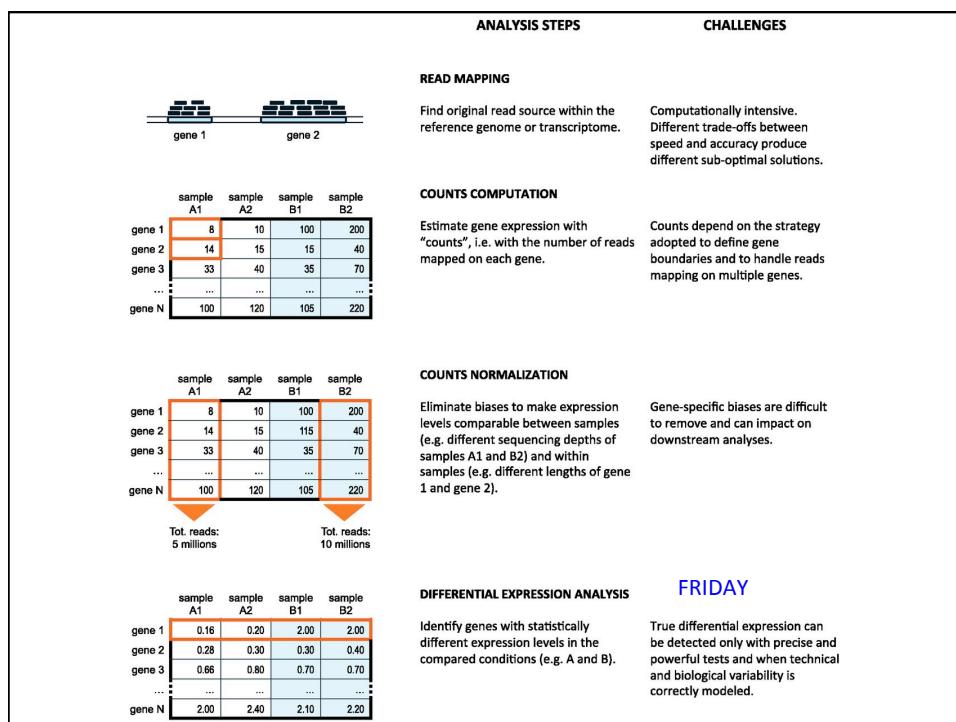


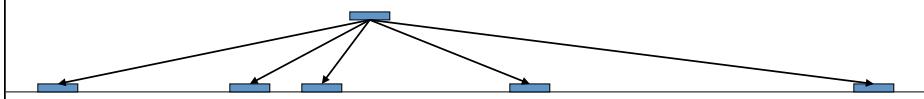
“Sometimes a writer wants to cut out unnecessary information from a *quotation* to place emphasis on the part of it that matters most or to eliminate unnecessary information.”

<https://quoting1.weebly.com/splicing-quotes.html>



Handling multi-map reads in counts

- a non-negligible fraction of RNA-seq reads are ‘multireads’: reads that map with comparable fidelity on multiple positions of the reference.



- Ignore them (earliest approaches)
- Weight them (“fragments”)
- Sophisticated approaches (“nearby” signal)

Recall: Normalization methods

- **RPKM: Reads Per Kilobase Million**
 - Count up the total reads in a sample and divide that number by 1,000,000 – this is our “per million” scaling factor.
 - Divide the read counts by the “per million” scaling factor. This normalizes for sequencing depth, giving you reads per million (RPM)
 - Divide the RPM values by the length of the gene, in kilobases. This gives you RPKM.
- **FPKM: Fragments Per Kilobase Million**
 - With paired-end RNA-seq, two reads can correspond to a single fragment, or, if one read in the pair did not map, one read can correspond to a single fragment. FPKM takes into account that two reads can map to one fragment (and so it doesn’t count this fragment twice).
 - Also handles multi-read mapping better (“fragments”)

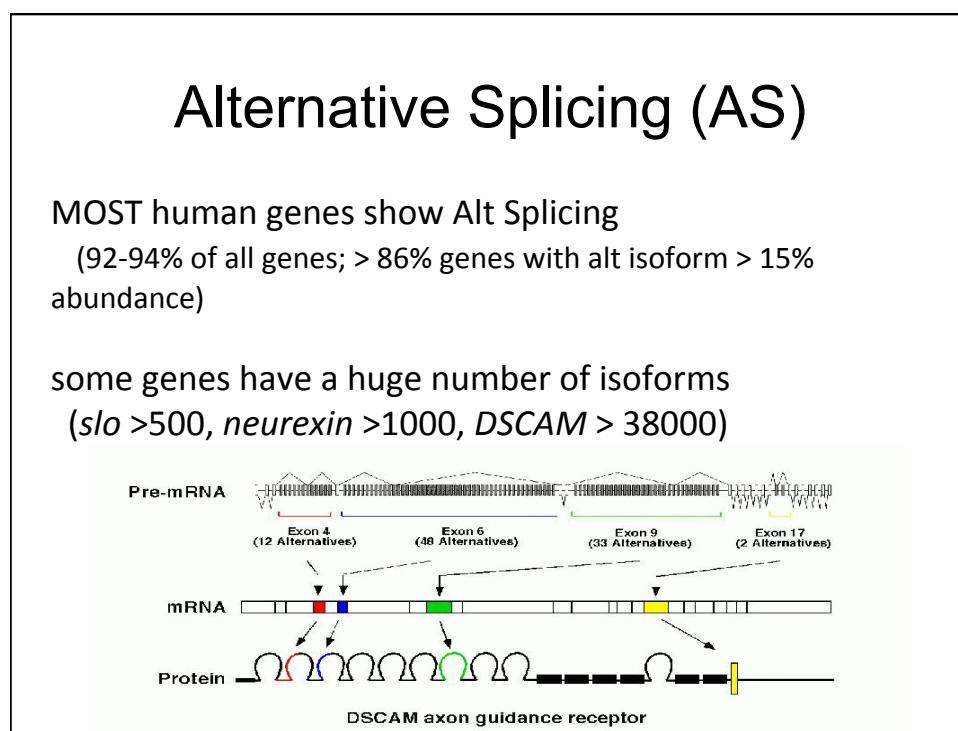
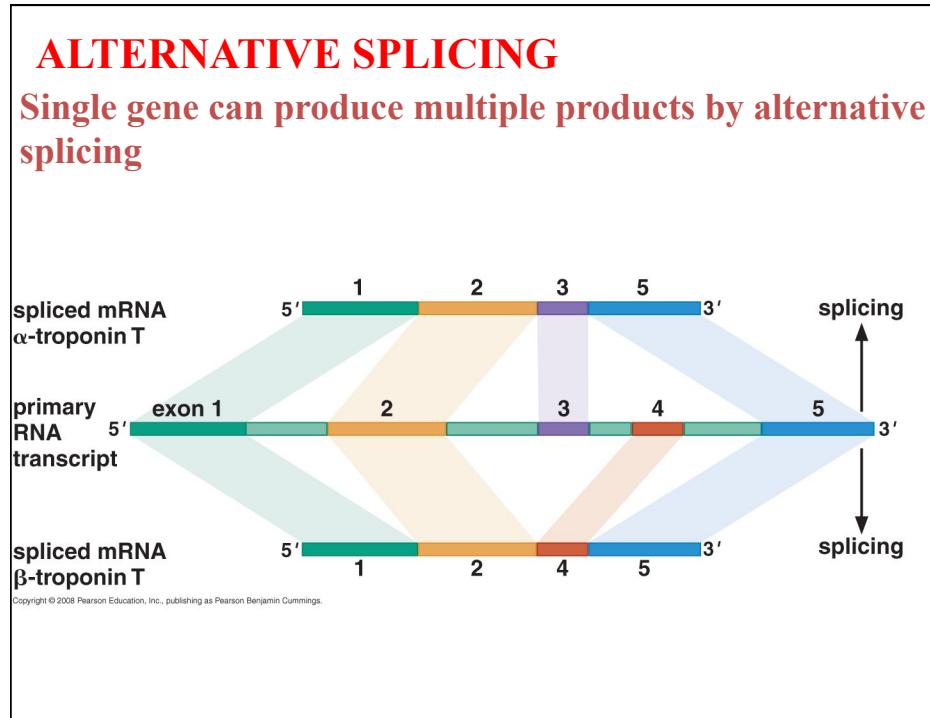
One more simple normalization

- **TPM: Transcripts Per Kilobase Million**
 - Divide the read counts by the length of each gene in kilobases. This gives you reads per kilobase (RPK).
 - Count up all the RPK values in a sample and divide this number by 1,000,000. This is your “per million” scaling factor.
 - Divide the RPK values by the “per million” scaling factor. This gives you TPM.
- The sum of all TPMs in each sample are the same. This makes it easier to compare the proportion of reads that mapped to a gene in each sample.
- However, requires that your “feature set” be more or less complete.

Simple abundance estimates are inappropriate for isoforms.

	Exon 1	Exon 2	Exon 3	abundance
Isoform 1				x_1
Isoform 2				x_2
Isoform 3				x_3
Length	l_1	l_2	l_3	
#reads	n_1	n_2	n_3	

In other words, RPKM/FPKM is too simplistic for isoforms because a read's origin is unclear. So how do we infer isoforms to properly identify a read's likely “origin”?



Large collections of splicing data (pre-short read sequencing era)

Curated databases

SWISS-PROT and RefSeq both support annotation of experimentally supported alternative splicing



cDNA Sequencing Projects



RIKEN sequenced >21000 full length mouse cDNAs

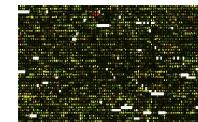
Many other projects underway (human, fly, plants,...)



Shinagawa et al. (2001) *Nature* 409:685–90

Microarray detection

Direct or indirect alternative splicing detection



Hu et al. (2001) *Genome Res* 11:1237–45
Yeailey et al. (2002) *Nat Biotech* 20:353–9

Boguski et al. (1993) *Nat Gen* 4:332–3

Public EST data sources (dbEST)

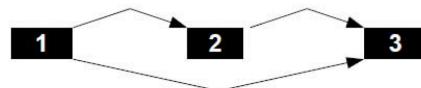
>4.5 million human EST sequences

>12 million total EST sequences

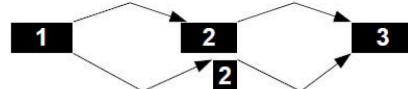
About 1000 new sequences per day

First step: detect splice junctions

Exon Skipping

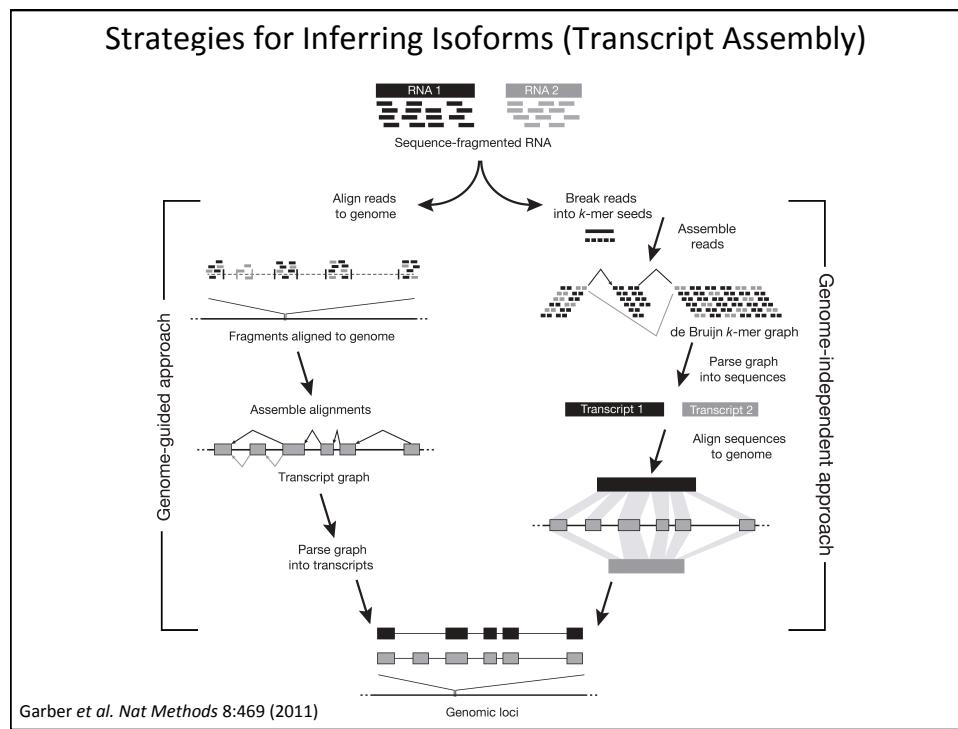
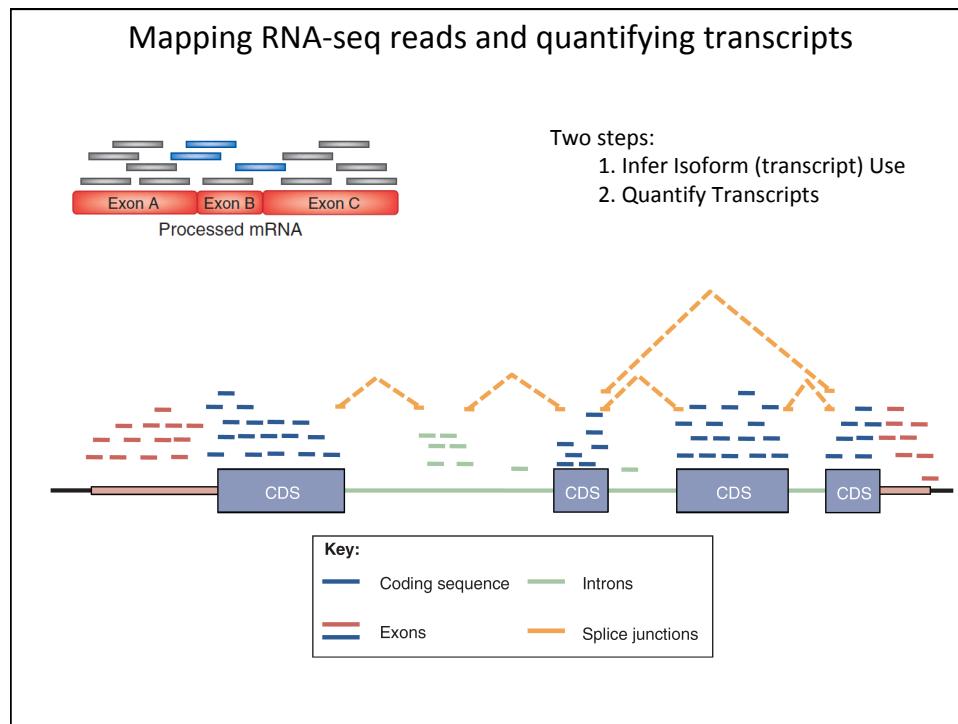


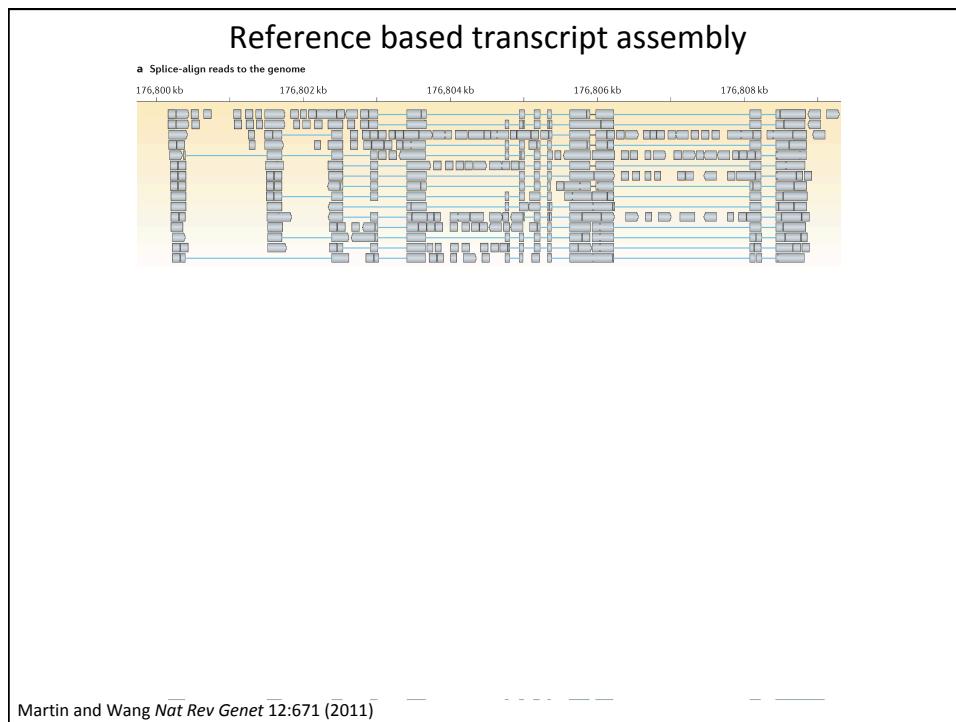
Alternate Exon Length



Mutually Exclusive Exons



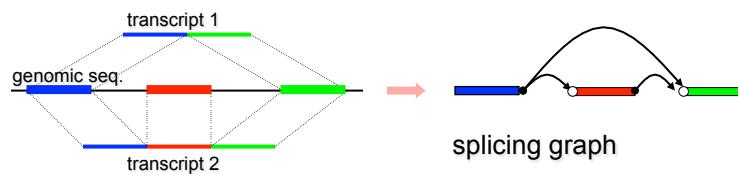




Splicing graph construction

If a reference genome is used:

- Map reads to the reference genome (short read aligner)
- Check alignment (splice sites, quality)
- Connect consecutive positions
- Build splicing graph



Splice graph approach

Replace the problem of finding a list of consensus sequences



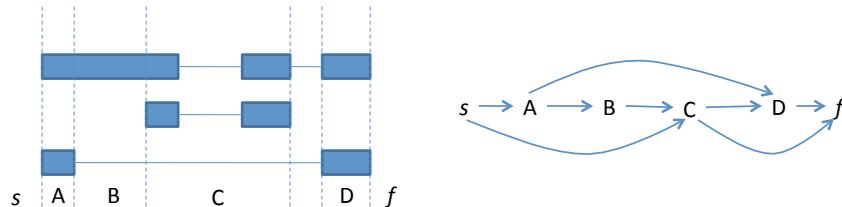
with ***Graph Reconstruction Problem***:

Given a set of expressed sequence, find a minimal graph (*splicing graph*) representing **all** transcripts as paths.



Heber, et. al. ISMB 2002

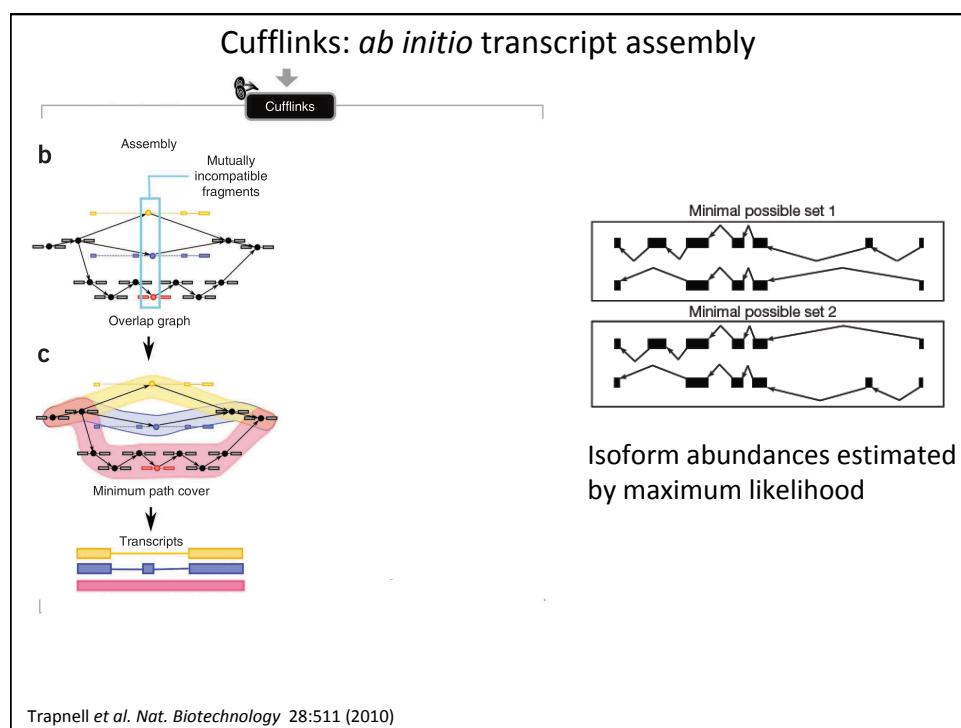
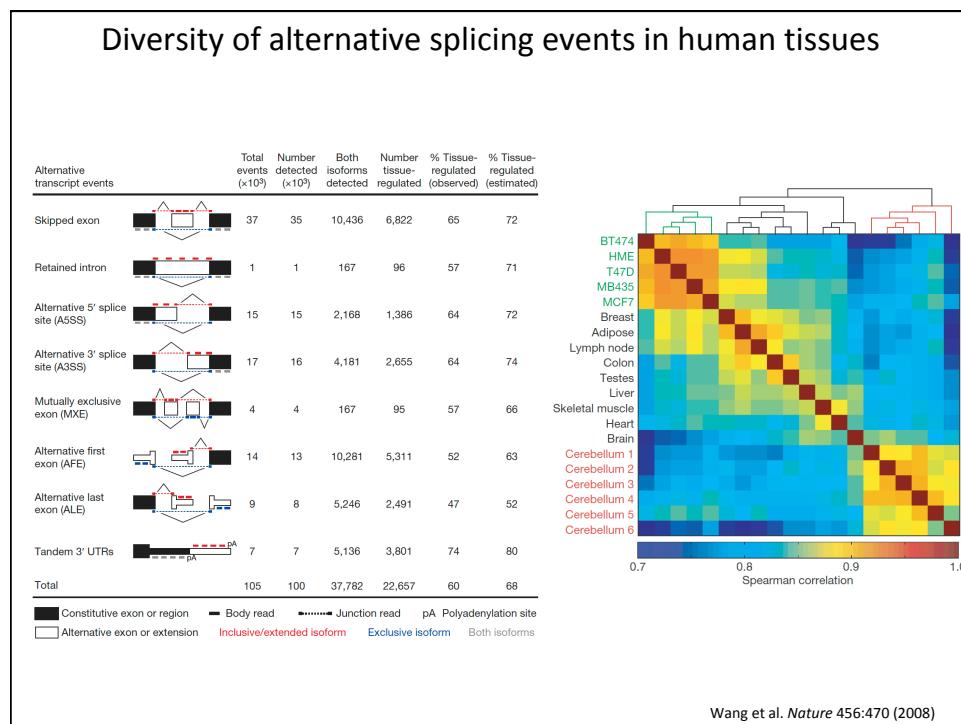
Splicing graph and splicing variants



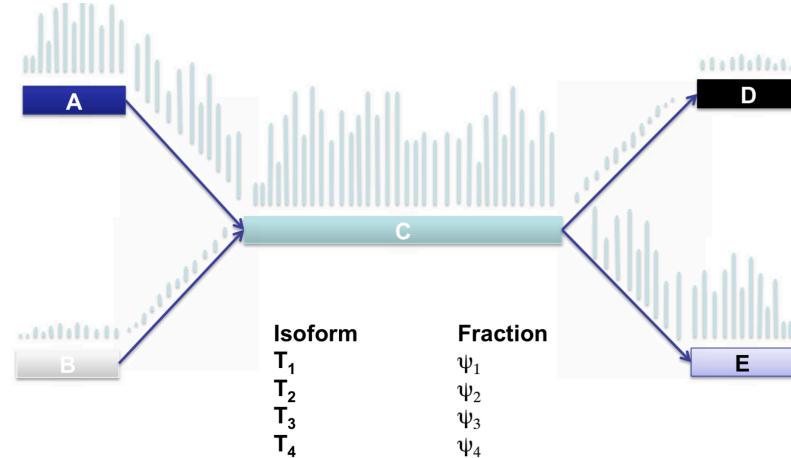
An edge in the splicing graph, called a *block*, represents a maximal sequence of adjacent exons or exon fragments that always appear together in a given set of splicing variants. Therefore, variants can be represented by sequence of blocks, e.g. {ABCD, C, AD}.

Vertices *s* and *f* are included into graph, and are linked to the 5' and 3' of each variant, respectively.

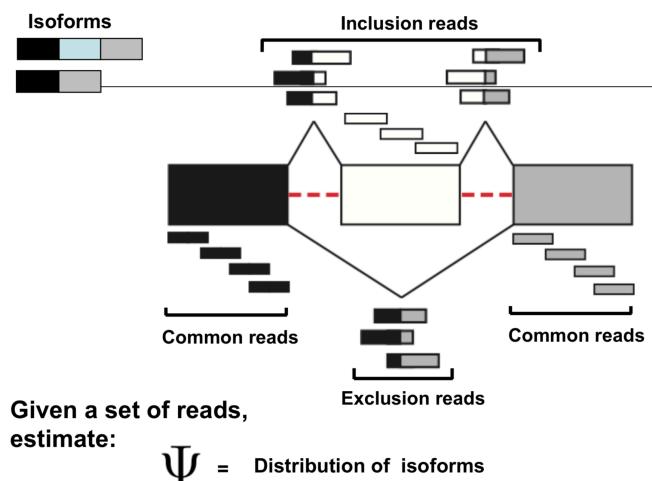
Lacroix, et. al. WABI, 2009



We can use mapped reads to learn the mixture of isoforms present.



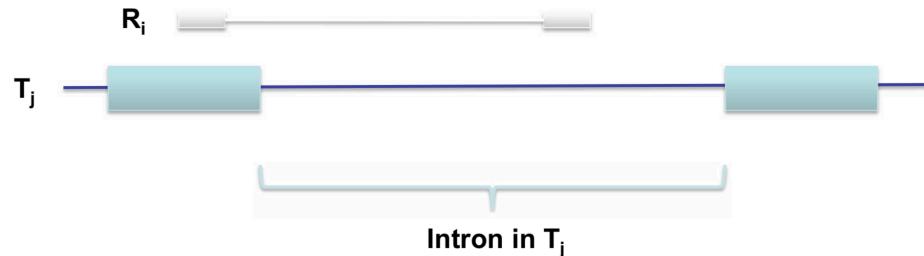
Our dataset contains a mixture of isoforms



Some reads can be excluded from some isoform possibilities.

If a single ended read or read pair R_i is structurally incompatible with transcript T_j , then

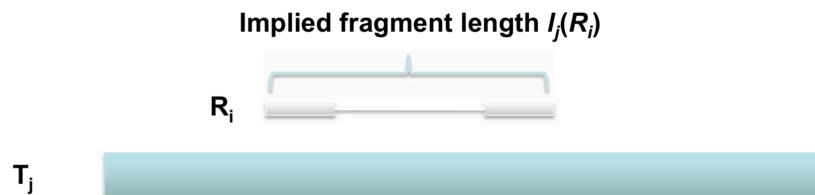
$$P(R = R_i | T = T_j) = 0$$



Some reads have a probability of originating with a particular isoform.

Assume our library fragments have a length distribution described by a probability density F . Thus, the probability of observing a particular paired alignment to a transcript:

$$P(R = R_i | T = T_j) = \frac{F(l_j(R_i))}{l_j}$$



Our probabilistic model

- Find expression abundances ψ_1, \dots, ψ_n for a set of isoforms T_1, \dots, T_n
- Observations are the set of reads R_1, \dots, R_m

$$P(R | \Psi) = \prod_{i=0}^m \sum_{j=0}^n \Psi_j P(R_i = R_i | T = T_j)$$

This is weighted sum ..

$$L(\Psi | R) \propto P(R | \Psi) P(\Psi)$$

This is BAYES again!

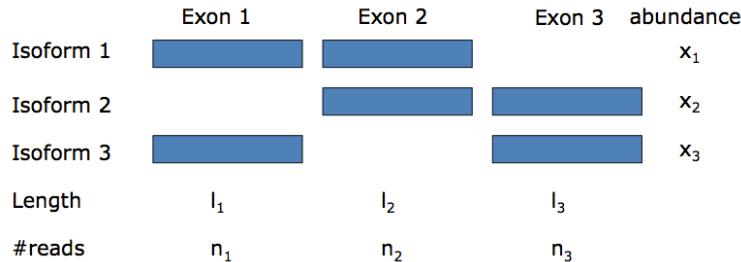
$$\Psi = \underset{\Psi}{\operatorname{argmax}} L(\Psi | R)$$

Our optimization.

- Can estimate mRNA expression of each isoform using total number of reads that map to a gene and ψ

Isoform Inference is akin to a “system of equations”

- If given known set of isoforms



- Estimate x to maximize the likelihood of observing n

Splicing graph construction

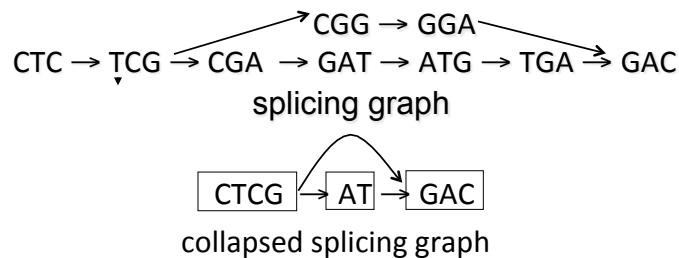
If a reference genome sequence is not used:

- Break sequences into k -mers (20-mers).
- Build graph using k -mers as vertices, connect them iff they occur consecutively in a sequence [Pevzner et al., 2001].

Example (3-mers):

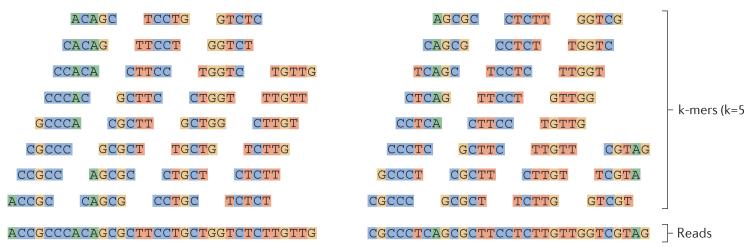
Sequences: CTCGATGAC, CTCGGAC

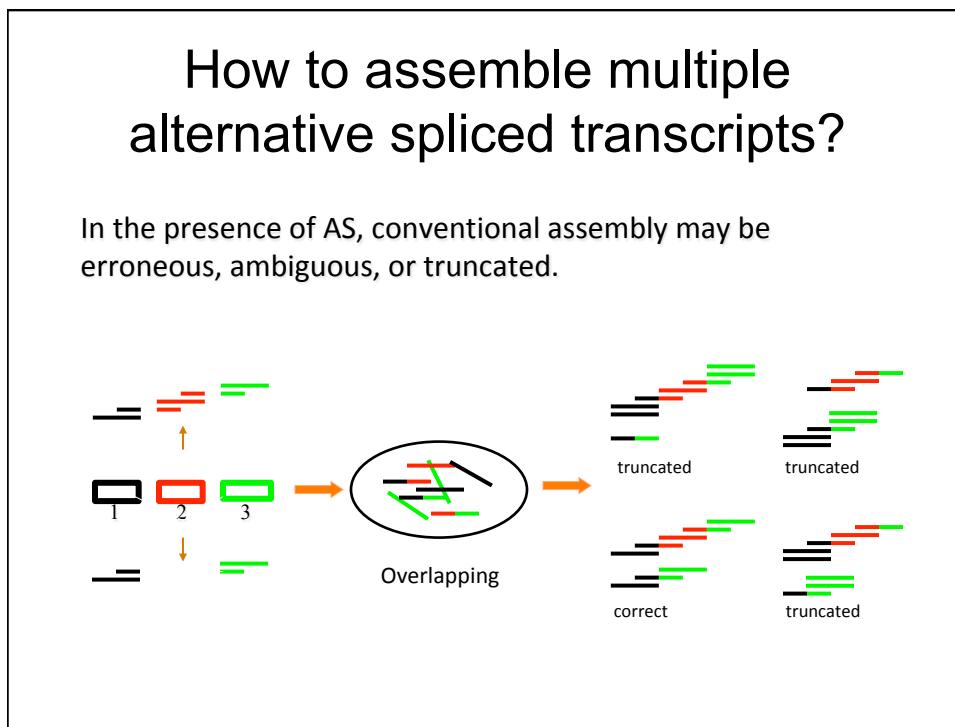
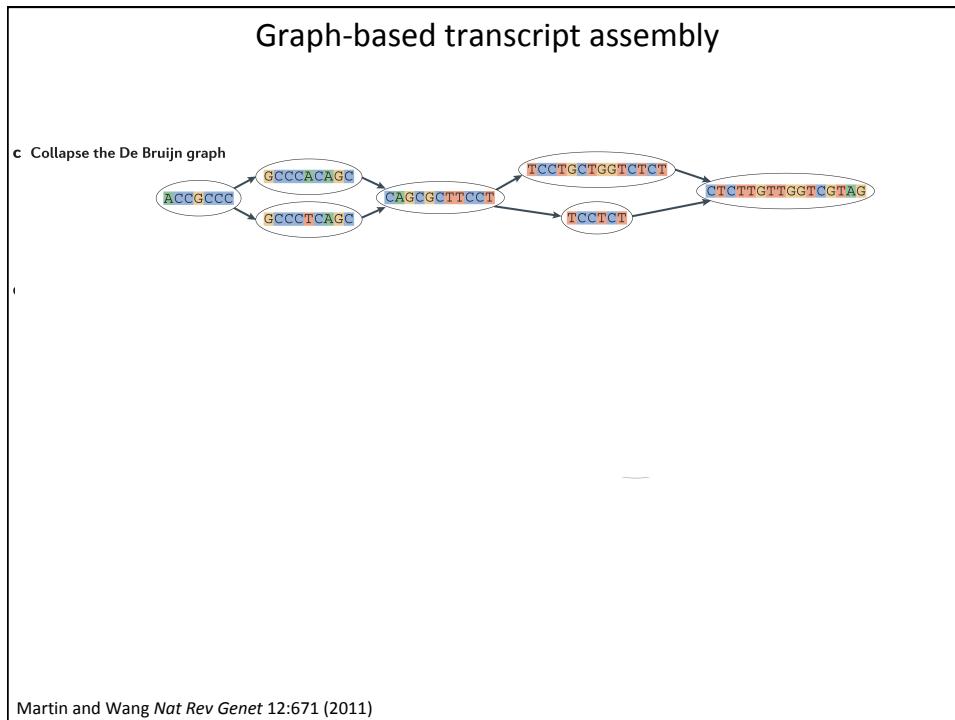
Vertices: {CTC, TCG, CGA, GAT, ATG, TGA, GAC, CGG, GGA}



Graph-based transcript assembly

a Generate all substrings of length k from the reads

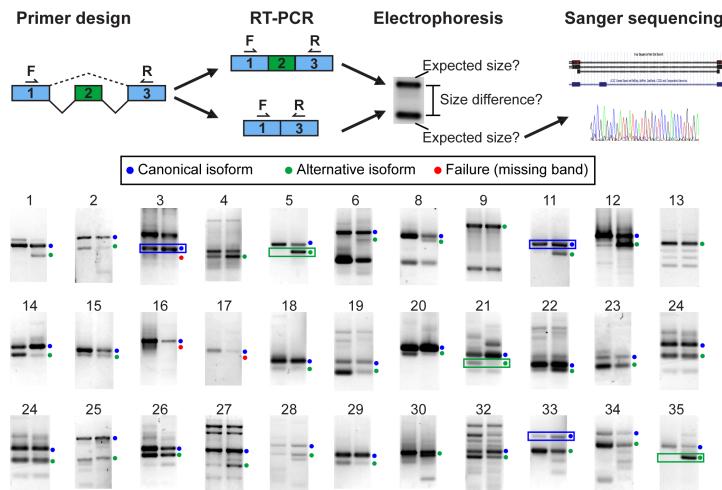




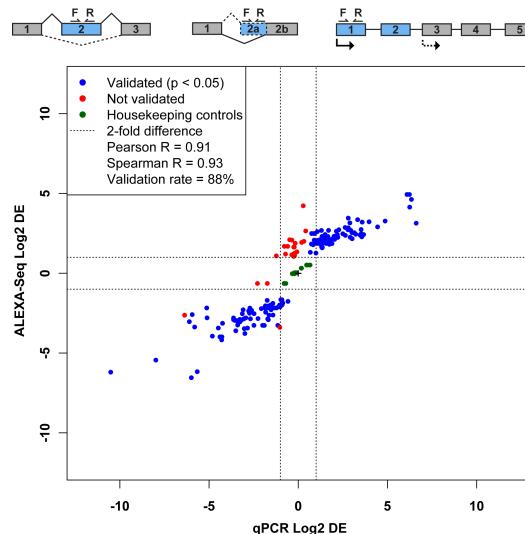
Common questions: how reliable are expression predictions from RNA-seq?

- Are novel exon-exon junctions real?
 - What proportion validate by RT-PCR and Sanger sequencing?
- Are differential/alternative expression changes observed between tissues accurate?
 - How well do DE values correlate with qPCR?
- 384 validations
 - qPCR, RT-PCR, Sanger sequencing
- See ALEXA-Seq publication for details:
 - Also includes comparison to microarrays
 - Griffith et al. *Alternative expression analysis by RNA sequencing*. Nature Methods. 2010 Oct;7(10):843-847.

Validation (qualitative)



Validation (quantitative)

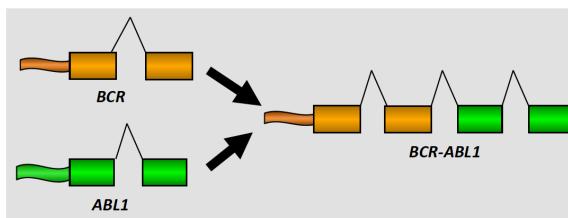


qPCR of 192
exons identified
as alternatively
expressed by
ALEXA-Seq

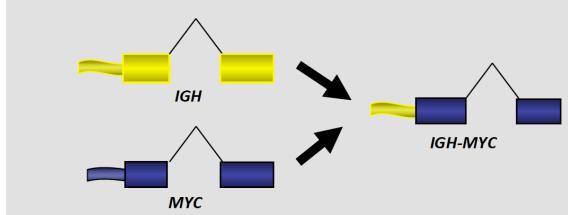
Validation rate = 88%

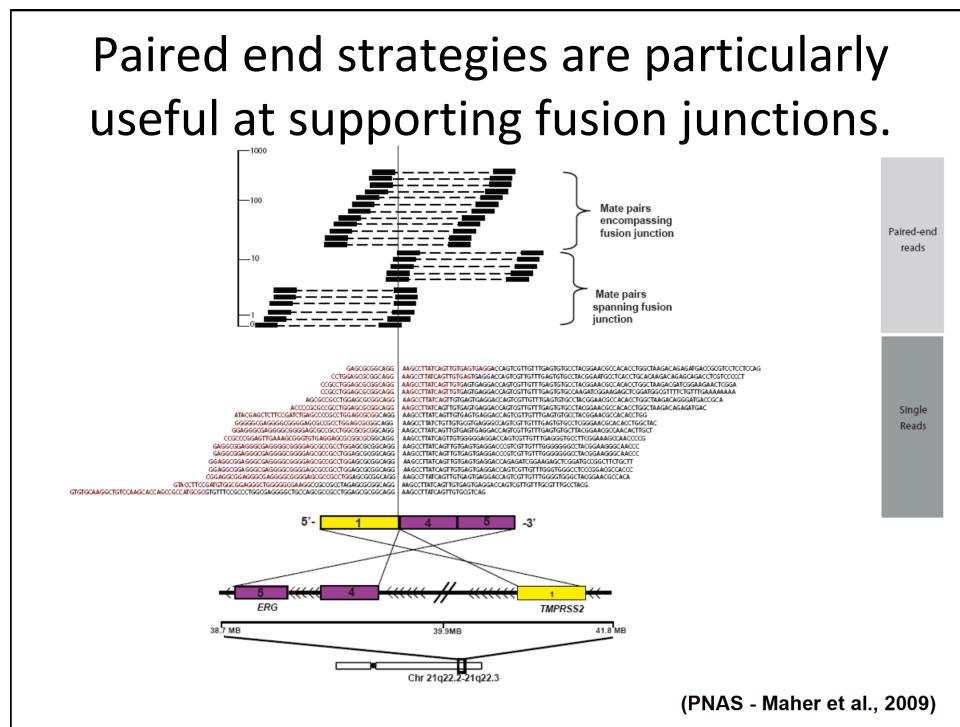
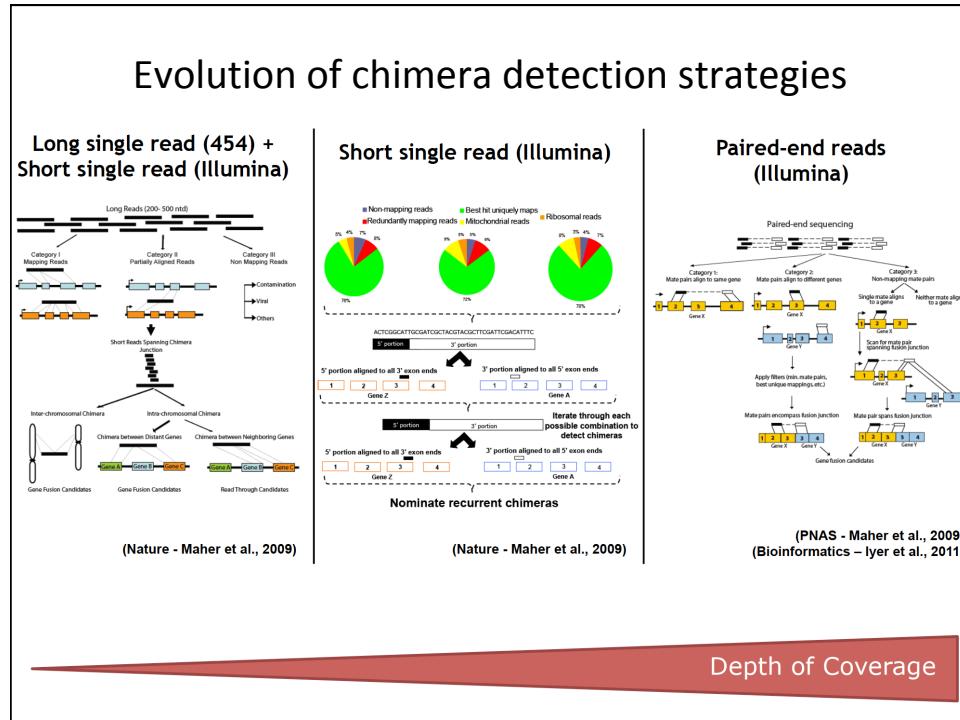
Gene fusions

I. Generation of a fusion gene



II. Activation of proto-oncogenes by relocation in the vicinity of active regulatory elements





Certain fusions have been detected in multiple cancers.

