

TopHat



Problem

- We map the 'reads' to the reference transcriptome
- Transcriptomes are incomplete even for well-studied species



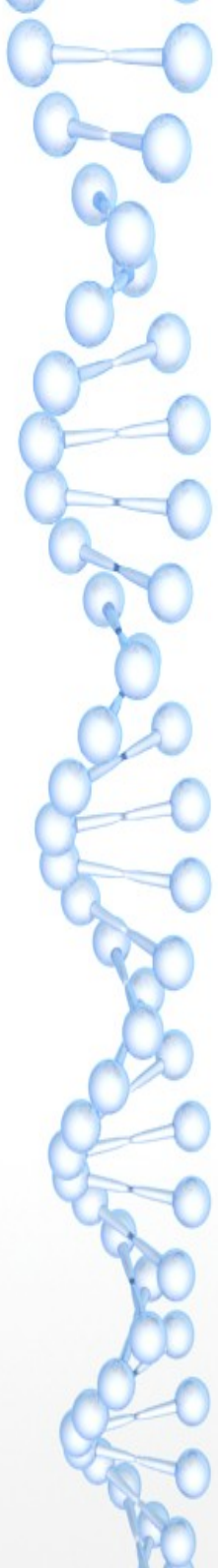
TopHat pipeline

- First map non-junction reads (those contained within exons) using Bowtie
- Extracts the sequences for the resulting islands of contiguous sequence from the sparse consensus, inferring them to be putative exons

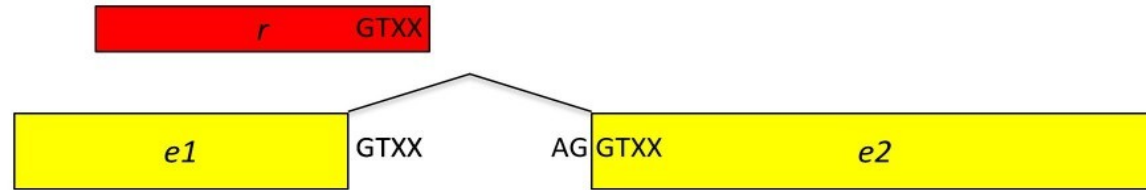


TopHat pipeline

- TopHat includes a small amount of flanking sequence from the reference on both sides of each island (default=45 bp)
- Consider all pairings of these sites that could form canonical (GT-AG)
- TopHat has a parameter that controls when two distinct but nearby exons should be merged into a single exon. (default is 6 bp)
- By default potential introns longer than 70 bp and shorter than 20 000 bp

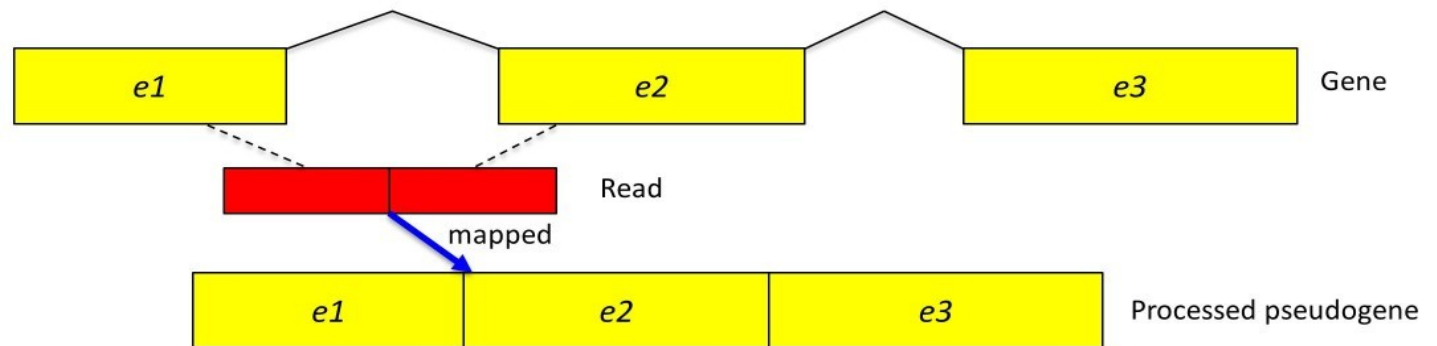


Incorrect mapping (non-gapped alignment)

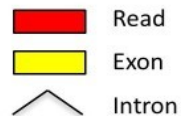


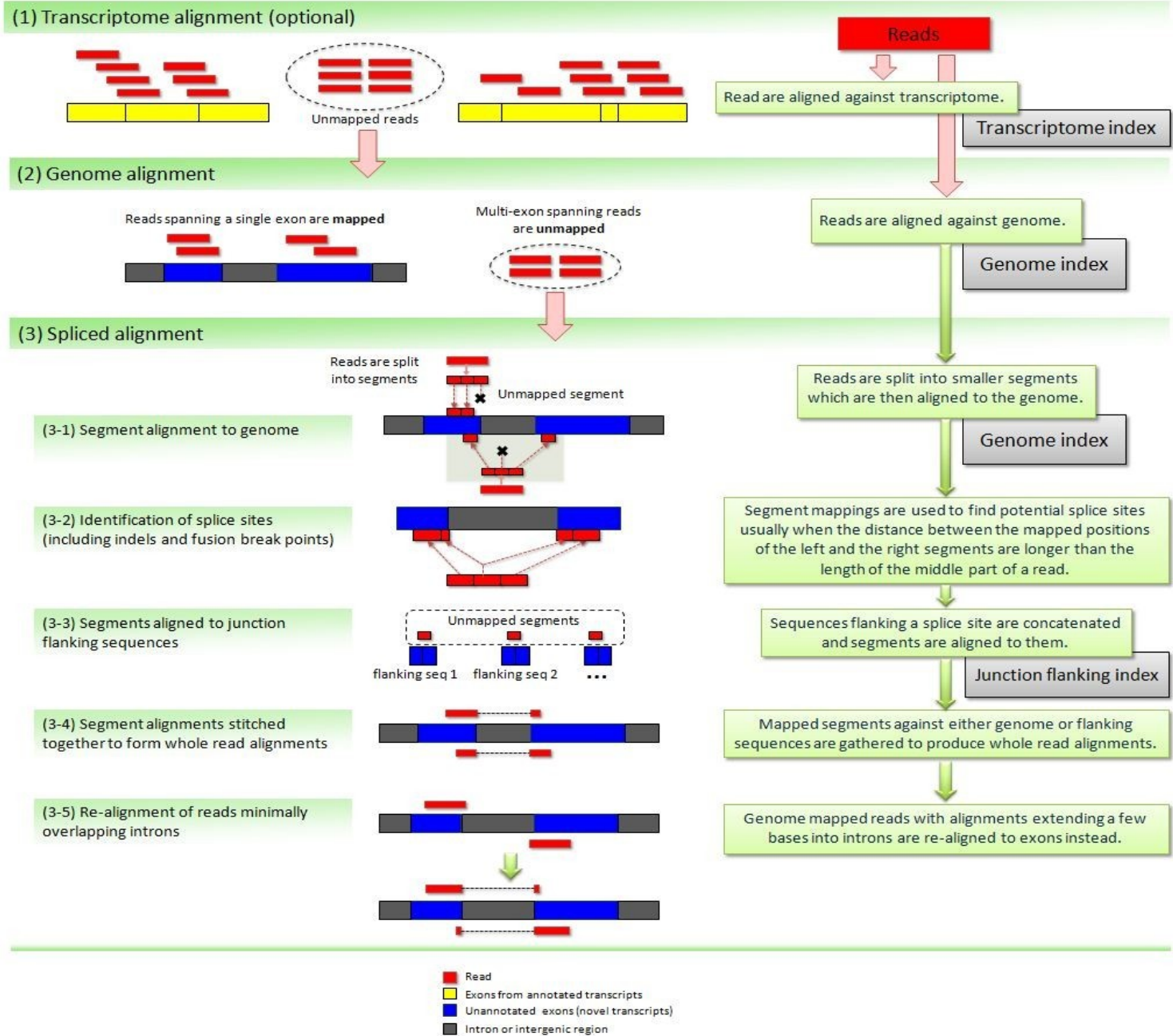
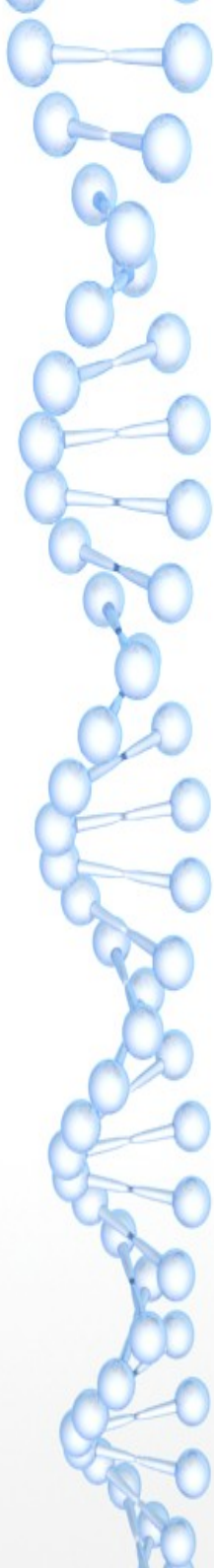
Correct mapping (spliced alignment)

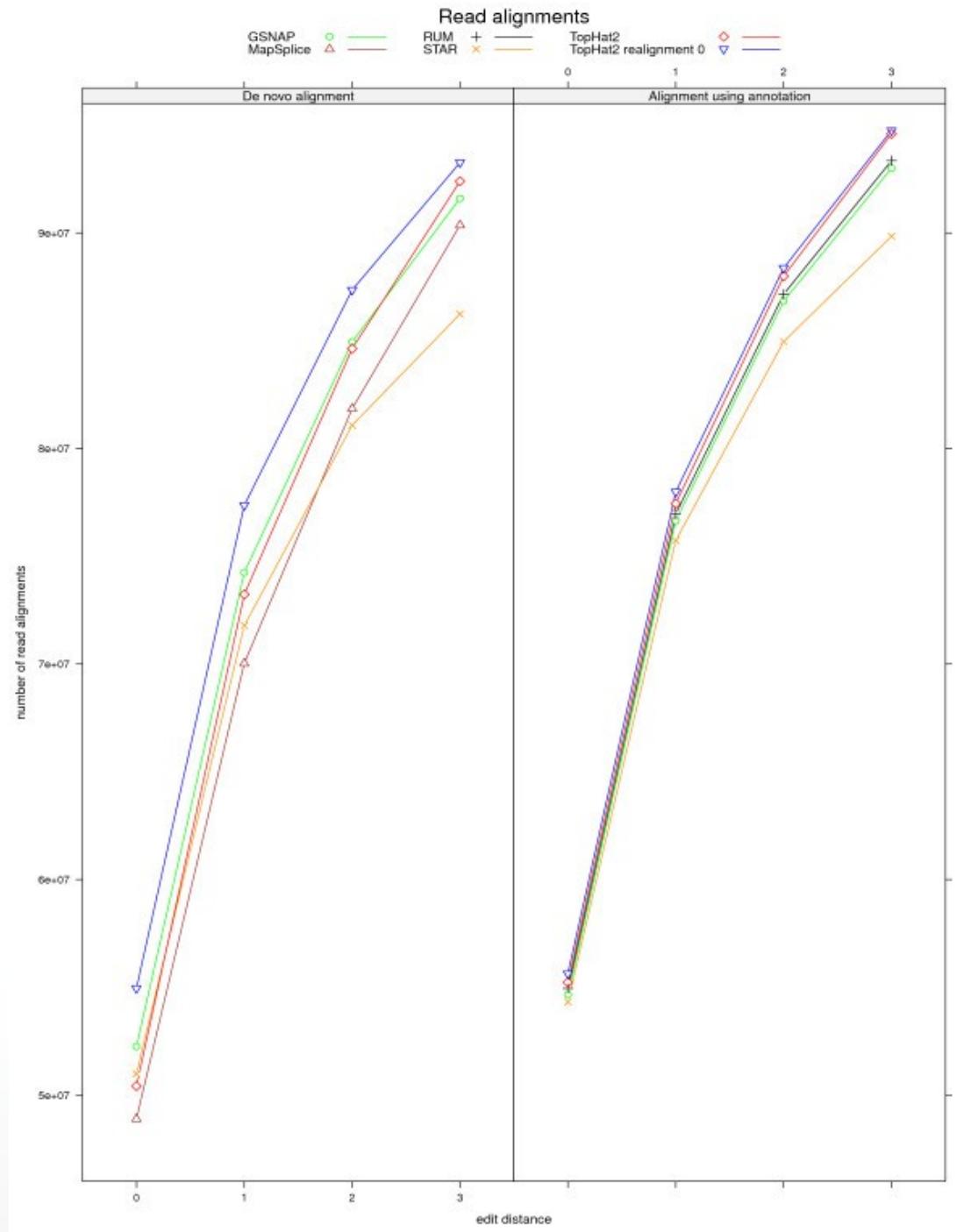
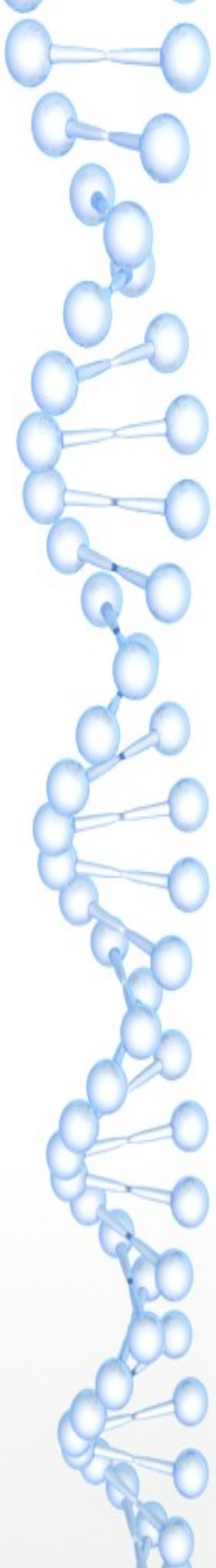
(1) Read r may be incorrectly mapped to the intron between exons $e1$ and $e2$.

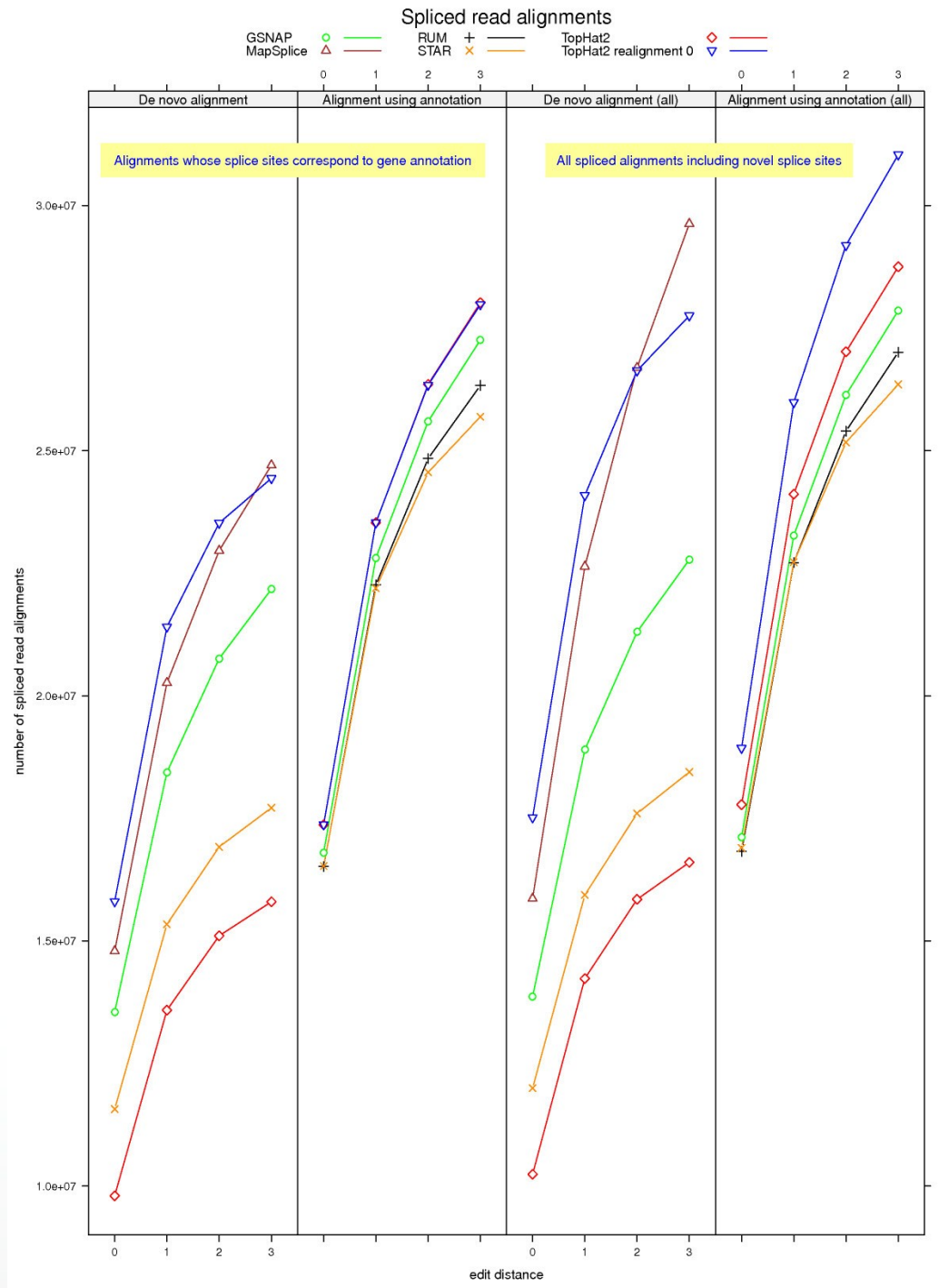
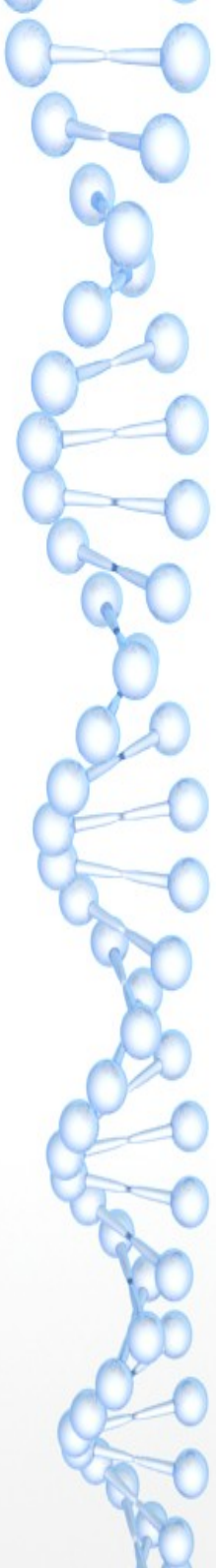


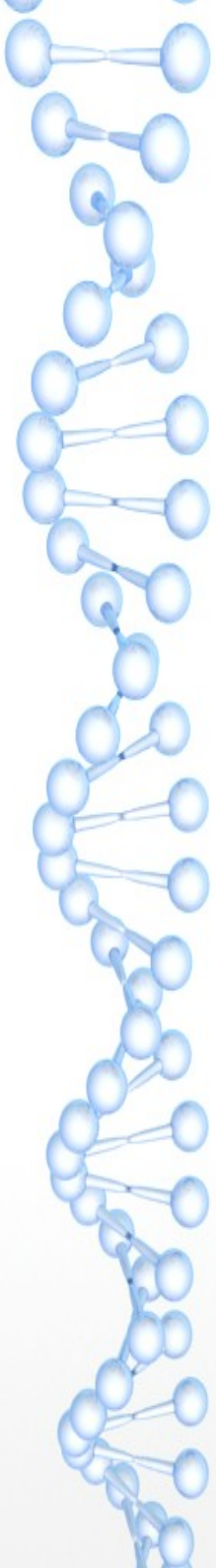
(2) Here, the read shown in red, which spans a splice junction, can be aligned end-to-end to a processed pseudogene.











Cufflinks



Problem

- Transcript assembly
- Transcript abundance estimation



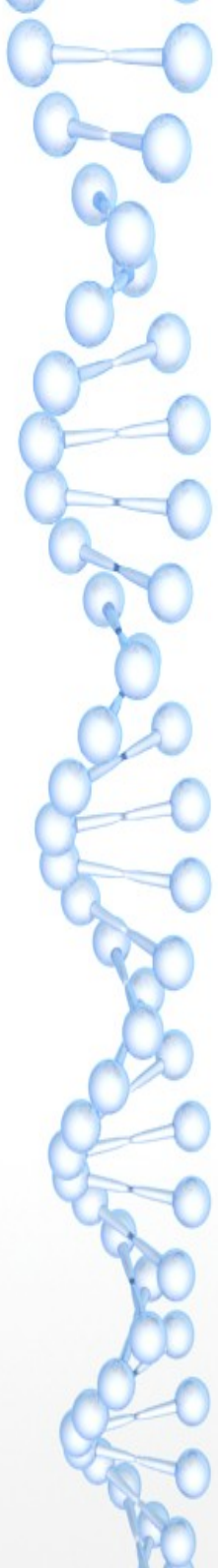
Assembling

- 1) Every fragment is consistent with at least one assembled transcript.
- 2) Every transcript is tiled by reads.
- 3) The number of transcripts is the smallest required to satisfy requirement (1).

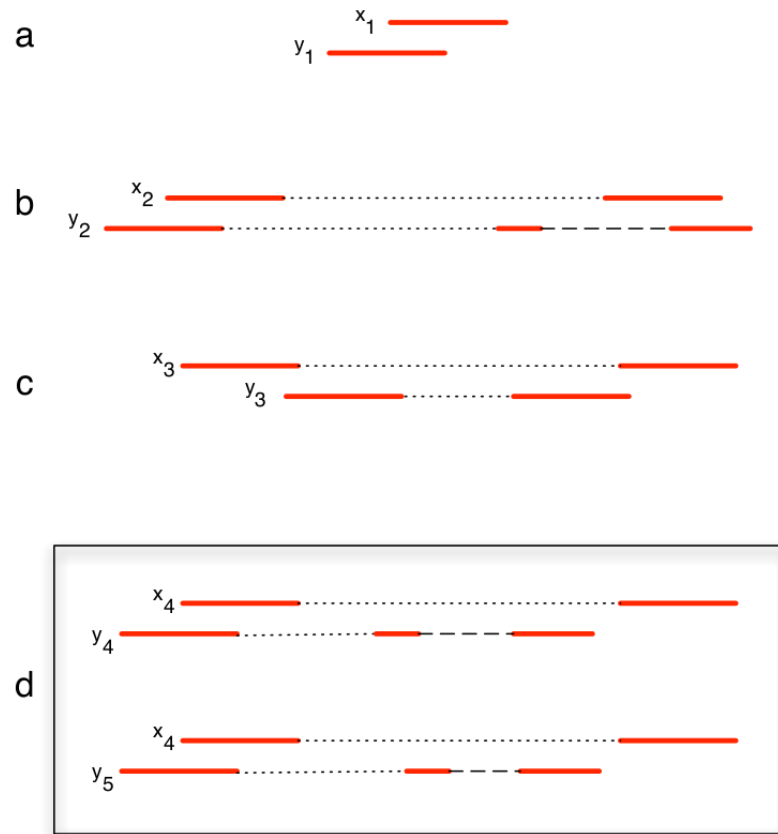


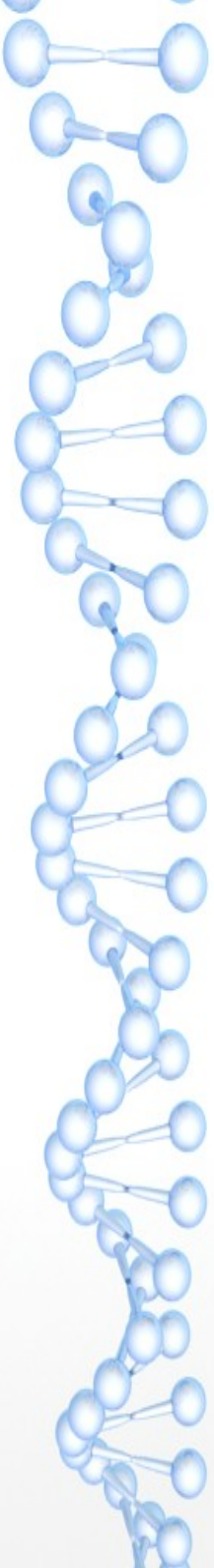
Assembling

- Fragments were first divided into non-overlapping loci, and each locus was assembled independently of the others
- Two reads are compatible if their overlap contains the exact same implied introns (or none)
- Each fragment alignment was assigned a node in an “overlap graph” G .
- A directed edge (x, y) was placed between nodes x and y when the alignment for x started at a lower coordinate than y , the alignments overlapped in the genome, and the fragments were “compatible”

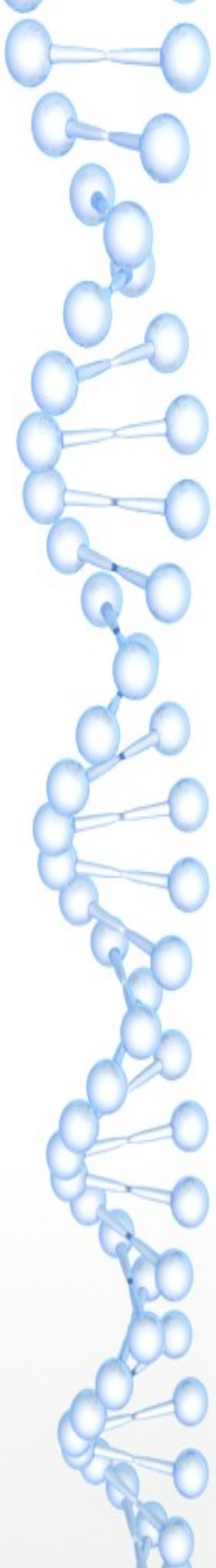


Compatible reads

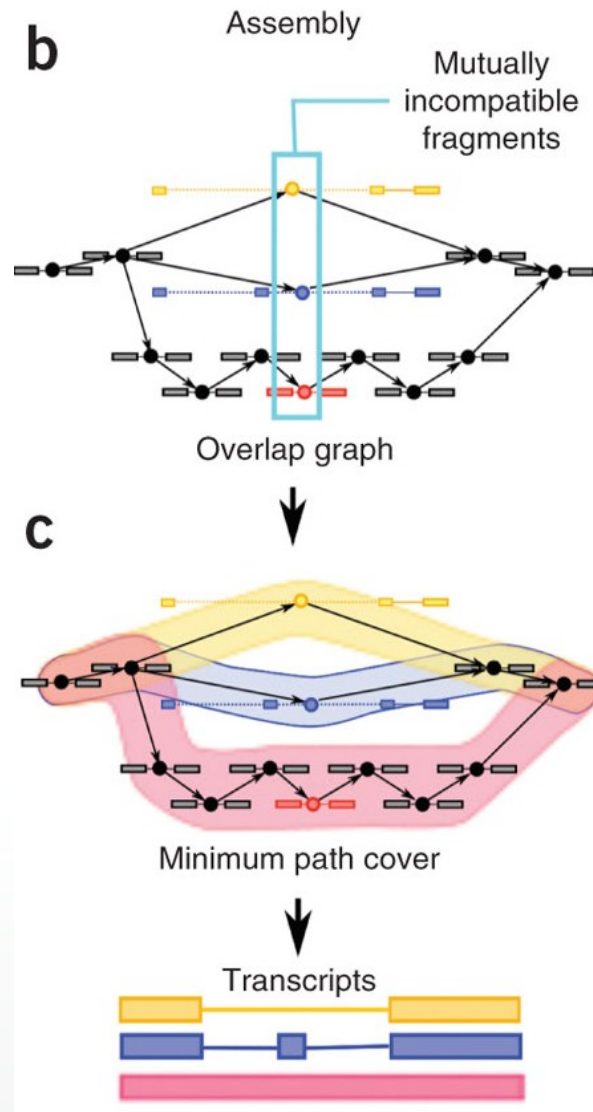




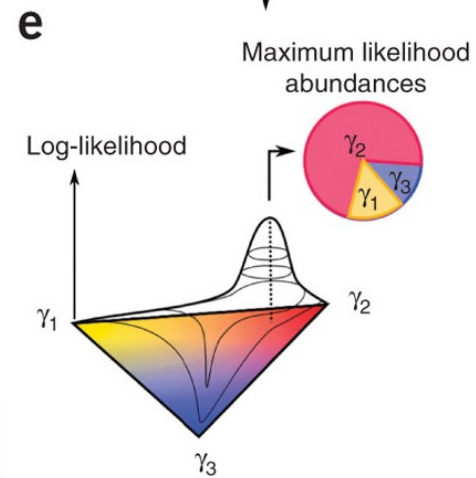
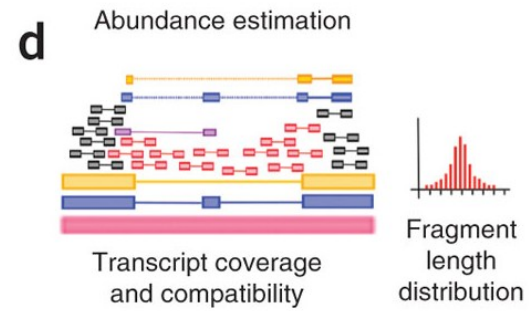
- Cufflinks then found a minimum path cover of G , meaning that every fragment node was contained in some path in the cover, and the cover contained as few paths as possible.



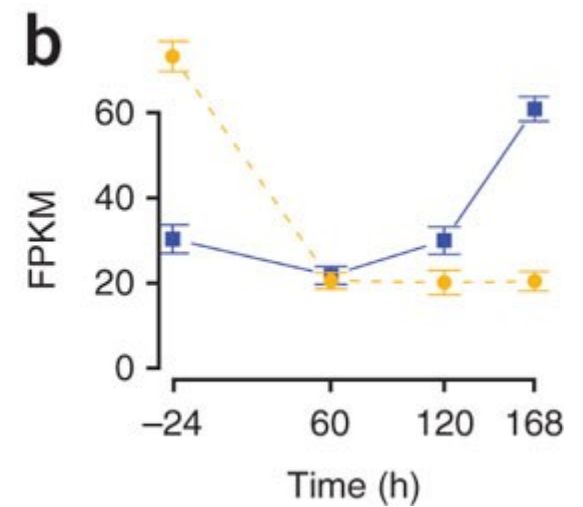
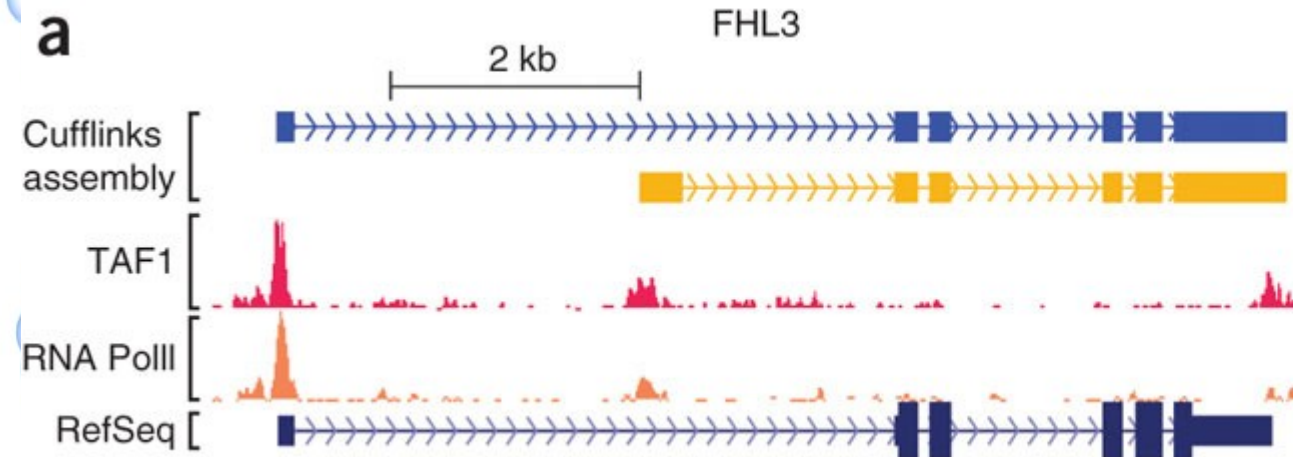
Assembly



Abundance



New FHL3 isoform





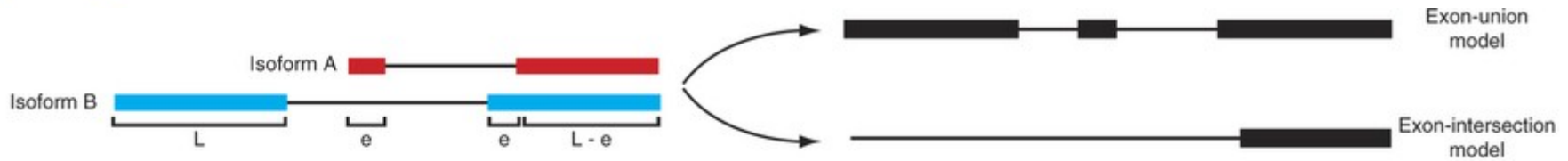
Cuffdiff 2

Current RNA-seq differential analysis methods focus on tackling one of two major challenges:

- The first is accurately deriving gene and isoform expression values from raw sequencing reads, which requires statistical computations at isoform-level resolution.
- The second is accounting for variability in measurements across biological replicates of an experiment.

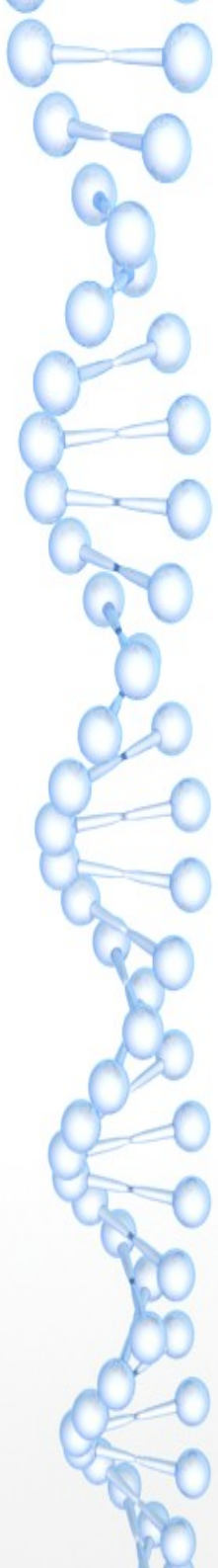
Read counting problem

a

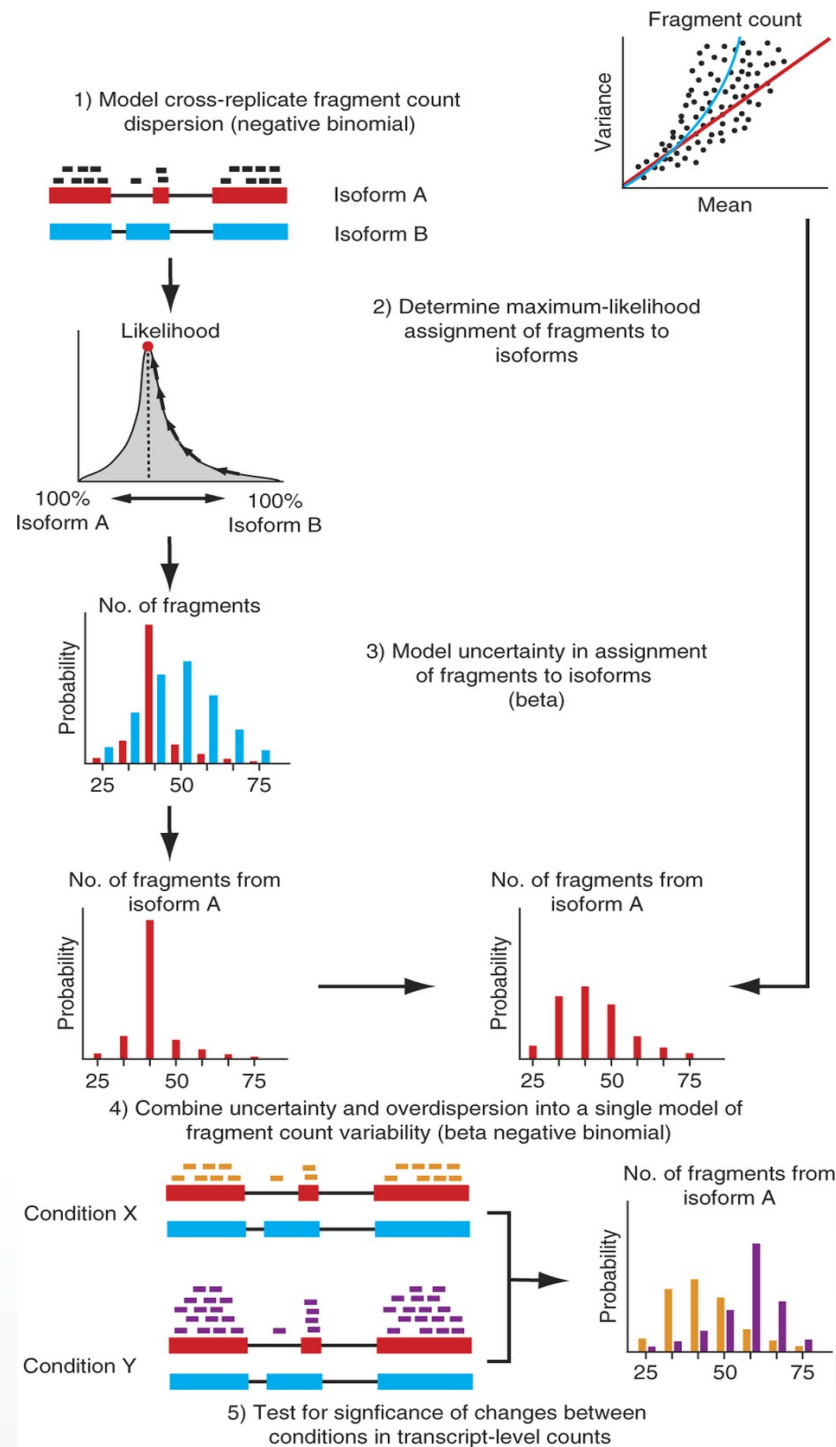
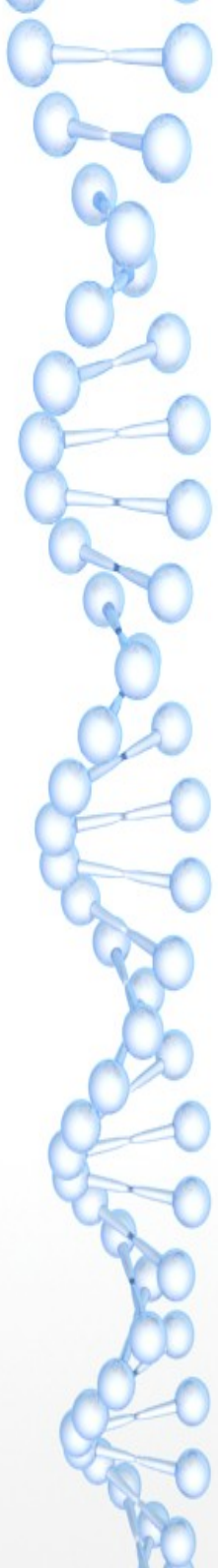


b

Condition A	Condition B	Log fold-change (union count)	Log fold-change (intersect count)	Log fold-change (true expression)
		$\log_2\left(\frac{10}{10}\right) = 0$	$\log_2\left(\frac{8}{7}\right) = 0.19$	$\log_2\left(\frac{10/L}{6/L + 4/2L}\right) = 0.32$
		$\log_2\left(\frac{6}{8}\right) = -0.41$	$\log_2\left(\frac{5}{5}\right) = 0$	$\log_2\left(\frac{6/L}{8/2L}\right) = 0.58$
		$\log_2\left(\frac{5}{10}\right) = -1$	$\log_2\left(\frac{4}{5}\right) = -0.1$	$\log_2\left(\frac{5/L}{10/2L}\right) = 0$



- Count uncertainty refers to the observation that in RNA-seq experiments it is common for up to 50% of reads to map ambiguously to different transcripts





Cuffdiff 2

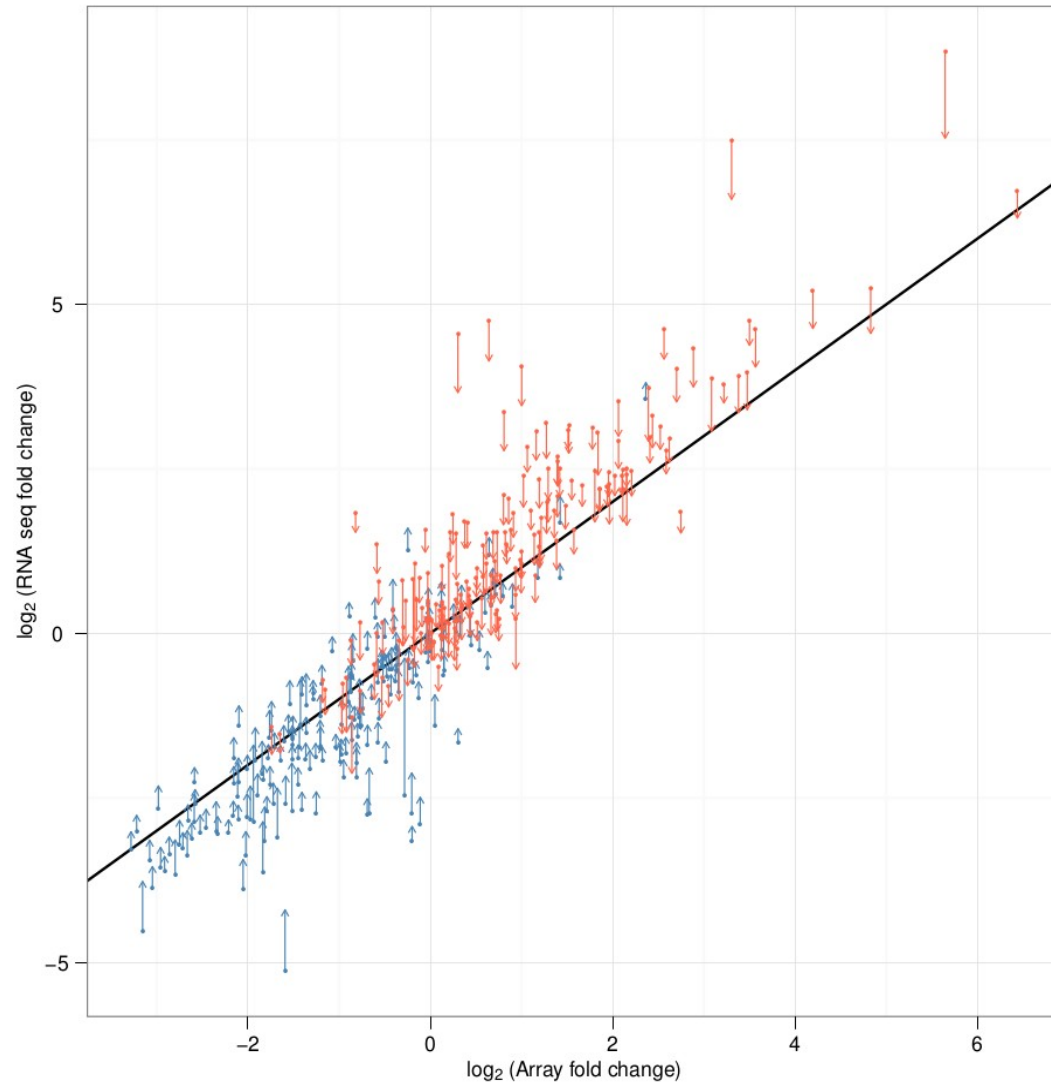
- Cuffdiff 2 estimates expression at gene- and transcript-level resolution
- Variance in the expression levels
- covariances between isoforms of the same gene from replicate experiments.
- This allows it to accurately estimate gene expression and perform differential analysis at gene-level



Response to loss of HOXA1

- Cuffdiff 2 – derived changes in gene expression in response to HOXA1 knockdown strongly agreed with values from microarrays
- Cuffdiff 2 improved concordance with the array measurements by 15% compared with the change in raw count

Response to loss of HOXA1





Gentleman's set

- Samtools
- Bowtie 2
- Bowtie 2 index
- Annotations known transcripts (GFF)
- TopHat 2
- Cufflinks 2



Articles

- TopHat: discovering splice junctions with RNA-Seq
- TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions
- Transcript assembly and abundance estimation from RNA-Seq reveals thousands of new transcripts and switching among isoforms
- Differential analysis of gene regulation at transcript resolution with RNA-seq