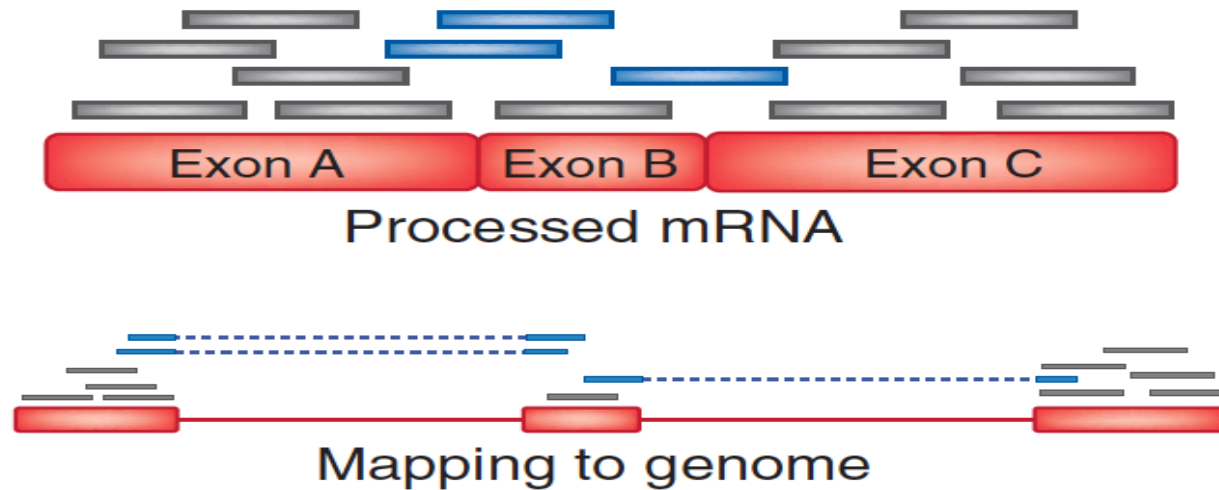


# STEP 2 - Map to Reference





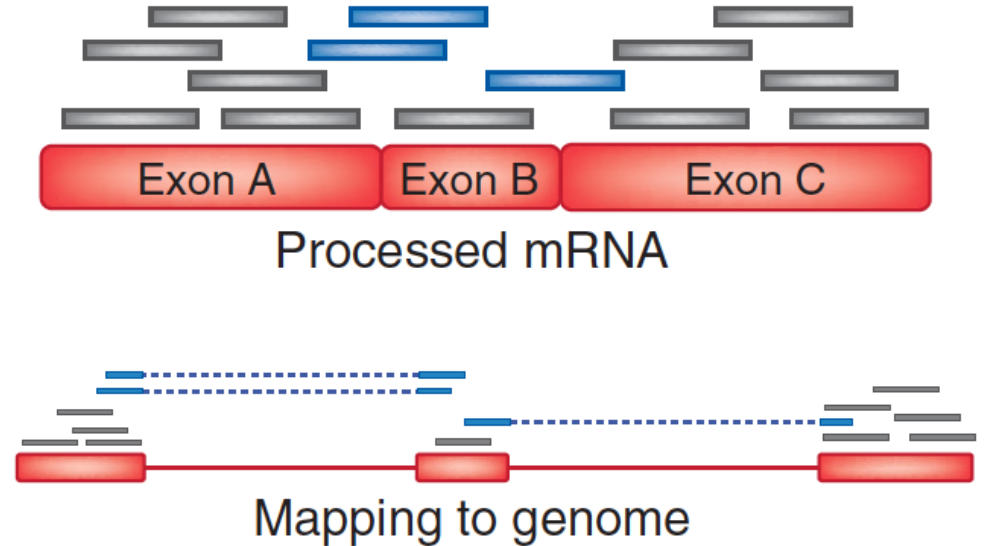
# Unspliced Mapping



Class	Category	Package	Notes
Read mapping	Seed methods	Short-read mapping package (SHRiMP) <sup>41</sup>	Smith-Waterman extension
		Stampy <sup>39</sup>	Probabilistic model
	Burrows-Wheeler transform methods	Bowtie <sup>43</sup> BWA <sup>44</sup>	Incorporates quality scores

# Spliced mapping

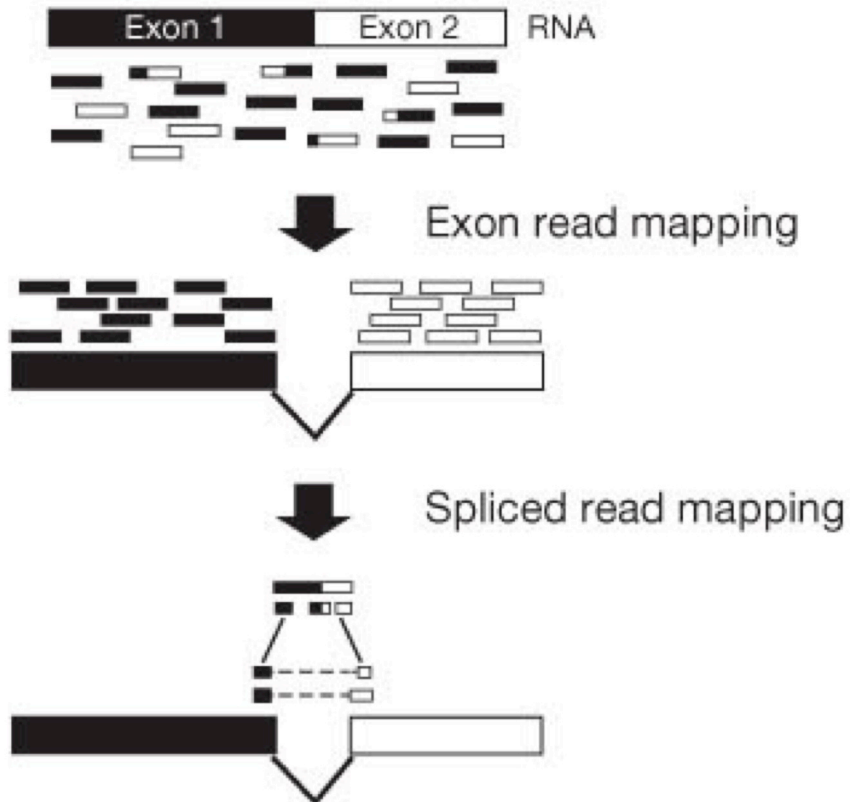
- Needed for identifying and quantifying splice variants from RNA Seq data.
- Tools:
  - Tophat
  - SpliceMap
  - MapSplice
  - STAR
  - RUM



Trapnell, C. & Salzberg, S. L. How to map billions of short reads onto genomes. *Nature Biotech.* **27**, 455–457 (2009).

# Spliced mapping

## a Exon-first approach



## b Seed-extend approach

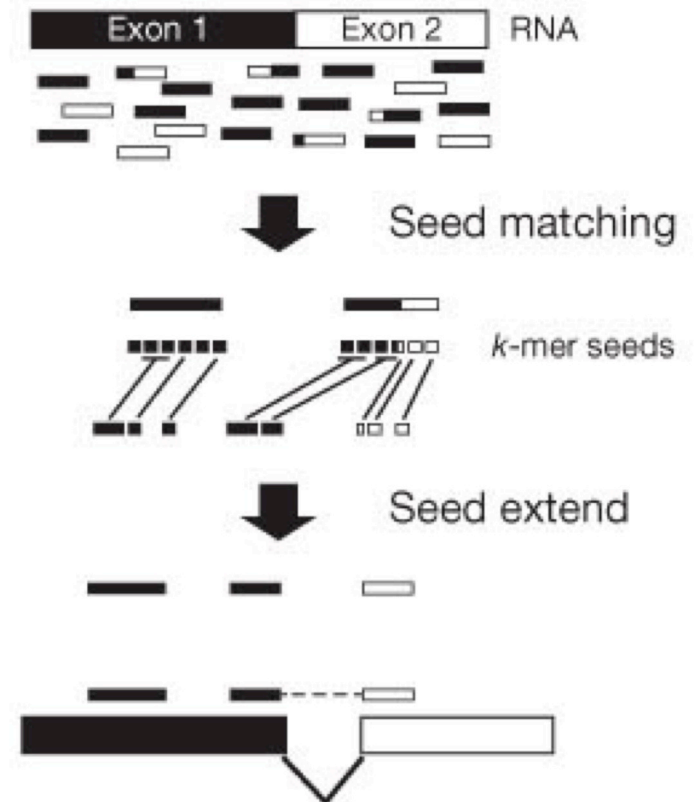


Figure :  
Garber et al, Nature Methods, 2011

# What to know about your data before mapping?

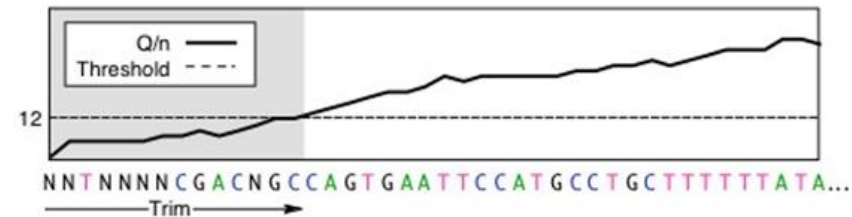
## KNOW YOUR DATA!

- Paired end? Single end?
- Traditional RNA-Seq? 3' tag ?
- Insert size estimate?



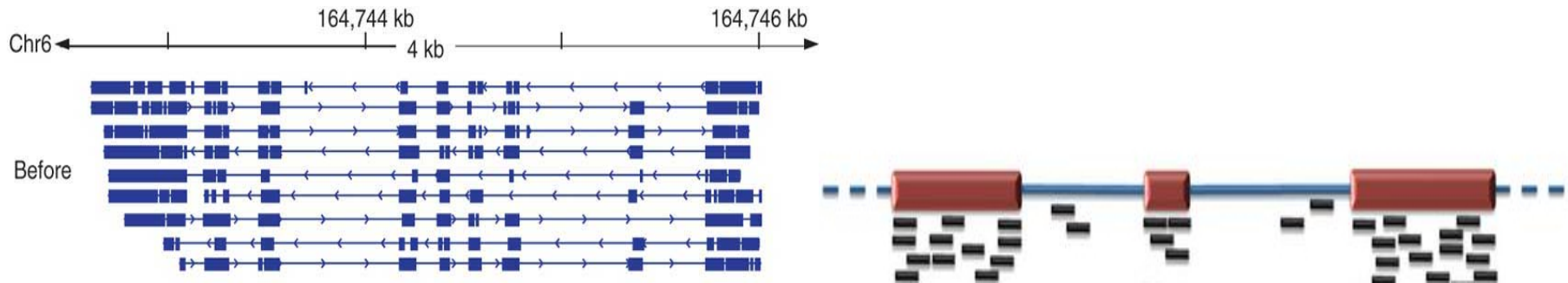
## PREPROCESSING

- Adaptor sequences trimmed?
- Primer sequences/barcodes removed?
- Poor quality regions trimmed?



# What to know about your reference before mapping?

- Mapping to genome vs transcriptome?



- Is your reference the right version?
- Does your annotation match your reference?

# What will your reference look like?

- FASTA Format

```
>gi|254160123|ref|NC_012967.1| Escherichia coli B str. REL606  
agcttttcattctgactgcaacgggcaatatgtctctgtgtggattaaaaaaagagtgctc  
tgatagcagcttctgaactggttacctgccgtgagtaaattaaattttattgacttagg  
tcactaaatactttaaccaatataggcatagcgcacagacagataaaaattacagagtac  
acaacatccatgaaacgcattagcaccaccattaccaccaccatcaccattaccacaggt  
....
```

- Using complex reference sequence names is a common problem during analysis. Might rename:

```
>REL606  
agcttttcattctgactgcaacgggcaatatgtctctgtgtggattaaaaaaagagtgctc  
tgatagcagcttctgaactggttacctgccgtgagtaaattaaattttattgacttagg
```

# What will your annotation look like?

- GFF3 Format

- seqname - The name of the sequence.
- source - The program that generated this feature.
- feature – Examples: "CDS", "start\_codon", "stop\_codon", and "exon".
- start - The starting position of the feature in the sequence.
- end - The ending position of the feature (inclusive).
- score - A score between 0 and 1000.
- strand - Valid entries include '+', '-', or '.' (for don't know/don't care).
- Frame – reading frame
- group – ID and other information about the entry

Example:

```
Rel606 refseq cds 1450 1540 500 + . Gene_id=« test_gene »
```

- Make sure the GFF3 file matches your reference fasta file.

# Mapping with BWA

- BWA is a fast short read aligner that uses the burrows-wheeler transform to perform alignment in a time and memory efficient manner.
- BWA Variants
  - For reads upto 100 bp long
    - BWA-backtrack: BWA aln/samse/sampe
  - For reads upto 1 Mbp long
    - BWA-SW
    - **BWA-MEM**: Newer! Typically faster!



# Mapping with BWA

- Create an index of your reference – `bwa index`
- Run mapping - `bwa mem`
  
- Help! I have a large number of reads. Make BWA go faster!
  - Use threading option (`bwa -t <threads>`)
  - Split one data file into smaller chunks, run multiple, parallel BWA instances, concatenate results.
    - Wait! We have a pipeline for that on lonestar – `runBWA_mem.sh` in `$BI/bin`

# Mapping with Tophat

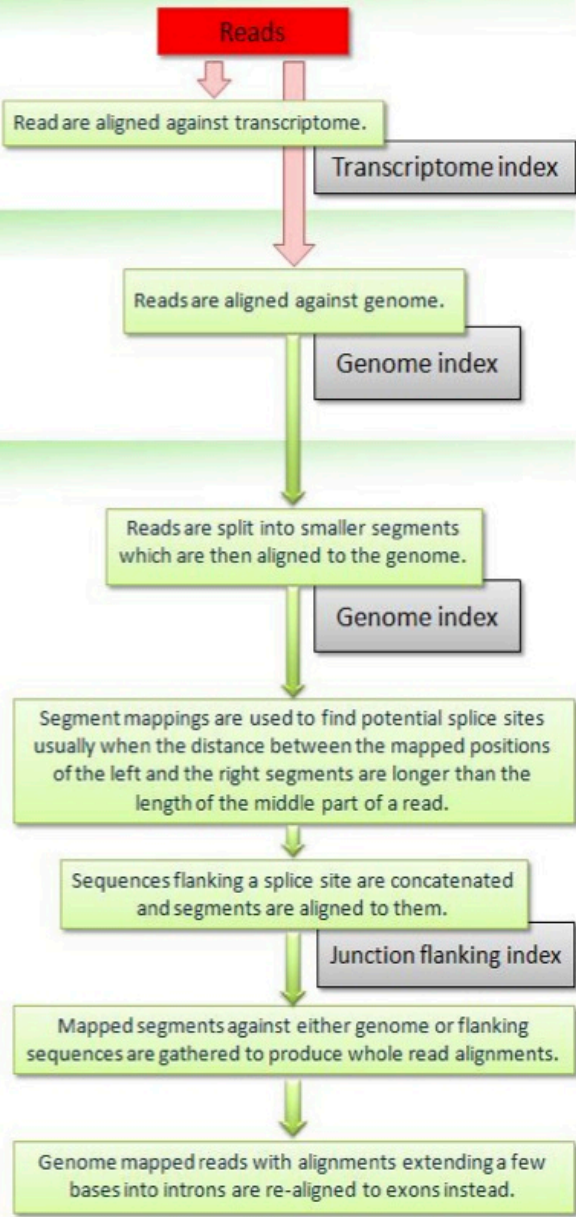
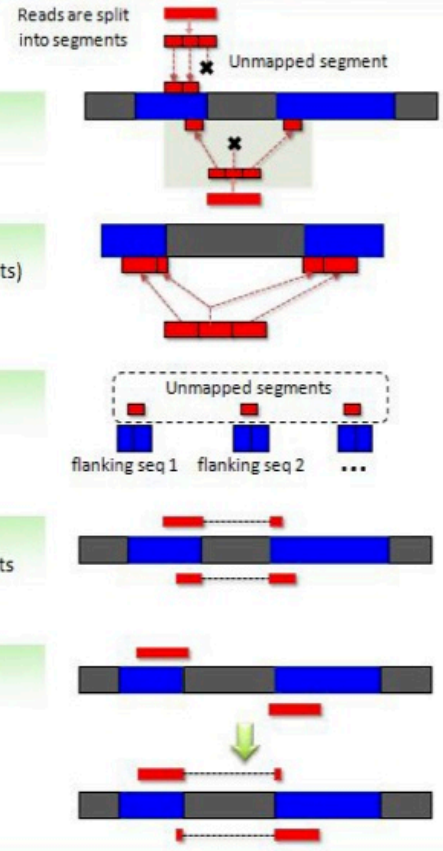
## (1) Transcriptome alignment (optional)



## (2) Genome alignment



## (3) Spliced alignment



# Mapping with Tophat

Steps:

1. Index the genome using bowtie
2. Map using tophat

Let's look at the command.

Help! I have a large number of reads. Make tophat go faster!

Use threading option (tophat -p <threads>)

Split one data file into smaller chunks, run multiple, parallel tophat instances, concatenate results.

Wait! We have a pipeline for that on lonestar – **fastTophat.sh** in \$BI/bin

# Mapping Output: SAM file format

- Alignment results generated in Sequence Alignment/Map format
- Tab delimited, with fixed columns followed by user-extendable key:data values.
- Most mappers also output unmapped reads in SAM file.
- SAMTOOLS – toolkit to manipulate, parse SAM files.

# Mapping Output: SAM File Format

SAM fixed fields:

<http://samtools.sourceforge.net/>

Col	Field	Type	Regexp/Range	Brief description
1	QNAME	String	[!-?A-~]{1,255}	Query template NAME
2	FLAG	Int	[0,2 <sup>16</sup> -1]	bitwise FLAG
3	RNAME	String	\*  [!-( )+-<>-~] [!-~]*	Reference sequence NAME
4	POS	Int	[0,2 <sup>29</sup> -1]	1-based leftmost mapping POSition
5	MAPQ	Int	[0,2 <sup>8</sup> -1]	MAPping Quality
6	CIGAR	String	\*  ([0-9]+[MIDNSHPX=])+	CIGAR string
7	RNEXT	String	\* =  [!-( )+-<>-~] [!-~]*	Ref. name of the mate/next segment
8	PNEXT	Int	[0,2 <sup>29</sup> -1]	Position of the mate/next segment
9	TLEN	Int	[-2 <sup>29</sup> +1,2 <sup>29</sup> -1]	observed Template LENgth
10	SEQ	String	\*  [A-Za-z=.]+	segment SEQUENCE
11	QUAL	String	[!-~]+	ASCII of Phred-scaled base QUALity+33

```
SRR030257.264529    99  NC_012967   1521    29  34M2S   =    1564
79  CTGGCCATTATCTCGGTGGTAGGACATGGCATGCC
AAAAAA;AA;AAAAAA??A%.;?&'3735',()0*,
XT:A:M NM:i:3 SM:i:29 AM:i:29 XM:i:3 XO:i:0 XG:i:0 MD:Z:23T0G4T4
```

# Mapping Output: Mapping Quality

- Mapping quality is the probability that a read is aligned to the wrong place.

$$p = 10^{**(-q/10)}$$

- BWA mapping quality calculated by considering:
  - Repeat structure of reference
  - Read base quality
  - Read alignment quality (mismatches etc)
  - Number of mappings

# Mapping Output: CIGAR score

Ref CTGGCCATTATCTC--GGTGGTAGGACATGGCATGCCC  
Read aaATGTCGCGGTG.TAGGAaggatcc



2S5M2I4M1D4M6S

Op	BAM	Description
M	0	alignment match (can be a sequence match or mismatch)
I	1	insertion to the reference
D	2	deletion from the reference
* N	3	skipped region from the reference
S	4	soft clipping (clipped sequences present in SEQ)
* H	5	hard clipping (clipped sequences NOT present in SEQ)
* P	6	padding (silent deletion from padded reference)
* =	7	sequence match
* X	8	sequence mismatch

\*Rarer / newer

CIGAR = "Concise Idiosyncratic Gapped Alignment Report"

# Mapping Output: BAM format

- SAM files are converted to BAM format through SAMTOOLS command:
  - `samtools view -b -S samfile > bamfile`
- BAM file is binary format.
- BAM file is compressed.
- BAM files are usually what you need for post mapping analysis and visualization.



# Assess Mapping Results - Samtools

- For parsing and manipulating mapping output files in SAM and BAM formats.
  - Sorting mapping output files
  - Merging multiple mapping output files
  - Converting from SAM to BAM and vice versa
  - Retrieving reads based on different criteria: reads mapping to a particular region, unmapped reads etc
  - Collecting statistics about your mapping results

# Assess Mapping Results - Samtools

1. Convert SAM file to BAM format  
samtools view
2. Sort and index newly created BAM file  
samtools sort  
samtools index
3. Mapping Statistics  
samtools flagstat  
samtools idxstats

# Assess Mapping Results - RNASEQC

## Transcript-associated Reads

Sample	Note	Intragenic Rate	Exonic Rate	Intronic Rate	Intergenic Rate	Expression Profiling Efficiency	Transcripts Detected	Genes Detected
K-562	v1.0 dUTPICell Line	0.897	0.538	0.359	0.103	0.411	79,585	18,663
GTEX-N7MS-2526	v1.0 dUTPIBrain 9.638445	0.888	0.446	0.442	0.111	0.327	87,101	20,970
GTEX-N7MT-0126	v1.0 dUTPILung 9.074045	0.907	0.464	0.443	0.092	0.276	90,362	21,217

## Coverage Metrics for Bottom 1000 Expressed Transcripts

The metrics in this table are calculated across the transcripts that were determined to have the highest expression levels.

Sample	Note	Mean Per Base Cov.	Mean CV	No. Covered 5'	5' 100Base Norm	No. Covered 3'	3' 100Base Norm	Num. Gaps	Cumul. Gap Length	Gap %
<a href="#">K-562</a>	v1.0 dUTPIFibroblast	7.17	0.84	739	0.90	791	0.833	2204	230166	15.6
<a href="#">GTEX-N7MS-2526</a>	v1.0 dUTPIBrain 9.638445	5.35	0.75	742	0.68	836	0.954	2403	207728	13.8
<a href="#">GTEX-N7MT-0126</a>	v1.0 dUTPILung 9.074045	4.60	0.77	713	0.69	788	0.843	2792	227526	14.7

It is important to note that these values are restricted to the bottom 1000 expressed transcripts. 5' and 3' values are per-base coverage averaged across all top transcripts. 5' and 3' ends are 100 base pairs. Gap % is the total cumulative gap length divided by the total cumulative transcript lengths.

## Coverage Metrics for Middle 1000 Expressed Transcripts

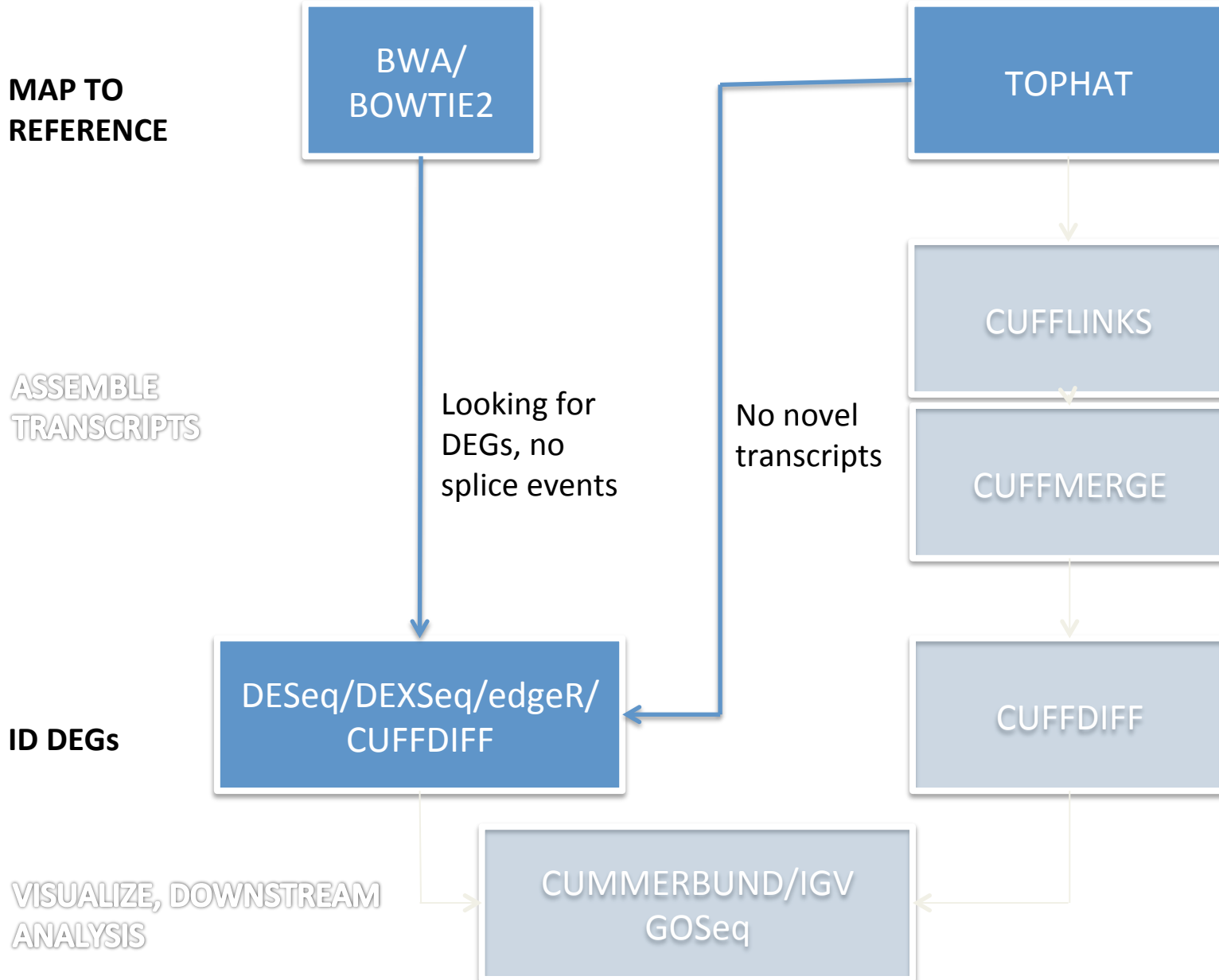
The metrics in this table are calculated across the transcripts that were determined to have the highest expression levels.

Sample	Note	Mean Per Base Cov.	Mean CV	No. Covered 5'	5' 100Base Norm	No. Covered 3'	3' 100Base Norm	Num. Gaps	Cumul. Gap Length	Gap %
<a href="#">K-562</a>	v1.0 dUTPIFibroblast	24.42	0.62	863	0.79	890	0.787	1045	83828	4.3
<a href="#">GTEX-N7MS-2526</a>	v1.0 dUTPIBrain 9.638445	14.61	0.61	854	0.59	943	0.949	972	69905	3.5
<a href="#">GTEX-N7MT-0126</a>	v1.0 dUTPILung 9.074045	11.90	0.63	852	0.63	877	0.841	1316	90803	4.5

# Mapping Summary

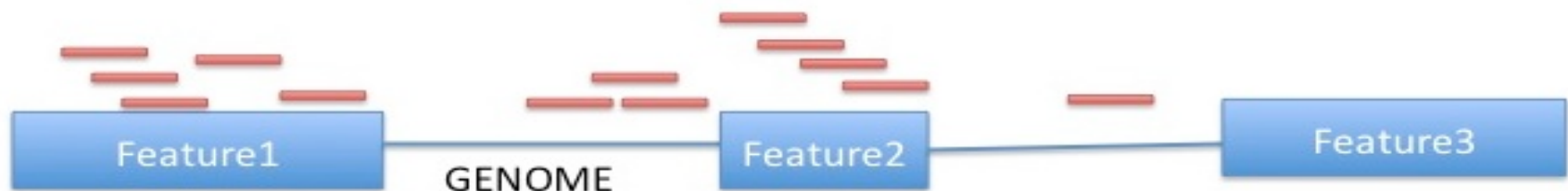
- Unspliced mappers (BWA, bowtie2) ok for quantifying known exons and junctions.
- Spliced mappers (tophat) can, in addition identify and quantify splice variants involving unknown exons.
- Samtools can be used to gather basic mapping statistics, RNASEQC for specific rna specific statistics

# STEP 4 and 5: Quantify Expression and ID DEGs



# STEP 4: Quantify Expression

- Bedtools
  - **Bedtools multicov** : Takes a feature file (GFF) and counts how many reads in the mapped output file (BAM) overlap the features.
  - Remember that the chromosome names in your gff file should match the chromosome names in the reference fasta file used in the mapping step.



# STEP 4 : Quantify Expression

	union	intersection_strict	intersection_nonempty
	gene_A	gene_A	gene_A
	gene_A	no_feature	gene_A
	gene_A	no_feature	gene_A
	gene_A	gene_A	gene_A
	gene_A	gene_A	gene_A
	ambiguous	gene_A	gene_A
	ambiguous	ambiguous	ambiguous

HTSeq –

- Gives you fine grained control over how to count genes, especially when a read overlaps more than one gene/feature.

# STEP 5: ID Differentially Expressed Genes

- Normalize of gene counts
- Represent the gene counts by a distribution that defines the relation between mean and variance (dispersion).
- Perform statistical test to compare this distribution between conditions.
- Provide fold change, P-value information, false discovery rate for each gene.



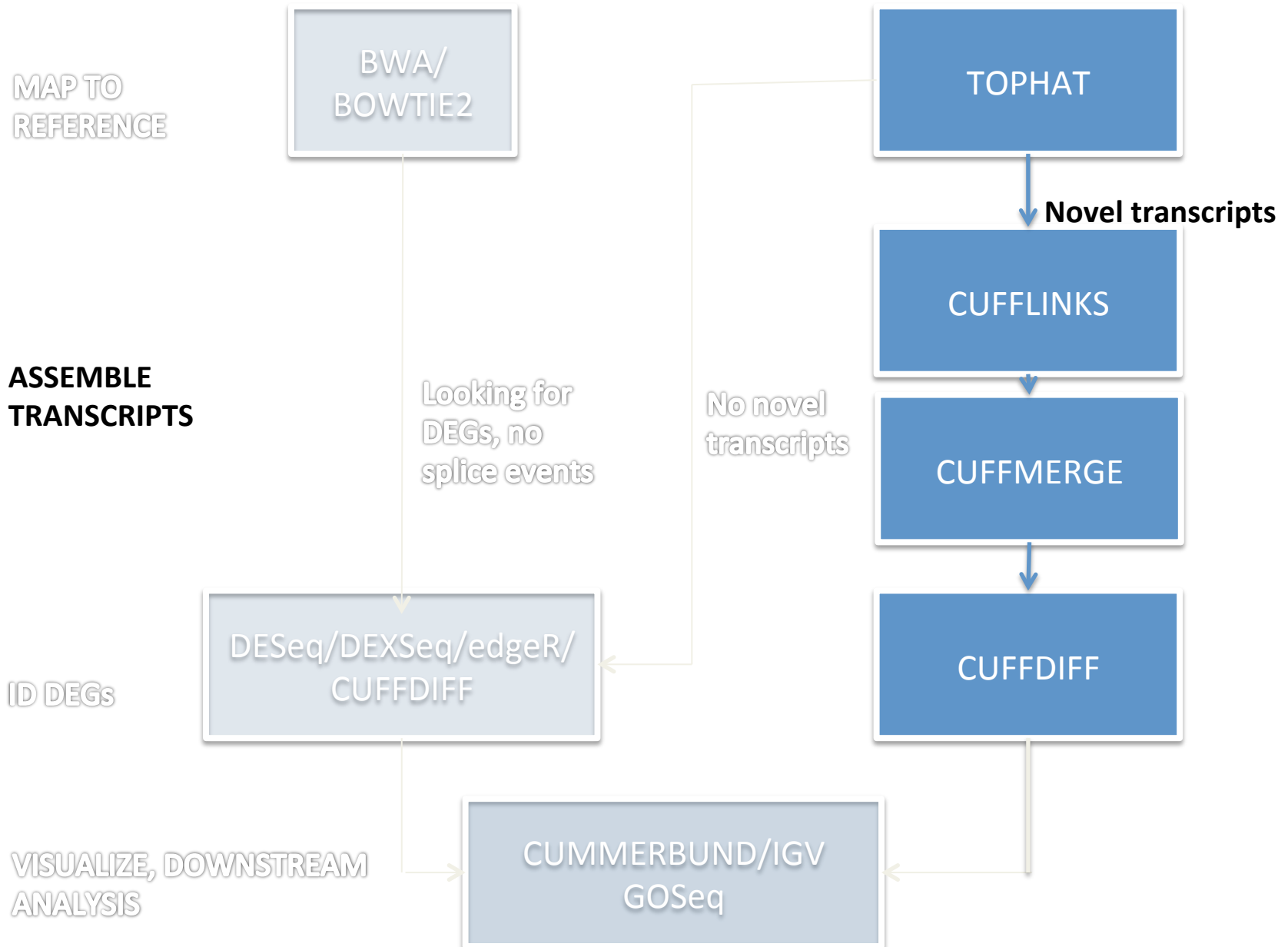
# STEP 5: ID Differentially Expressed Genes

	DESeq	edgeR	DEXSeq	Cuffdiff
Normalization	Median scaling size factor	Median scaling size factor/TMM	Median scaling size factor	FPKM
Distribution	Negative binomial	Negative binomial	Negative binomial	Negative binomial
DE Test	Negative binomial test	Fisher exact test	Modified T test	T test
Advantages	Straightforward, fast, has a method to work on data with no replicates	Straightforward, fast, good with small number of replicates.  Can handle comparisons across multiple conditions.	Good for identifying exon-usage changes	Good for identifying isoform-level changes, splicing changes, promotor changes.  Not as straightforward

# STEP 5: ID Differentially Expressed Genes

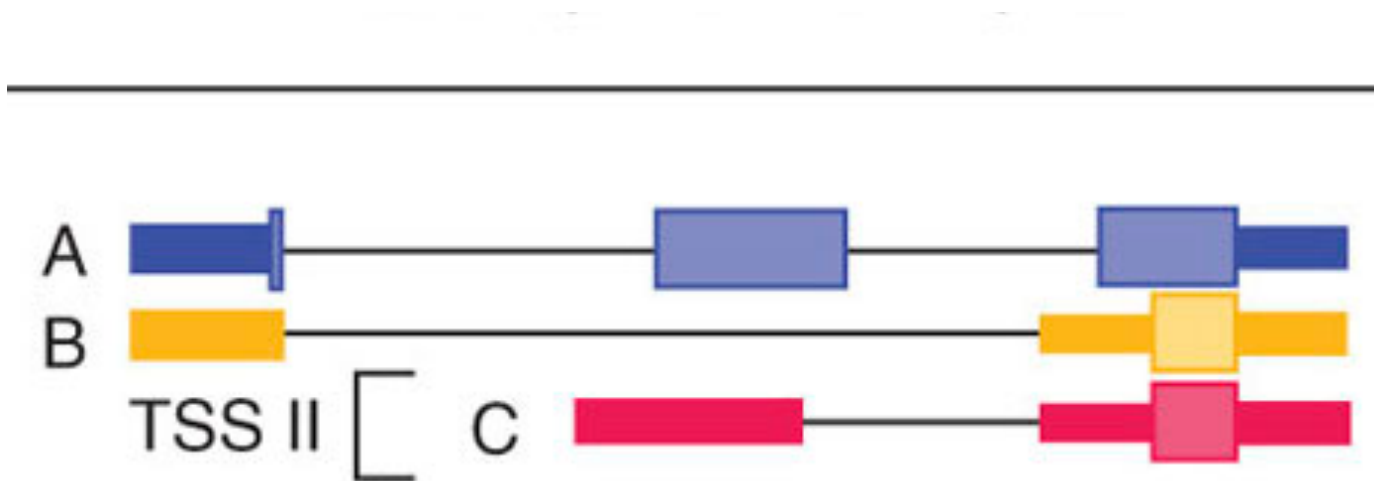
- **DESeq** Input: **RAW** count data , with each column representing a biological replicate/condition.
- DESeq R commands available at:  
<https://wikis.utexas.edu/display/bioiteam/Testing+for+differential+expression>
- Let's look at bedtools and DESeq results for now.
- **Cuffdiff** covered further down the line.

# STEP 3: Assemble Transcripts



# What is a gene? What is a transcript?

A gene can have multiple transcripts!



- We want to identify all these transcripts, whether annotated or not.

# What? Why? How?

- What are we trying to do?
  - Novel gene and transcript discovery!!
- Why?
  - If our sample condition has novel transcripts not found in annotated transcriptome, that tells us something about what may make our sample unique.
  - If across two conditions, novel transcripts are differentially expressed, those transcripts may be the underlying reasons for the biological differences between our two conditions.
- How?
  - Using the tuxedo pipeline.

# TUXEDO PIPELINE

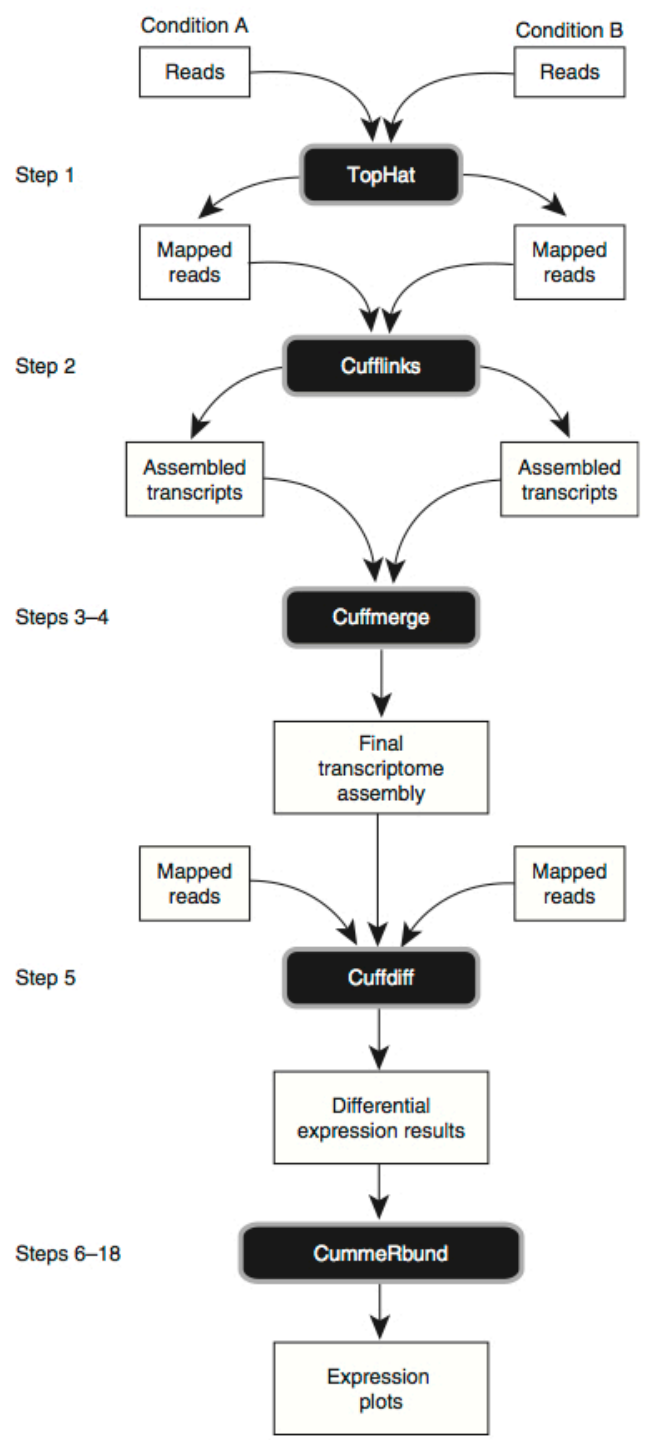
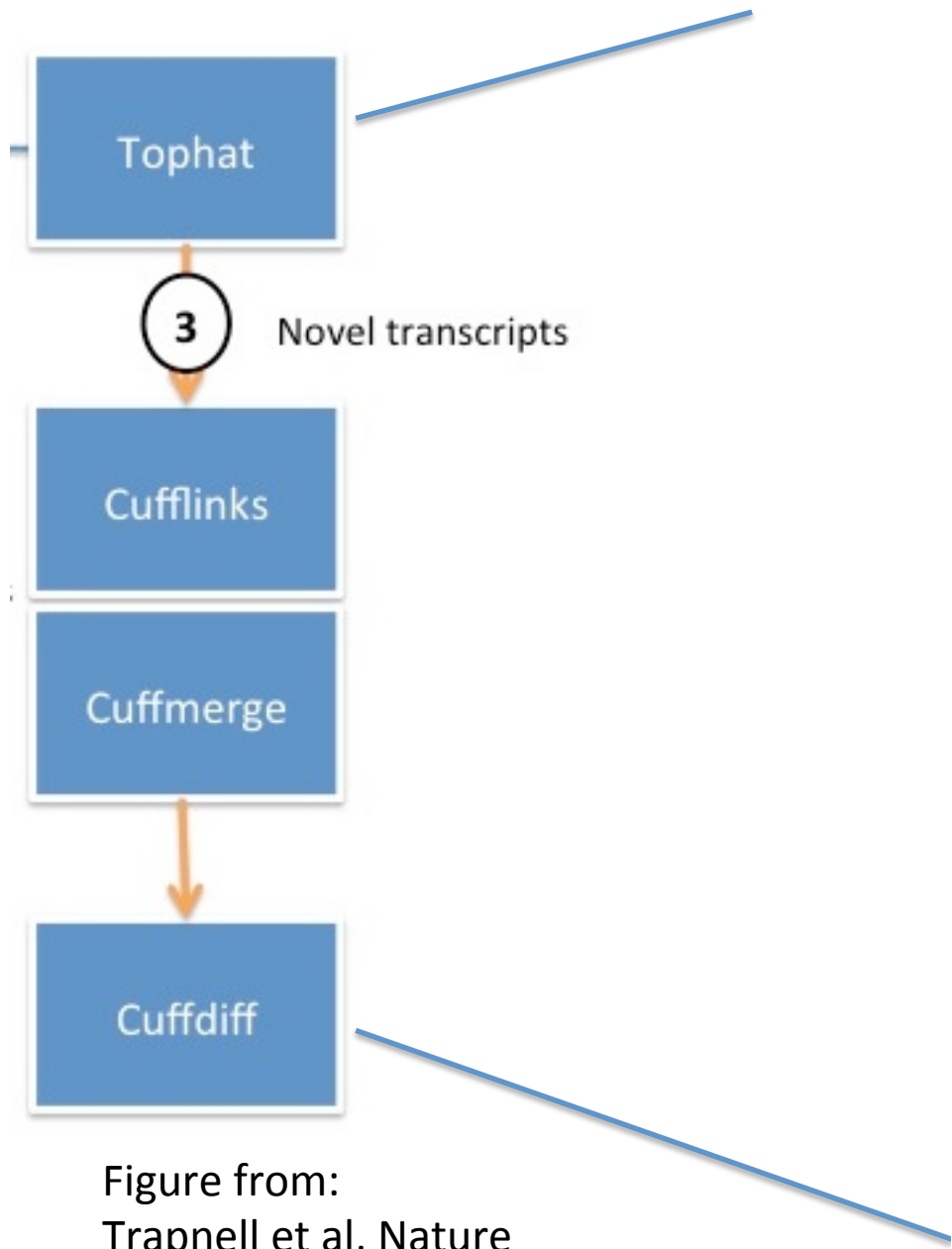


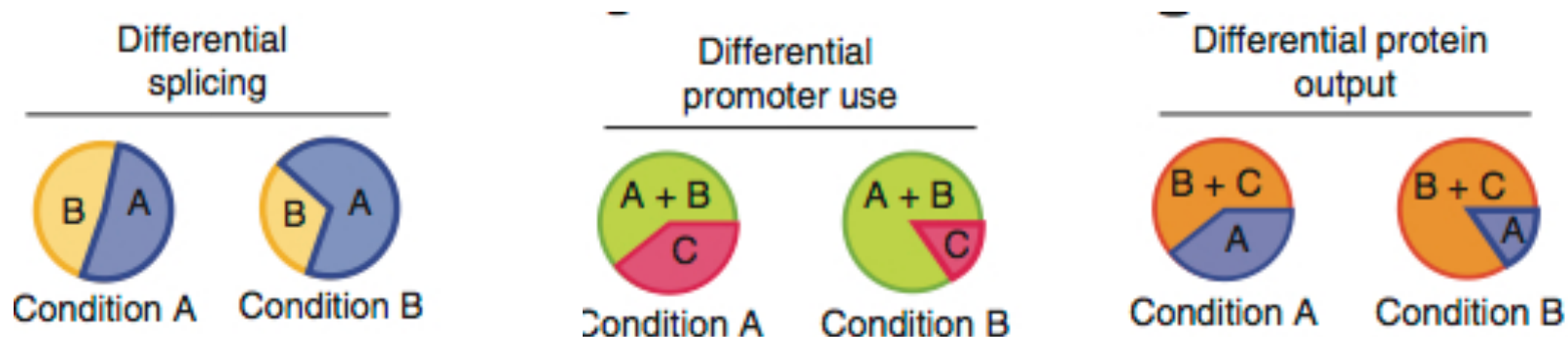
Figure from:  
Trapnell et al, Nature  
protocols, 2012.

# What do we get at the end of running this pipeline?

A view of how the transcriptome is different between condition C1 and condition C2

- Both in terms of annotated genes and transcripts.
- And novel genes and transcripts

Differential gene expression and so much more...



# A. TOPHAT

Tophat maps your data to your reference in a splice-aware manner, that will also identify junctions. We've already looked at to run it.

Output: Mapped output in bam format



# B. CUFFLINKS

Reconstructs/assembles transcript for each sample.

## **Why is transcript assembly hard?**

Difficult to tell which read came from which transcript

- - Many short reads, many transcripts!
- Transcripts are expressed in different amounts. So, coverage of reads can be vastly different.
- Reads can come from mature mRNA (exons only) and precursor RNA (containing partial introns).

**Table 1** | Selected list of RNA-seq analysis programs

Class	Category	Package	Notes	Uses	Input
<b>Read mapping</b>					
Unspliced aligners <sup>a</sup>	Seed methods	Short-read mapping package (SHRiMP) <sup>41</sup> Stampy <sup>39</sup>	Smith-Waterman extension Probabilistic model	Aligning reads to a reference transcriptome	Reads and reference transcriptome
	Burrows-Wheeler transform methods	Bowtie <sup>43</sup> BWA <sup>44</sup>	Incorporates quality scores		
Spliced aligners	Exon-first methods	MapSplice <sup>52</sup> SpliceMap <sup>50</sup> TopHat <sup>51</sup>	Works with multiple unspliced aligners Uses Bowtie alignments	Aligning reads to a reference genome. Allows for the identification of novel splice junctions	Reads and reference genome
	Seed-extend methods	GSNAP <sup>53</sup> QPALMA <sup>54</sup>	Can use SNP databases Smith-Waterman for large gaps		
<b>Transcriptome reconstruction</b>					
Genome-guided reconstruction	Exon identification	G.Mor.Se	Assembles exons	Identifying novel transcripts using a known reference genome	Alignments to reference genome
	Genome-guided assembly	Scripture <sup>28</sup> Cufflinks <sup>29</sup>	Reports all isoforms Reports a minimal set of isoforms		
Genome-independent reconstruction	Genome-independent assembly	Velvet <sup>61</sup> TransABySS <sup>56</sup>	Reports all isoforms	Identifying novel genes and transcript isoforms without a known reference genome	Reads
<b>Expression quantification</b>					
Expression quantification	Gene quantification	Alexa-seq <sup>47</sup>	Quantifies using differentially included exons	Quantifying gene expression	Reads and transcript models
		Enhanced read analysis of gene expression (ERANGE) <sup>20</sup> Normalization by expected uniquely mappable area (NEUMA) <sup>82</sup>	Quantifies using union of exons Quantifies using unique reads		
	Isoform quantification	Cufflinks <sup>29</sup> MISO <sup>33</sup> RNA-seq by expectation maximization (RSEM) <sup>69</sup>	Maximum likelihood estimation of relative isoform expression	Quantifying transcript isoform expression levels	Read alignments to isoforms
Differential expression		Cuffdiff <sup>29</sup> DegSeq <sup>79</sup> EdgeR <sup>77</sup>	Uses isoform levels in analysis Uses a normal distribution	Identifying differentially expressed genes or transcript isoforms	Read alignments and transcript models
		Differential Expression analysis of count data (DESeq) <sup>78</sup> Myrna <sup>75</sup>	Cloud-based permutation method		

Figure :  
Garber et al, Nature Methods, 2011

**Most commonly used, if you have a genome.**

**Less resource-intensive**

We'll call this coverage islands method

<b>Transcriptome reconstruction</b>				
Genome-guided reconstruction	Exon identification	G.Mor.Se	Assembles exons	Identifying novel transcripts using a known reference genome
	Genome-guided assembly	Scripture <sup>28</sup> Cufflinks <sup>29</sup>	Reports all isoforms Reports a minimal set of isoforms	
Genome-independent reconstruction	Genome-independent assembly	Velvet <sup>61</sup> TransABySS <sup>56</sup>	Reports all isoforms	Identifying novel genes and transcript isoforms without a known reference genome

**If you don't have a genome.**

**If you believe your sample has major rearrangements**

**More CPU and RAM intensive**

Figure :

Garber et al, Nature Methods, 2011

# Genome guided transcript assembly

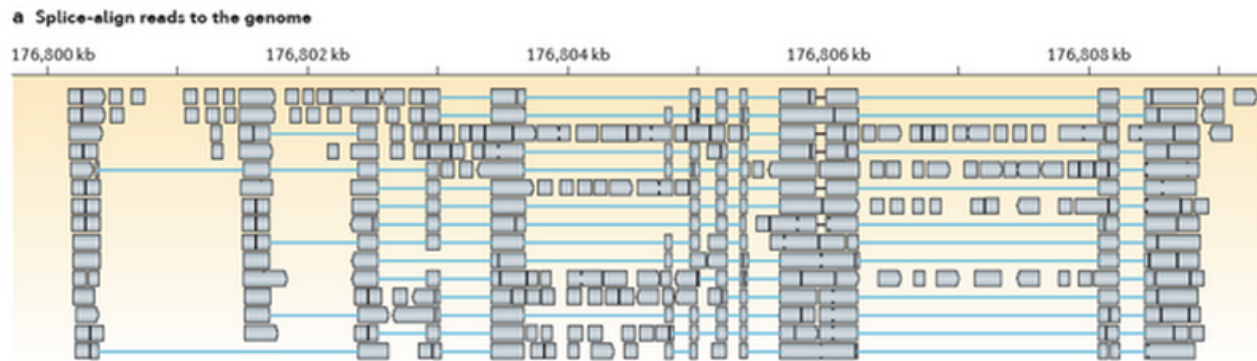
## Different assembly methods

- **Coverage islands**

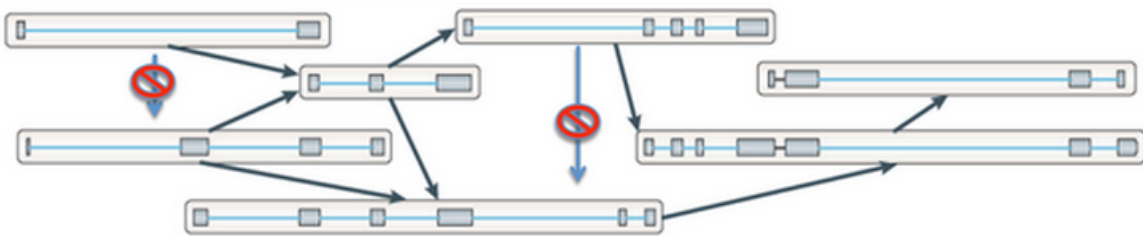
- ID putative exons by looking for coverage islands.
- Older method, were meant for shorter read lengths.
- G.MorSe

- **Exon first approach**

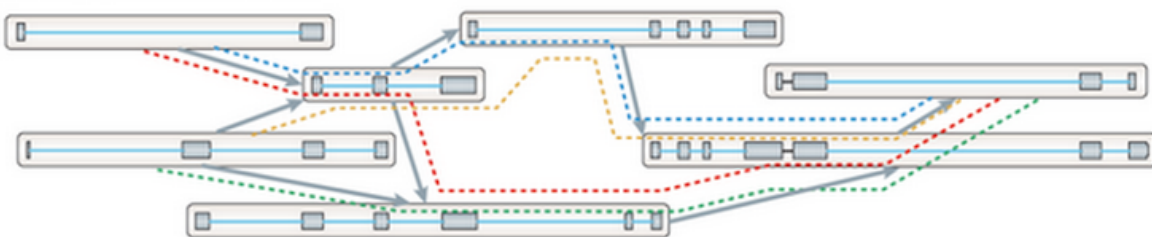
- Directly uses mappings of spliced reads to reconstruct transcriptome.
- Uses graph topology.
- **Cufflinks (part of tuxedo suite)**, scripture



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



**c Traverse the graph to assemble variants**



**d Assembled isoforms**

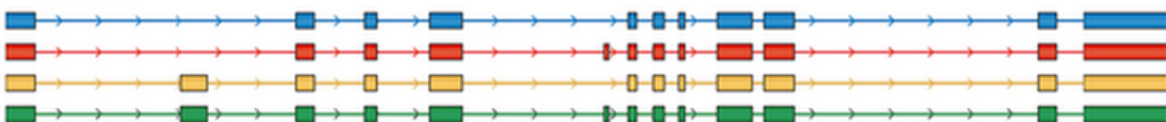
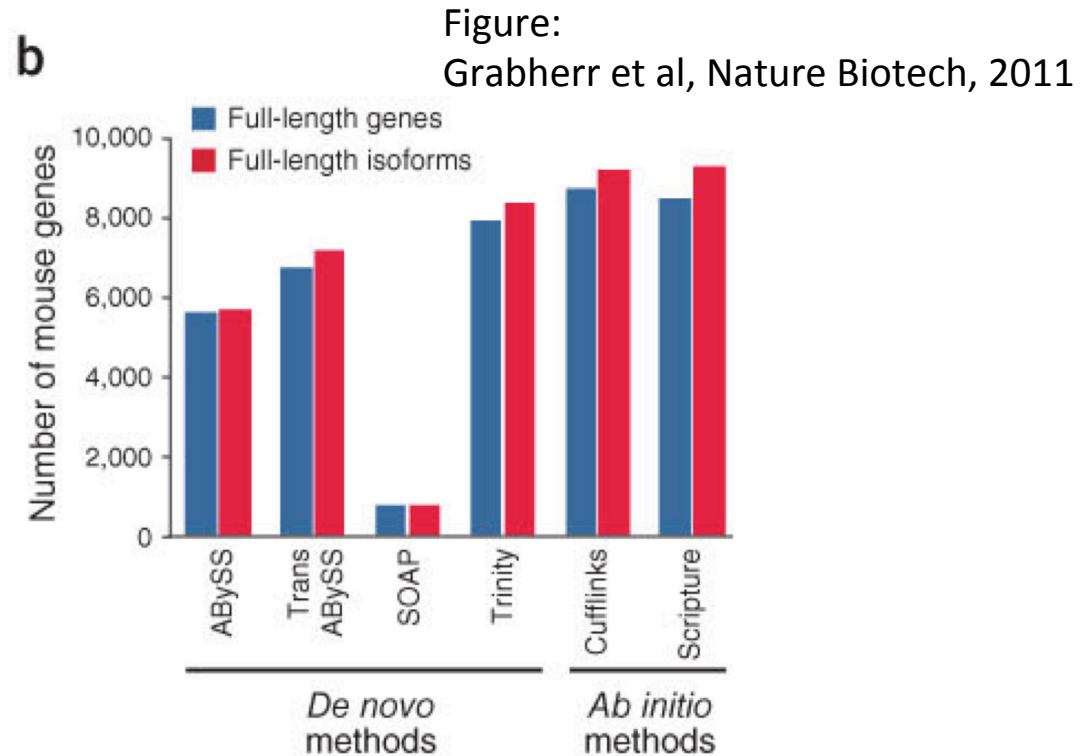
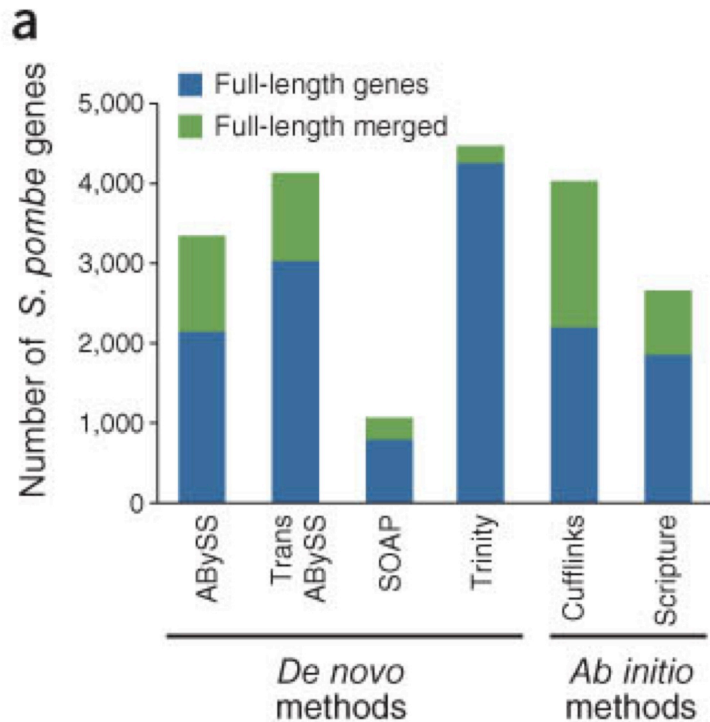


Figure :  
[http://sourceforge.net/projects/trinityrnaseq/files/misc/RNASEQ\\_WORKSHOP/rnaseq\\_workshop\\_slides.pdf](http://sourceforge.net/projects/trinityrnaseq/files/misc/RNASEQ_WORKSHOP/rnaseq_workshop_slides.pdf)

# How do these tools compare?



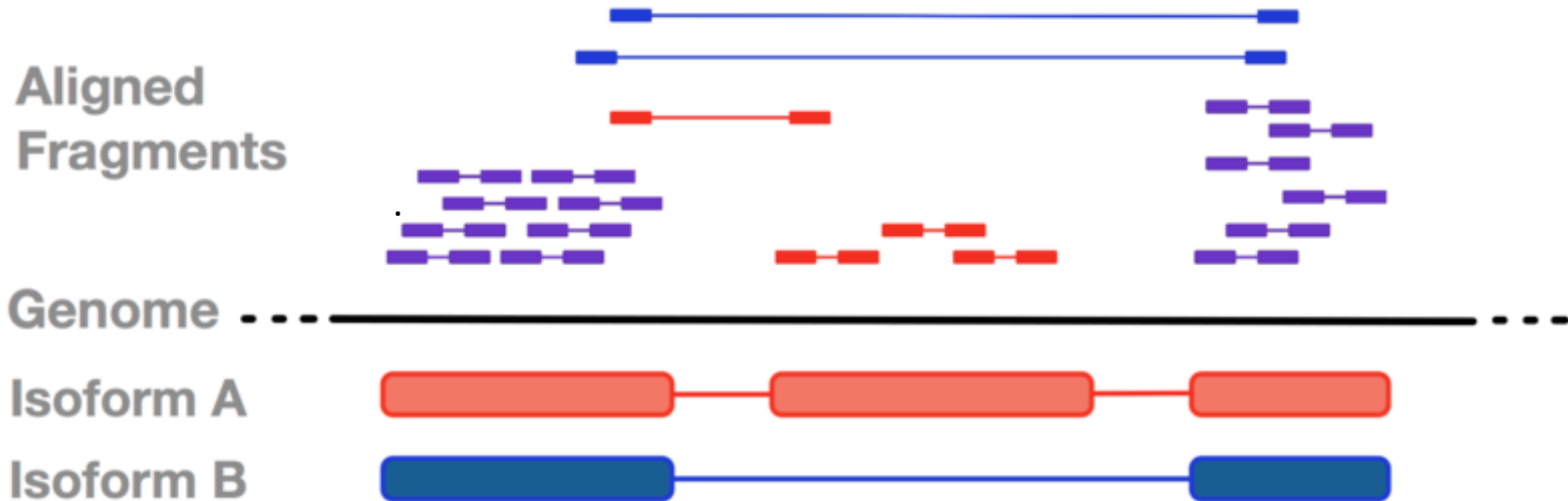
Program combination	vs. orthology annotation		vs. EST annotation	
	Base-level accuracy (%) <sup>1</sup>	Confirmed junctions (%) <sup>1</sup>	Base-level accuracy (%) <sup>1</sup>	Confirmed junctions (%) <sup>1</sup>
TopHat + Cufflinks	83.9	75.8	68.9	63.0
GSNAP + Cufflinks	79.4	71.2	65.7	58.4
GSNAP + Cufflinks (subsample <sup>2</sup> )	80.3	72.7	60.2	66.3
TopHat + Scripture	70.3	67.9	60.8	62.5

Figure:  
Palmieri et al,  
PLOS One, 2012

<sup>1</sup>Base level accuracy and percentage of confirmed junctions with different combinations of mapper and assembler on the sample ps94 males compared to the orthology annotation and the EST annotation (<sup>2</sup>based on 48 M reads).

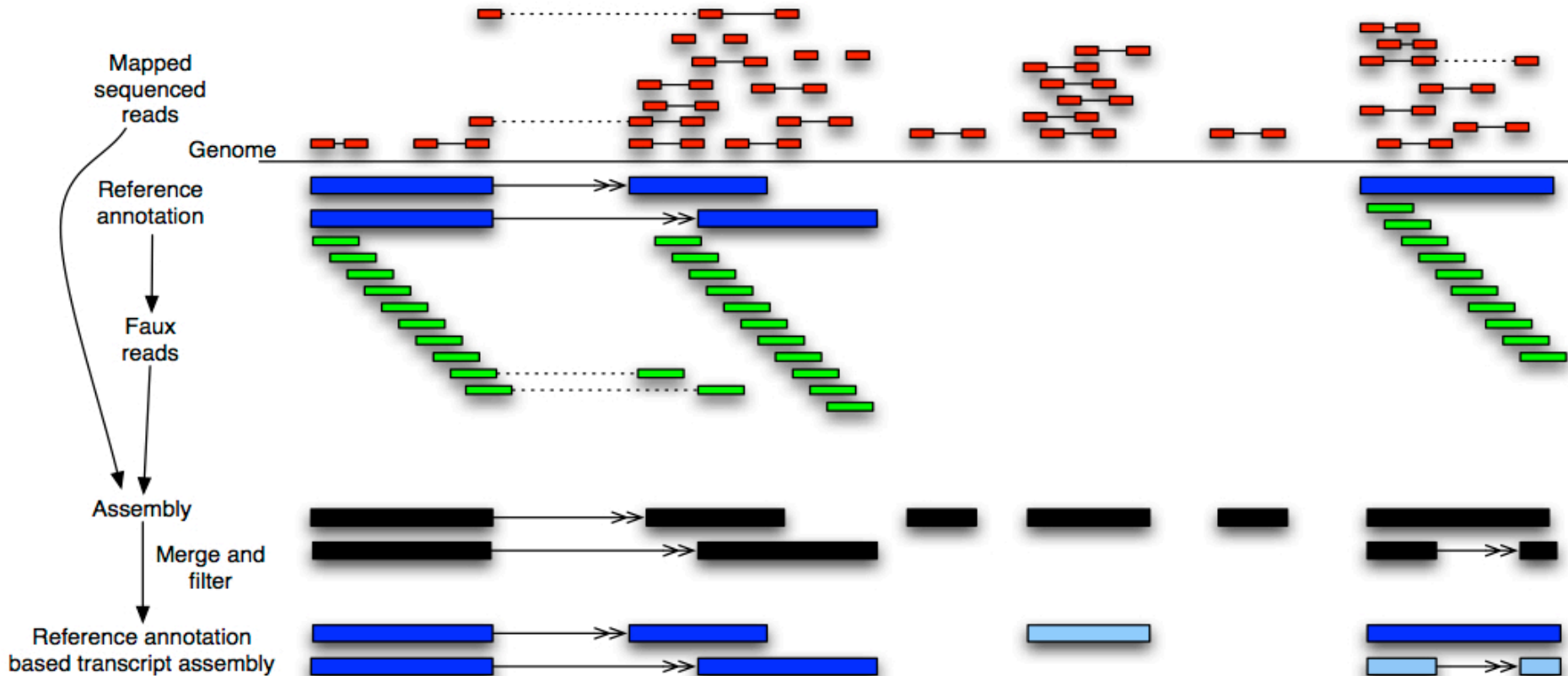
# How does Cufflinks do transcript assembly

Exon first method!



# RABT

- Reference annotation based transcript assembly (RABT)
  - Uses existing annotation to guide assembly of transcripts.





# After assembly

- Calculates abundance for these assembled transcripts.
- Normalized using FPKM (Fragments Per Kilobase of Exon Per Million) (variation of RPKM)
  - RPKM normalizes for **transcript length variations** and **sequencing depth**.
  - $RPKM = (\text{No. of Mapped reads} * 10^9) / (\text{length of transcript} * \text{total no. of reads})$
  - FPKM just exchanges reads with fragments.

# General syntax for cufflinks command

```
cufflinks [options] <accepted_hits.bam>
```

Some of the important options:

- p/--num-threads

- G/--GTF (quantify only annotated transcripts)

- g/--GTF-guide (both annotated and novel transcripts)

- b/--frag-bias-correct

- u/--multi-read-correct

# General syntax for cufflinks command

-b/--frag-bias-correct

When quantifying abundance, corrects for sequence-specific bias at the ends of reads by 'learning' from the data.

-u/--multi-read-correct

By default, if a read maps to 2 genes it will count as 50% (half a read) towards each gene.

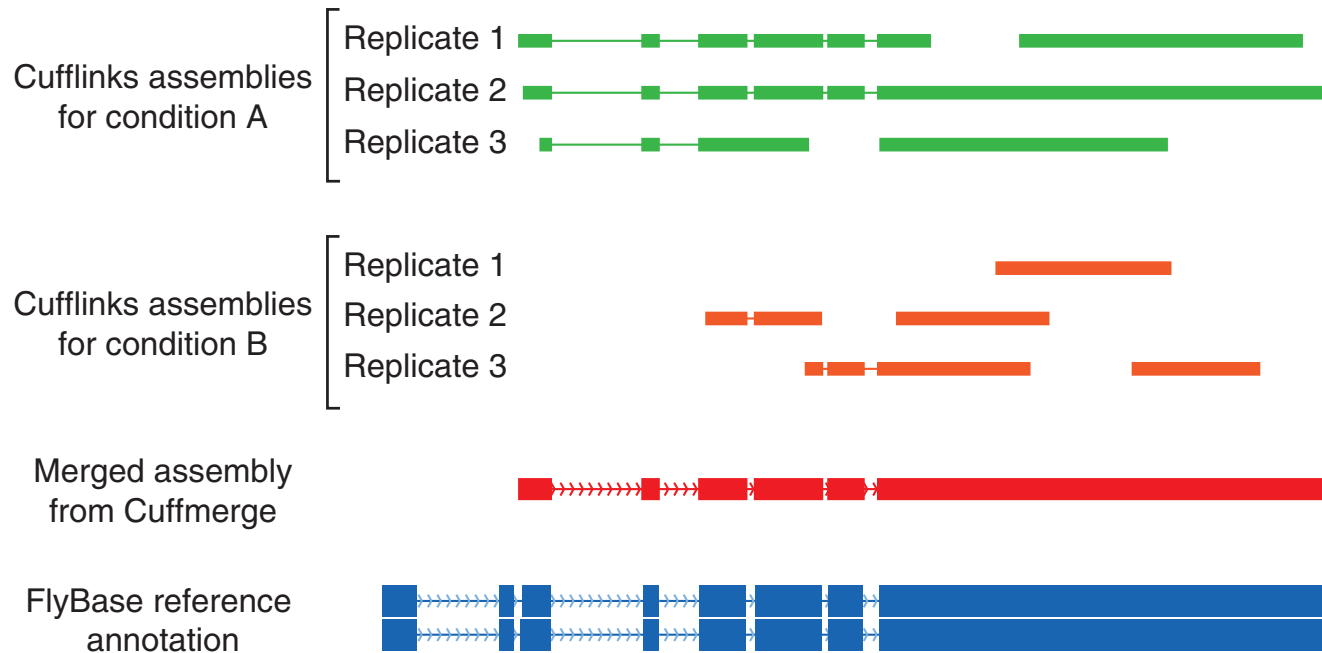
With this flag, it handles this question in a more fine-tuned manner.

# Let's look at some results from a cufflinks transcript assembly

- Input:
  - Tophat mapped results (bam files)
  - Transcriptome annotation (genes.gtf)
- Let's look at the [wiki](#) and the output files.

# C. CUFFMERGE

- Cuffmerge is used to merge all the transcripts that cufflinks assembled into one file.



# C. CUFFMERGE

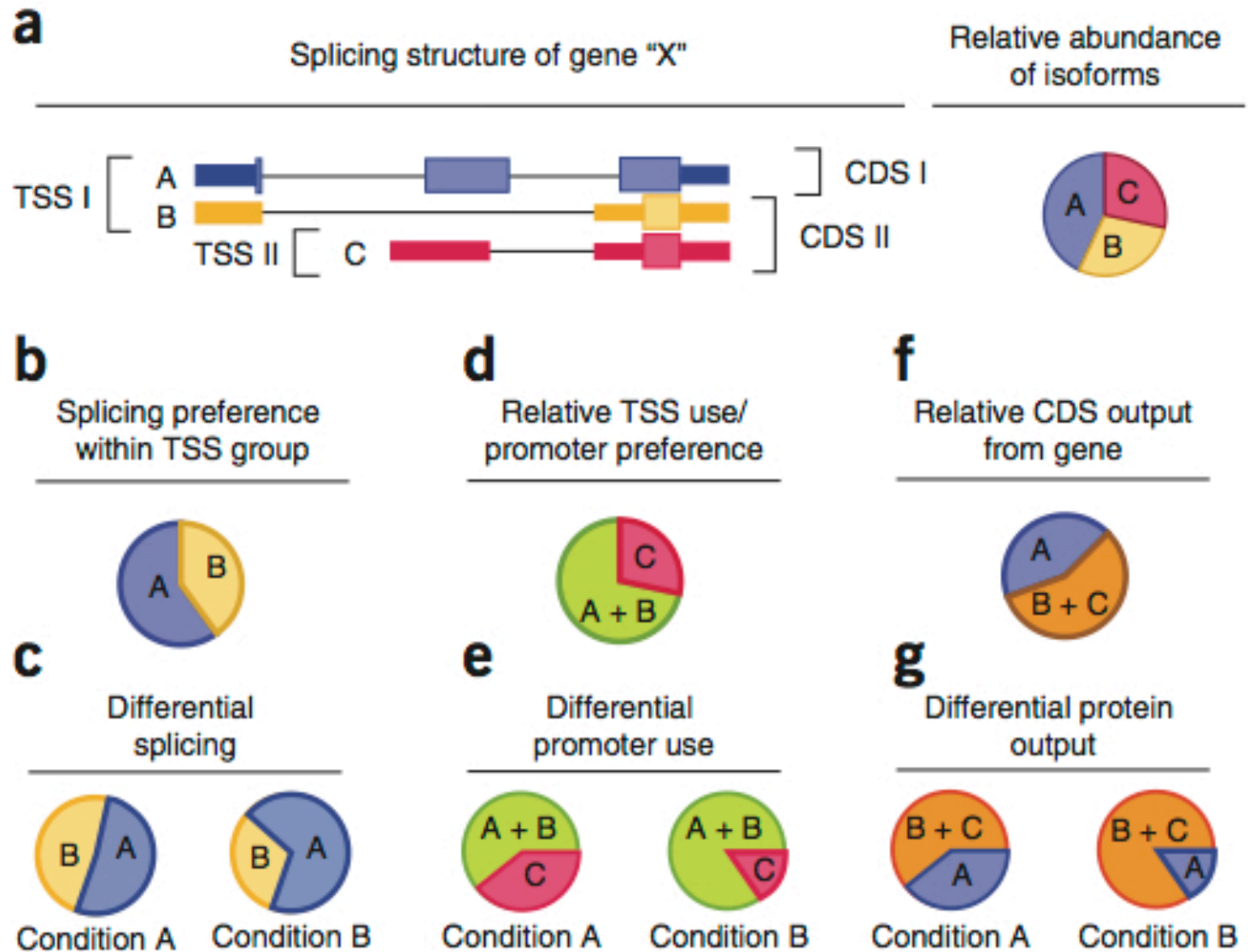
- Input: All cufflinks assembly files (in gtf format)
- Output: merged.gtf
  - Your very own gtf file, containing all the transcripts found in your samples (both novel and otherwise).
  - Also information about how the novel transcripts relate to the known transcripts
- **SWITCH TO THE WIKI** for instructions on viewing these results

# D. CUFFDIFF

- Calculates differential expression!
- Input:
  - Our newly created merged.gtf file or A gtf file we downloaded (genes.gtf)
  - Mapped bam files
- Calculates difference in isoform-level expression among conditions.
- If the chance of seeing this difference is small enough under the chosen statistical model, it is deemed significantly differentially expressed.

# D. CUFFDIFF

Figure from: Differential analysis of gene regulation at transcript resolution with rNA-seq, Trapnell et al, Nature Biotechnology, 2013





## D. CUFFDIFF

- SWITCH TO THE WIKI for instructions on viewing these results

# DESeq/edgeR output vs Tuxedo pipeline output

- We generated differential expressed genes using DESeq too. So, why the big fuss?
  - They were all from annotated genes. So, they all has flybase ids.
  - Now our output has genes with ids ‘CUFF...’ - they are novel.
  - In addition to differential gene expression, we also have results for differential regulation.
  - We also have results telling us where our novel transcripts are with respect to the annotated ones.

# Limitations of the Tuxedo Pipeline

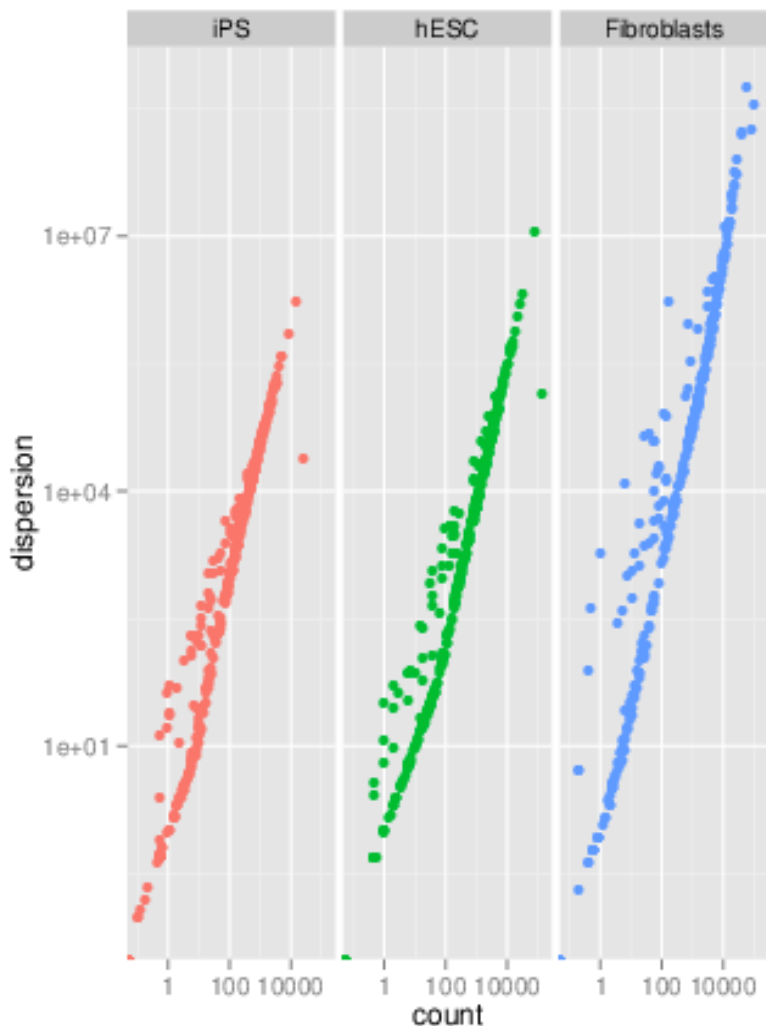
- A Reference is needed.
- Kind of a black box.
- Not quick.
  - Step 1, align the RNA-seq reads to the genome: ~6 h
  - Steps 2–4, assemble expressed genes and transcripts: ~6 h
  - Step 5, identify differentially expressed genes and transcripts:~6 h

If you don't have a genome:

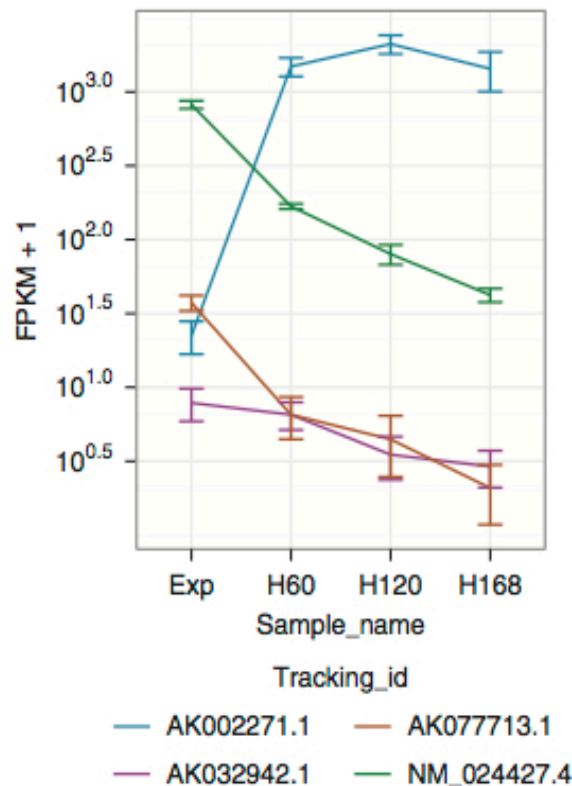
- De novo transcriptome assembly using trinity.
- Map your data to this to calculate gene expression changes.

# STEP 6: Visualize and Perform Other Downstream Analysis

- Visualize using Cummerbund



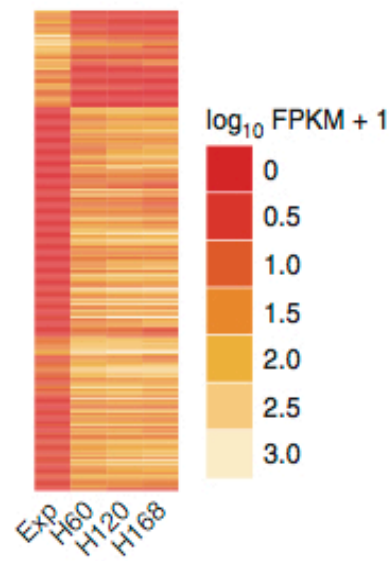
**a** `expressionPlot(isoforms(tpn1), logMode=T)`



**b**

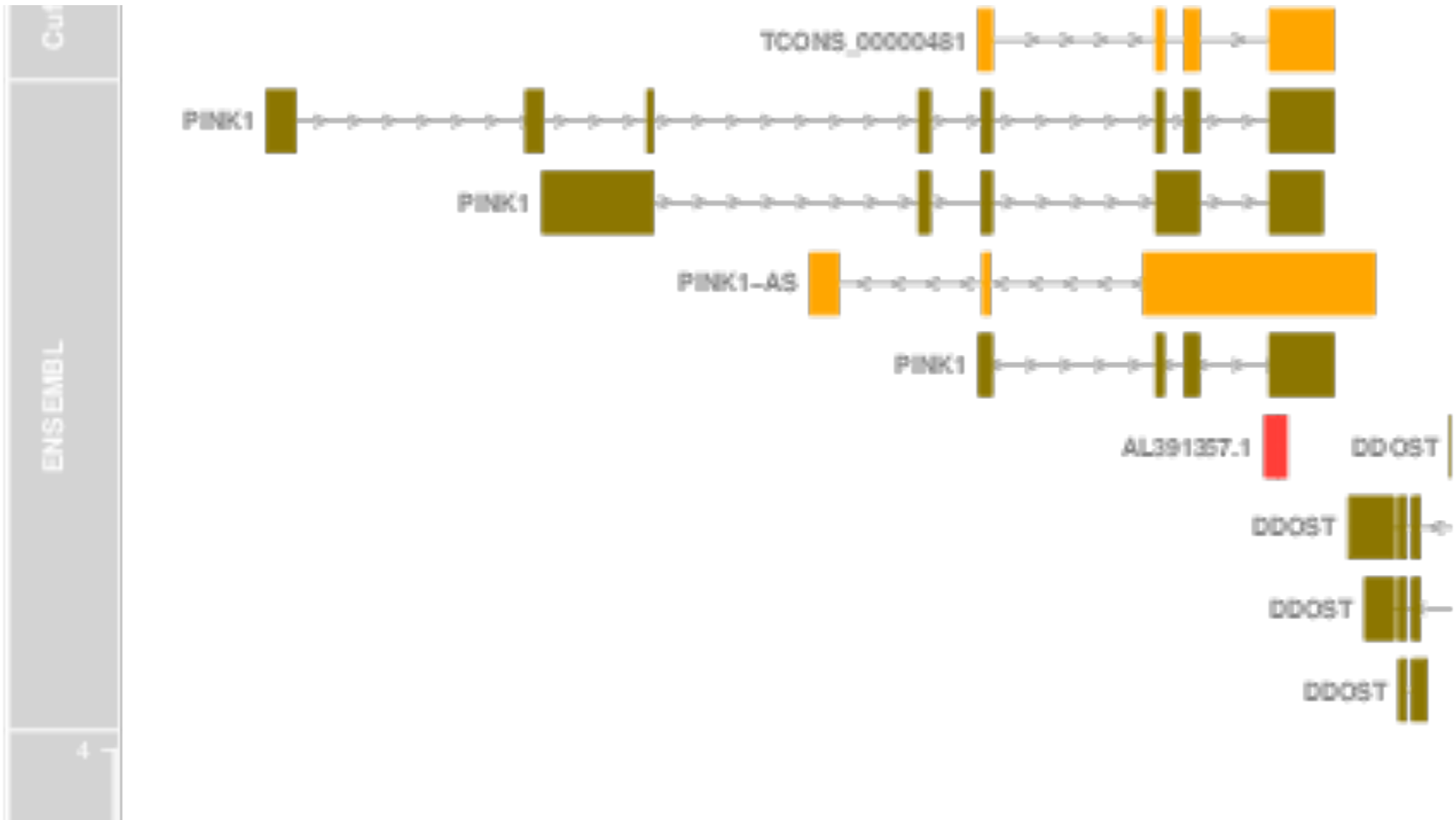
```
sig_genes <- getGenes(cd, geneIdList)
```

```
csHeatmap(sig_genes,  
clustering="row",  
labRow=F)
```



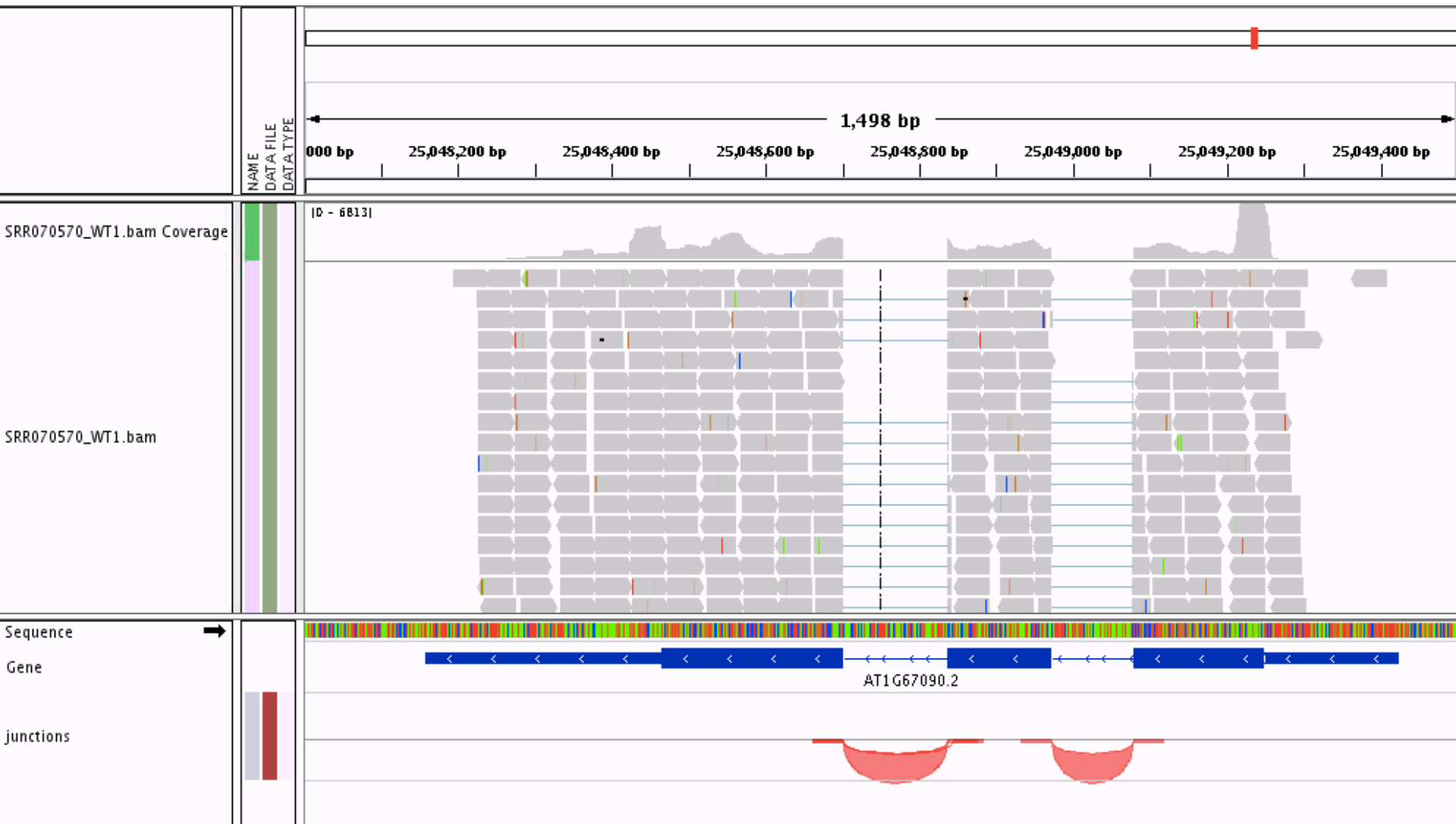
# STEP 6: Visualize and Perform Other Downstream Analysis

Visualize using Cumberbund



# STEP 6: Visualize and Perform Other Downstream Analysis

- Visualization using IGV



# STEP 6: Perform Other Downstream Analysis

- ID enriched gene ontology (GO) terms in our DEGs using GOSeq
- Commands and the examples on the wiki.
- For GO enrichment, we take the following things into account:
  - A. Total number of genes we are looking at.
  - B. Number of genes of interest, that is, in our DEG list.
  - C. Total number of genes in the GO term
  - D. Number of genes from our genes of interest that are also in the GO term.

If the number of genes from our list that belong to GO term (D) is significant compared to the total number of genes in that GO term (C) and the total number of genes in our experiment (A), we consider that GO term to be enriched in our data.

# Thank you!

- Visit the Bioinformatics Consultants at GDC
- Come to Byte Club meetings
  - Join UT Lists-bioiteam



# APPENDIX: Submitting Jobs to Lonestar

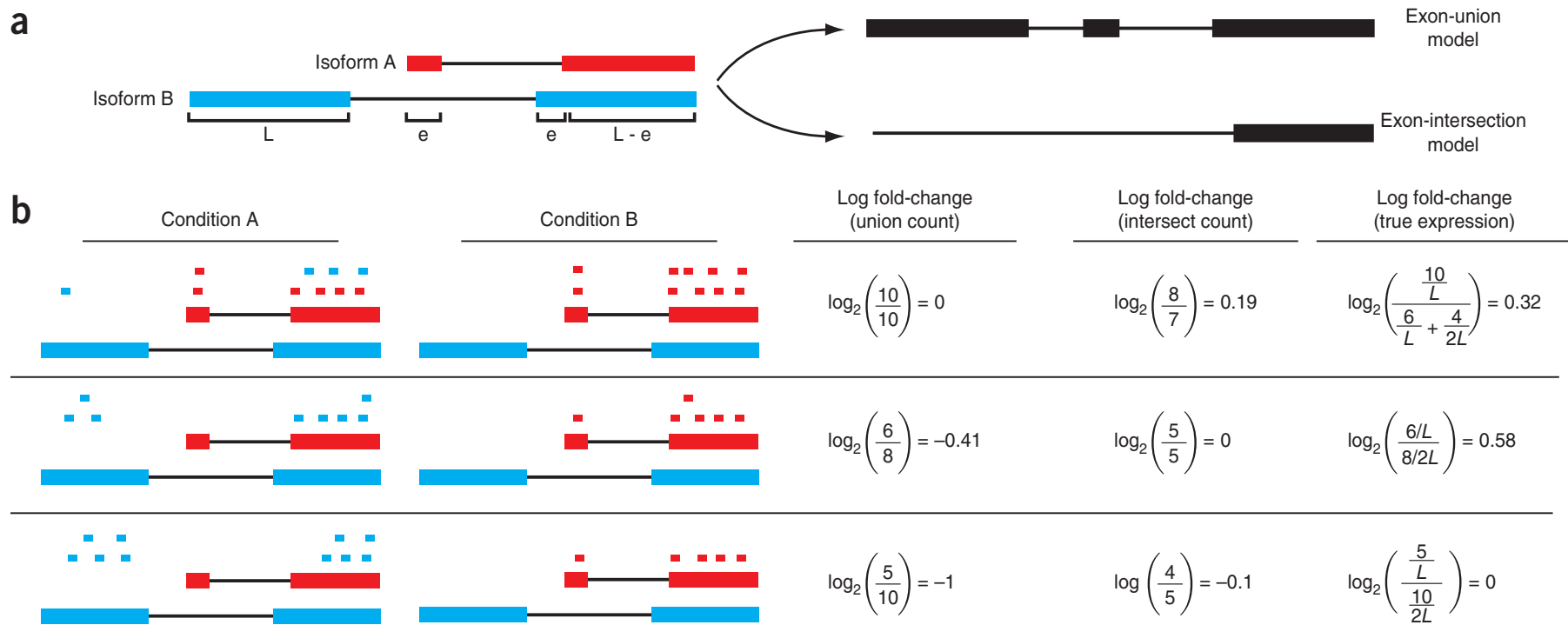
- <https://wikis.utexas.edu/display/bioiteam/Submitting+Jobs+to+Lonestar>

# Of course, Tuxedo Pipeline can be run without looking for novel events

- NO NOVEL JUNCTIONS: Simple differential gene expression analysis against a set of known transcripts.
  - User provides a gff/gtf file containing annotated features. Quantify only the annotated features and id DEGs.
- NOVEL JUNCTIONS ALSO: In addition to known transcripts, novel transcripts should be explored.
  - User provides a gff/gtf file containing annotated features. But you also allow the search for novel variants as well. Both annotated and novel variants are quantified and DEGs are identified.
- ONLY NOVEL/DE NOVO JUNCTIONS: No gff/gtf file is provided. Using just the read data and the genome reference, construct *de novo* transcripts, quantify them and id DEGs.

# Other differential expression tools vs cuffdiff

Others	Cuffdiff
<b>Raw count method</b> for assigning counts to genes	<b>Isoform deconvolution method</b> for assigning counts to genes
Count the reads mapping to <b>exons</b> of each gene/normalization factor = expression for gene	Count the reads that map to each <b>isoform</b> of the gene/normalization factor = expression for gene
If all isoforms of the gene are up/down, works fine	If all isoforms of the gene are up/down, works fine
If some isoforms of the gene are up and some are down, <b>inaccurate results</b>	If some isoforms of the gene are up and some are down, <b>works fine</b>



**Figure 1** Changes in fragment count for a gene does not necessarily equal a change in expression. **(a)** Simple read-counting schemes sum the fragments incident on a gene's exons. The exon-union model counts reads falling on any of a gene's exons, whereas the exon-intersection model counts only reads on constitutive exons. **(b)** Both of the exon-union and exon-intersection counting schemes may incorrectly estimate a change in expression in genes with multiple isoforms. The true expression is estimated by the sum of the length-normalized isoform read counts. The discrepancy between a change in the union or intersection count and a change in gene expression is driven by a change in the abundance of the isoforms with respect to one another. In the top row, the gene generates the same number of reads in conditions A and B, but in condition B, all of the reads come from the shorter of the two isoforms, and thus the true expression for the gene is higher in condition B. The intersection count scheme underestimates the true change in gene expression, and the union scheme fails to detect the change entirely. In the middle row, the intersection count fails to detect a change driven by a shift in the dominant isoform for the gene. The union scheme detects a shift in the wrong direction. In the bottom row, the gene's expression is constant, but the isoforms undergo a complete