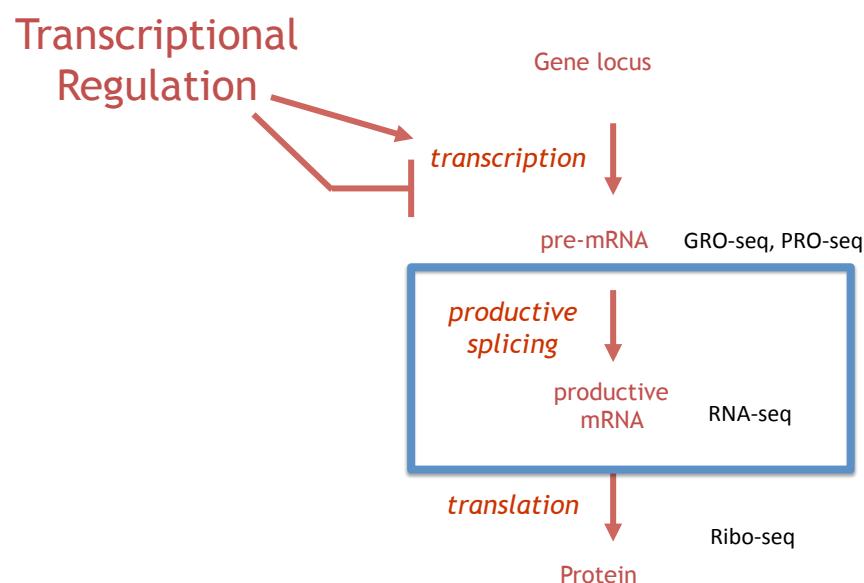
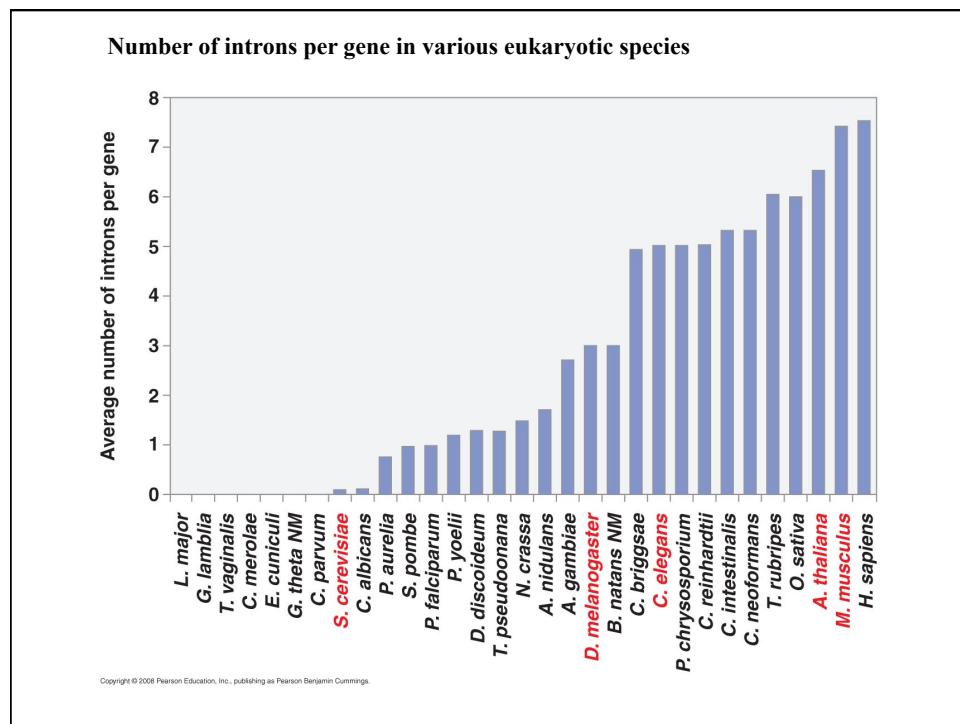
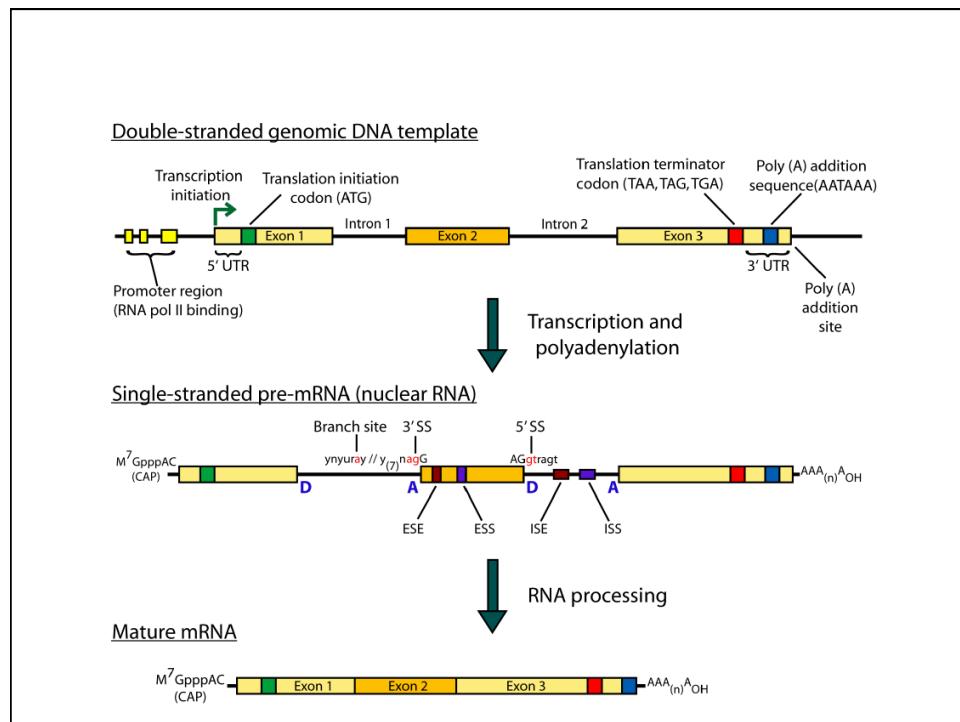
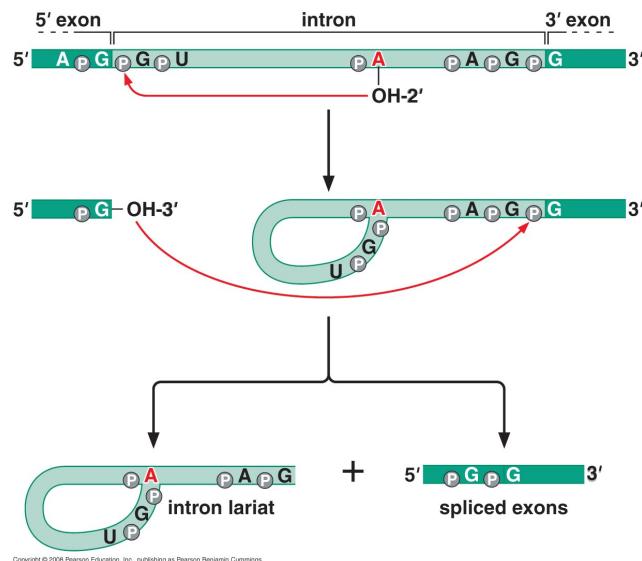


“Everybody is talented because everybody who is human has something to express.”  
— [Brenda Ueland](#)





**The intron is removed in a form called lariat and the flanking exons are joined**



The Nobel Prize in Chemistry 1989  
"for their discovery of catalytic properties of RNA"



**Sidney Altman**

1/2 of the prize

Canada and USA

Yale University  
New Haven, CT, USA

b. 1939



**Thomas R. Cech**

1/2 of the prize

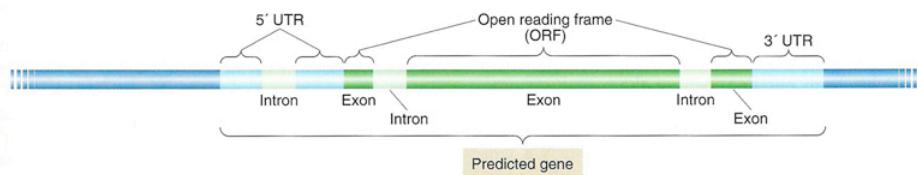
USA

University of Colorado  
Boulder, CO, USA

b. 1947

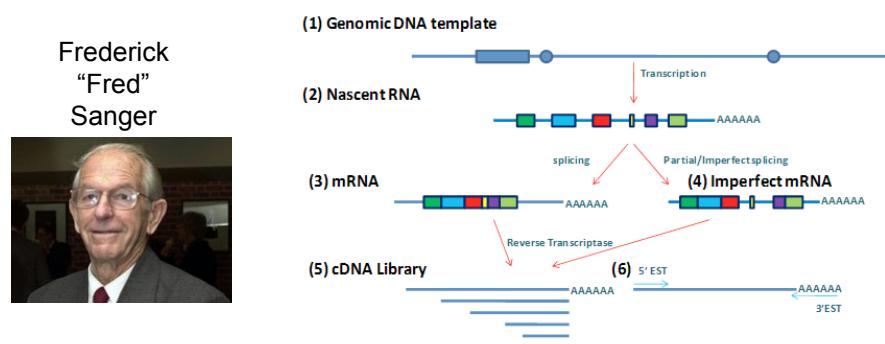
## Applications of RNA-seq

- Differential expression
- Gene fusion
- Alternative splicing
- Novel transcribed regions
- Allele-specific expression
- RNA editing
- Transcriptome assembly for non-model organisms



## The earliest approach to expression

Frederick  
“Fred”  
Sanger

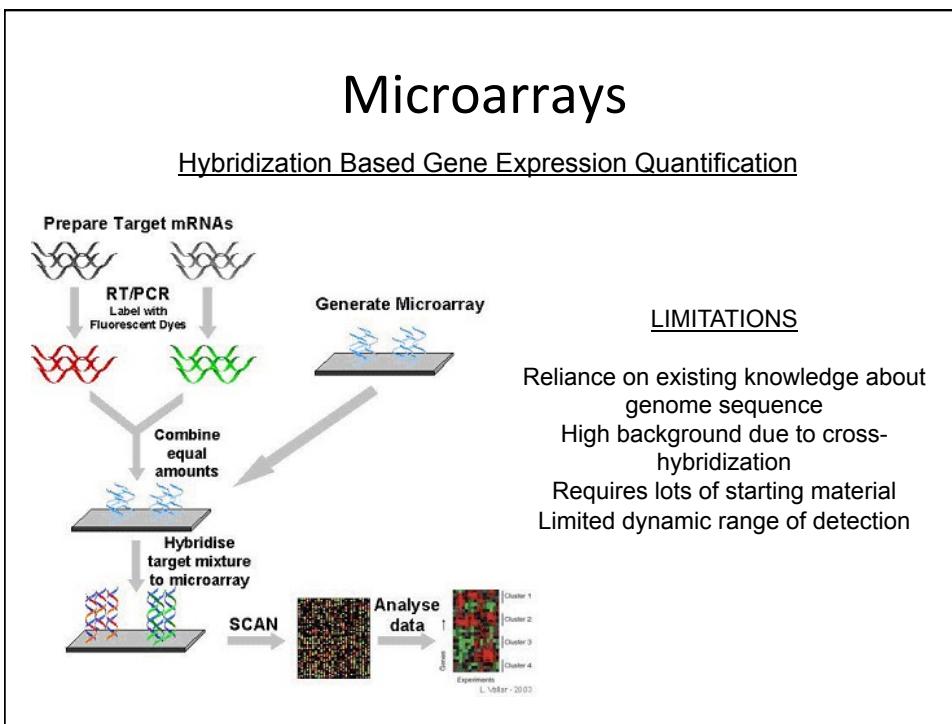
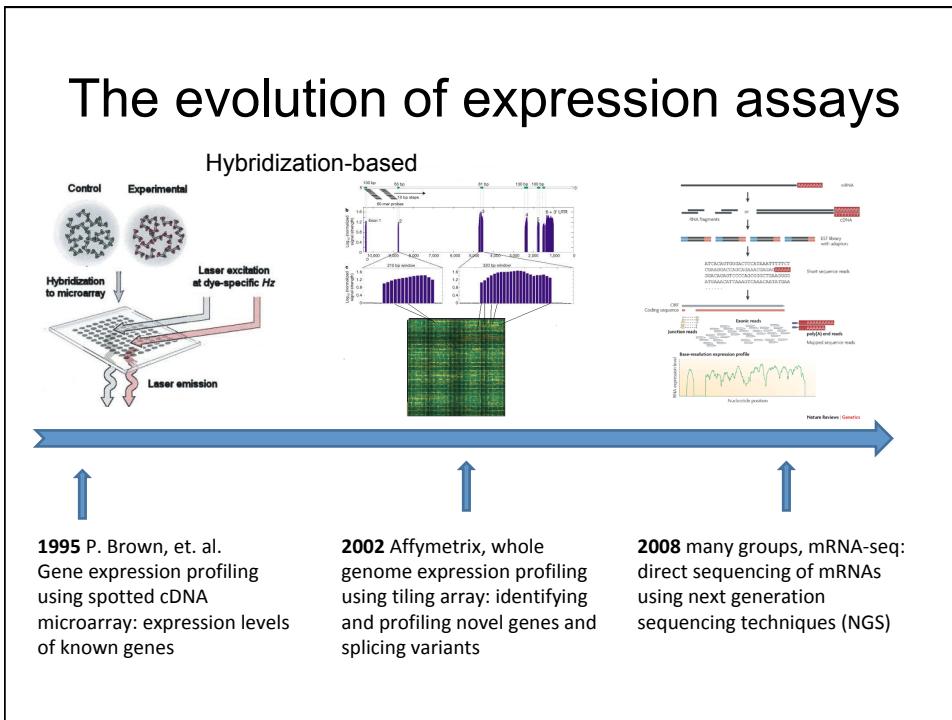


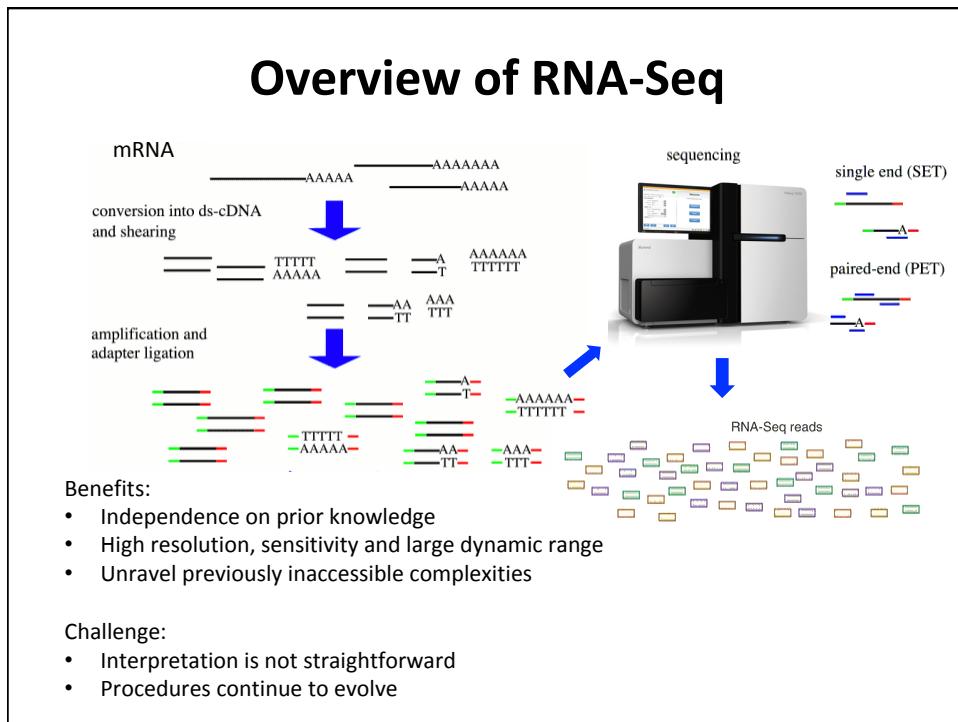
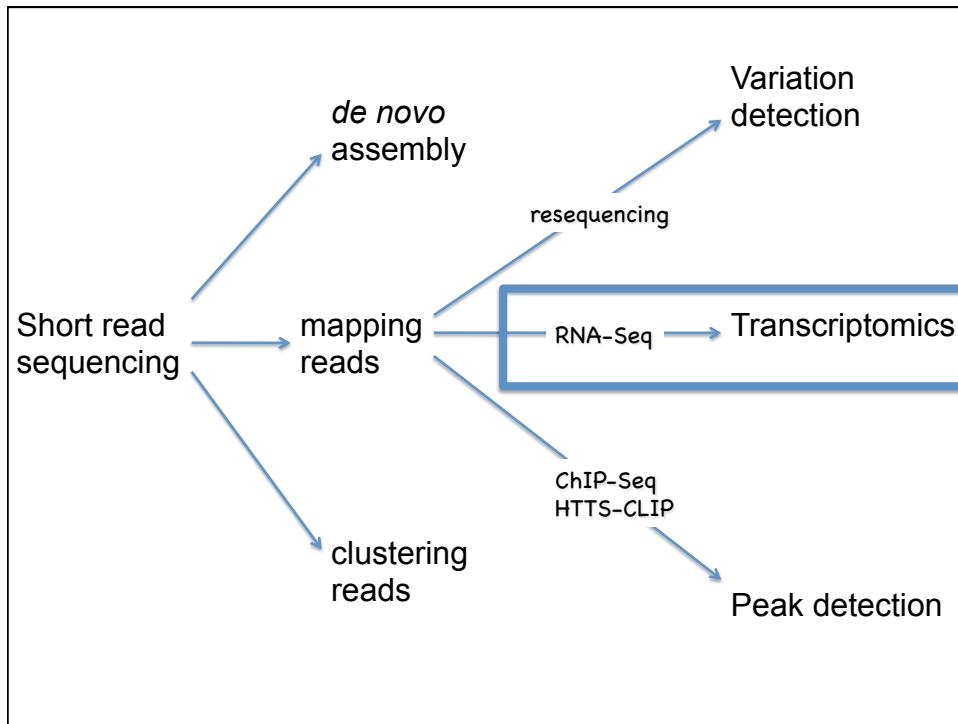
Tag Based Sequencing Approaches  
Serial Analysis of Gene Expression (SAGE)  
Cap Analysis of Gene Expression (CAGE)

} Expressed Sequence Tags (ESTs)

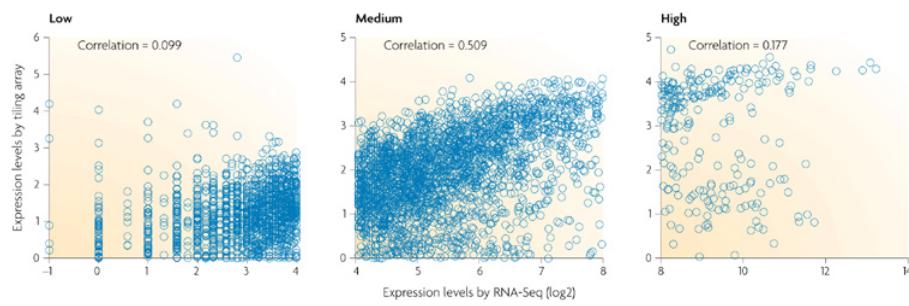
### Limitations of Sanger Sequencing

- Low throughput
- Expensive
- Not quantitative (no enough depth!)





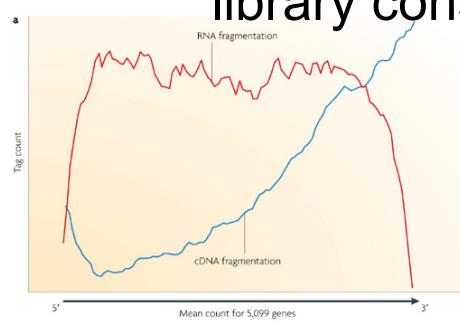
## RNA-seq and microarray agree fairly well only for genes with medium levels of expression



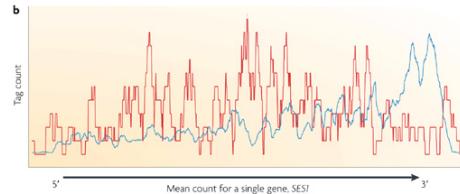
Nature Reviews | Genetics

*Saccharomyces cerevisiae* cells grown in nutrient-rich media. Correlation is very low for genes with either low or high expression levels because microarray saturates at high levels and both techniques struggle with very lowly expressed.

## Challenges for RNA-Seq: library construction



Fragmentation of oligo-dT primed cDNA (blue line) is more biased towards the 3' end of the transcript. RNA fragmentation (red line) provides more even coverage along the gene body, but is relatively depleted for both the 5' and 3' ends.



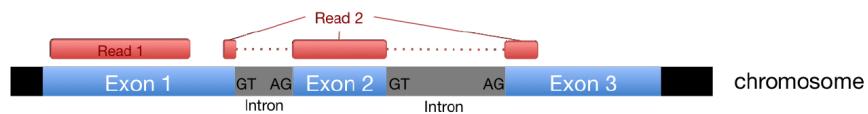
A specific yeast gene, SES1 (seryl-tRNA synthetase)

## Mature RNA presents a unique problem for read mapping ...

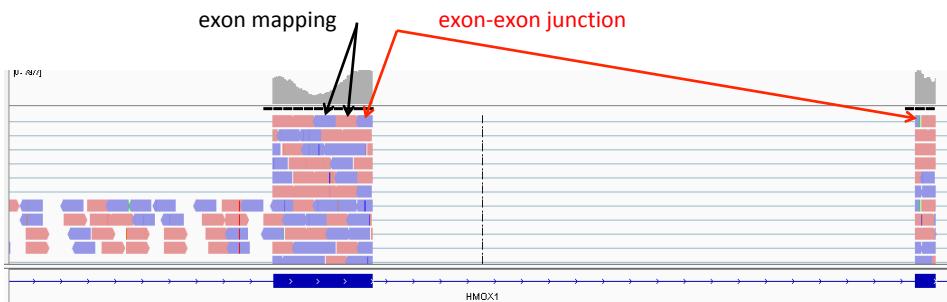
(a) Aligning to the transcriptome



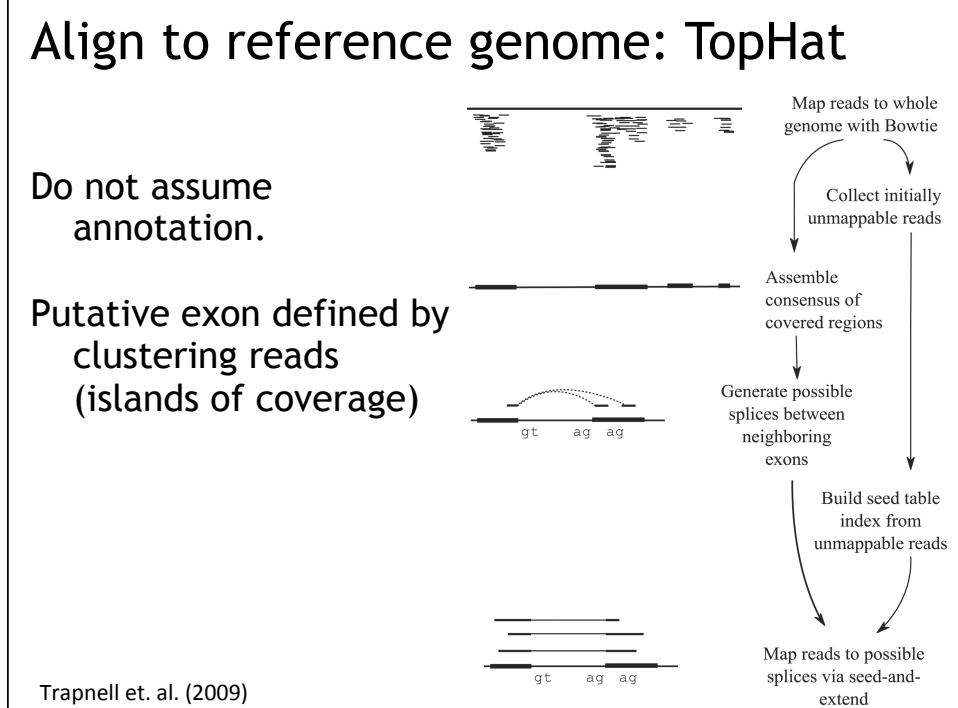
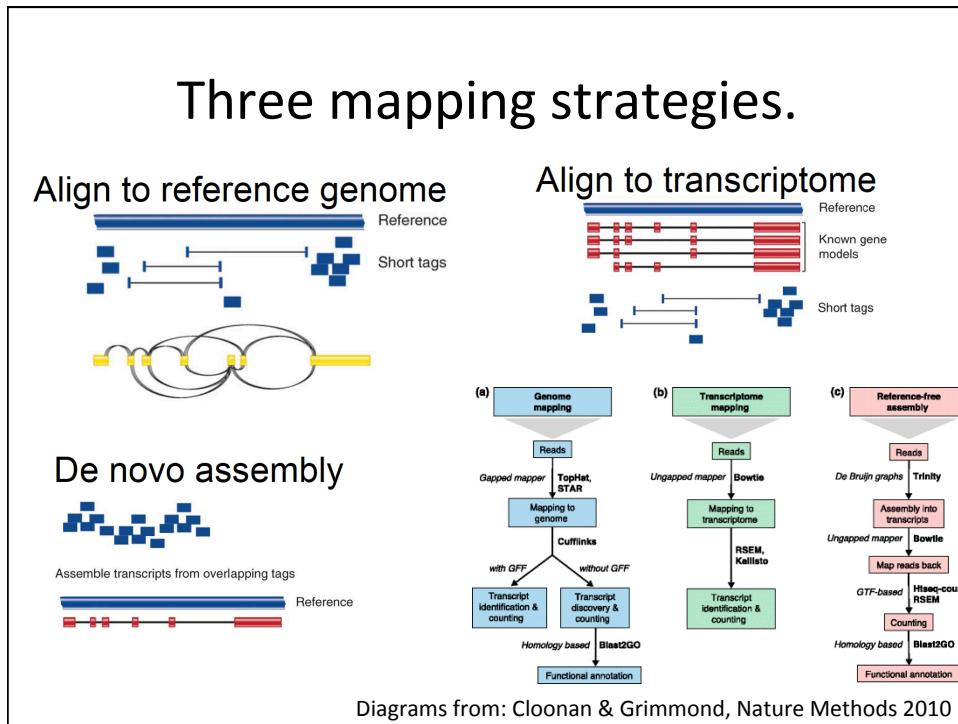
(b) Aligning to the genome



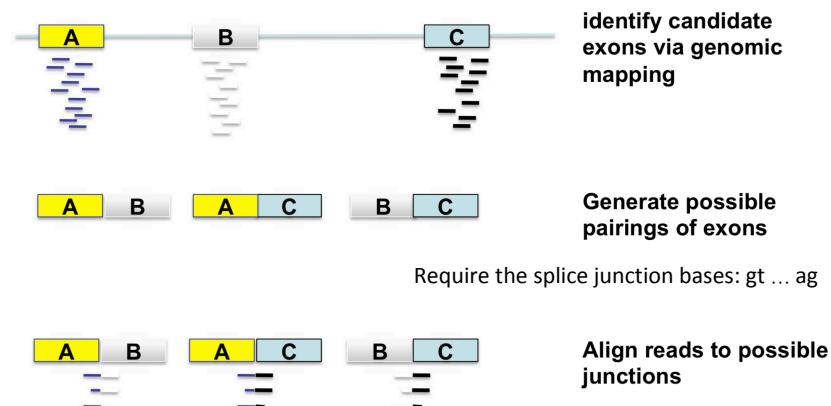
## Read mapping



Unlike DNA-Seq, when mapping RNA-Seq reads back to reference genome, we need to pay attention to **exon-exon junction reads**



Consider all possible pairings of exons.

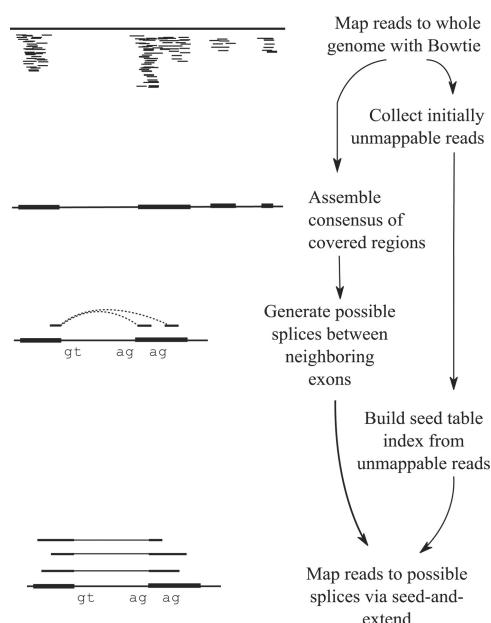


## TopHat: splice junctions

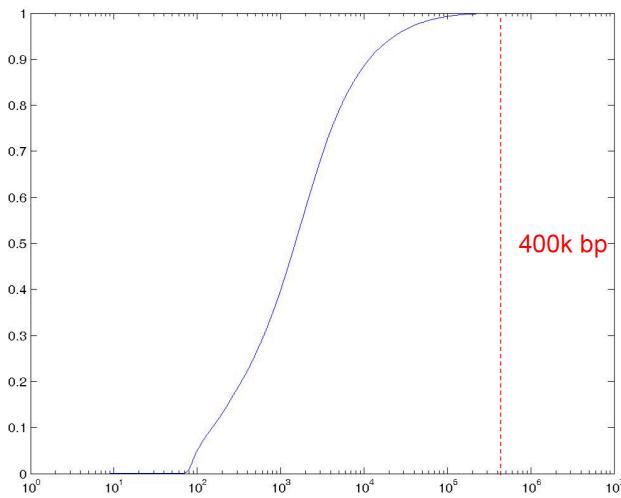
Find GT-AG pairing sites between neighboring islands

The distance between two sites should  $> 70\text{bp}$  and  $< 20\text{k bp}$ , as most intron length lies within this range

Build exon based “genome” for mapping split reads.



### The cumulative distribution function of the intron sizes

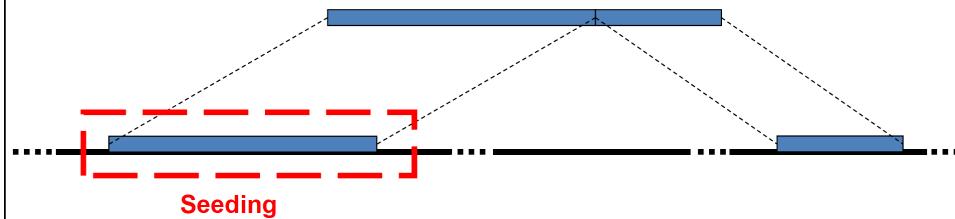


\* Based on hg19 human Refseq annotation.

### Alternative: Finding split reads

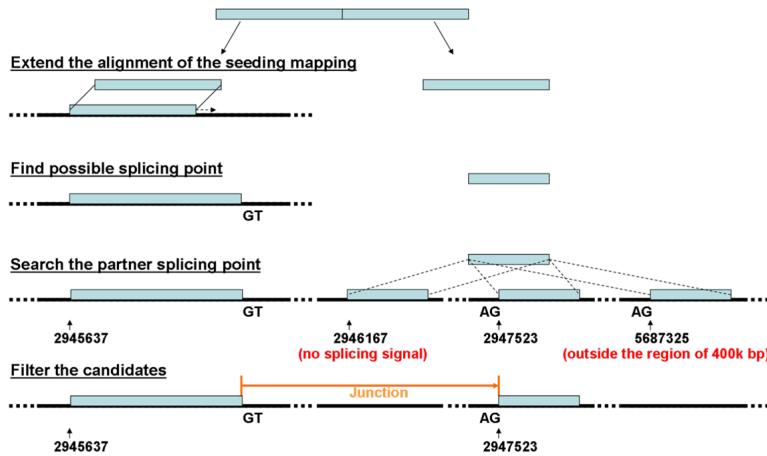
#### Basic concept

- Split map

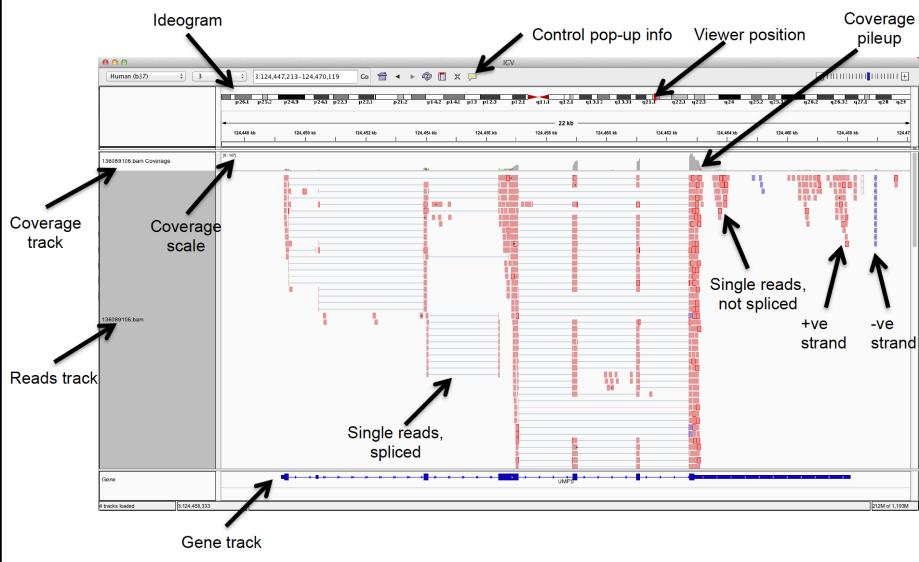


- if read length  $\geq 50$  bp
- at least one of the halves will have non-split map

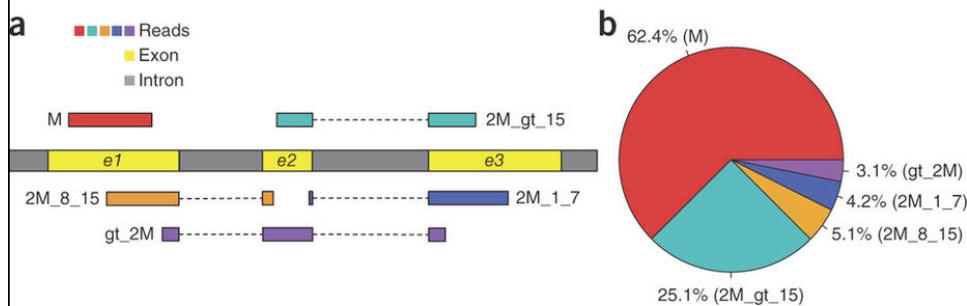
## Junction search:



## Many splice junctions per gene

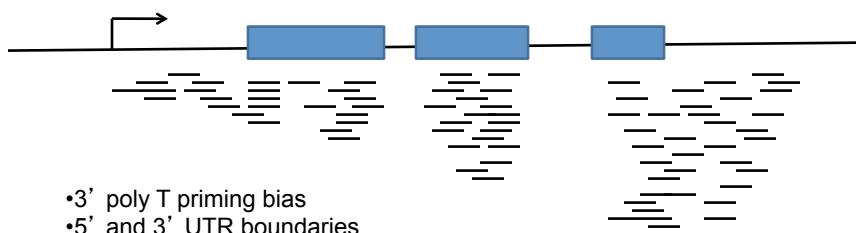


## But many types of reads possible



Reads with poor anchoring (orange, dark blue) or spanning multiple exons (purple) were largely missed by first generation splice junction aligners (TopHat) but are now being discovered by second generation aligners (HISAT).

## Expression quantification can be non-trivial



- 3' poly T priming bias
- 5' and 3' UTR boundaries
- Alternate splicing
- Cryptic exons
- Cryptic start sites
- Paired-end reads can help

## Expression quantification

Gene A: 200

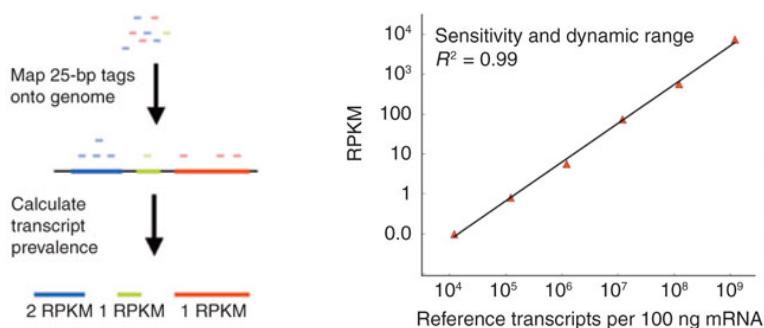
Gene B: 300

Is expression of Gene A < Expression of Gene B?

The number of reads is roughly proportional to

- the length of the gene
- the total number of reads in the library

## Quantification of known transcripts



- The expression levels of known transcripts (*exon model*) are measured by the number of reads per kilobase of transcript per million mapped reads (RPKM)

Mortazavi, A., Williams, B. A., McCue, K., Schaeffer, L. & Wold, B. Mapping and quantifying mammalian transcriptomes by RNA-Seq. *Nature Methods* 5, 621–628 (2008).

## Gene and isoform abundance

### Gene abundance

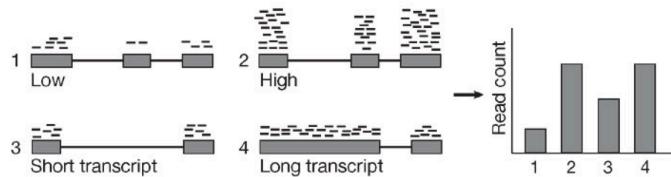
--**RPKM**: reads per kilobase of the transcript per million mapped reads to the transcriptome (Mortazavi et.al., 2008)

1 RPKM ~ 0.3 to 1 transcript per cell

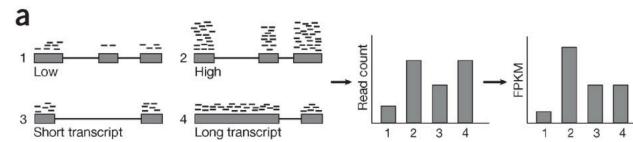
### Cons:

Could not be used to measure isoform abundance, because most reads are *shared by multiple isoforms*.

- Question: What are the RPKM-corrected expression values and why?



- Note especially normalization for fragment length (transcripts 3 and 4)



Graphic credit: Garber et al. (2011) *Nature Methods* 8:469–477. Note that the authors here use the related term **FPKM**, **Fragments per KB per million reads**, which is suitable for paired-end reads (we will not cover the details here).

