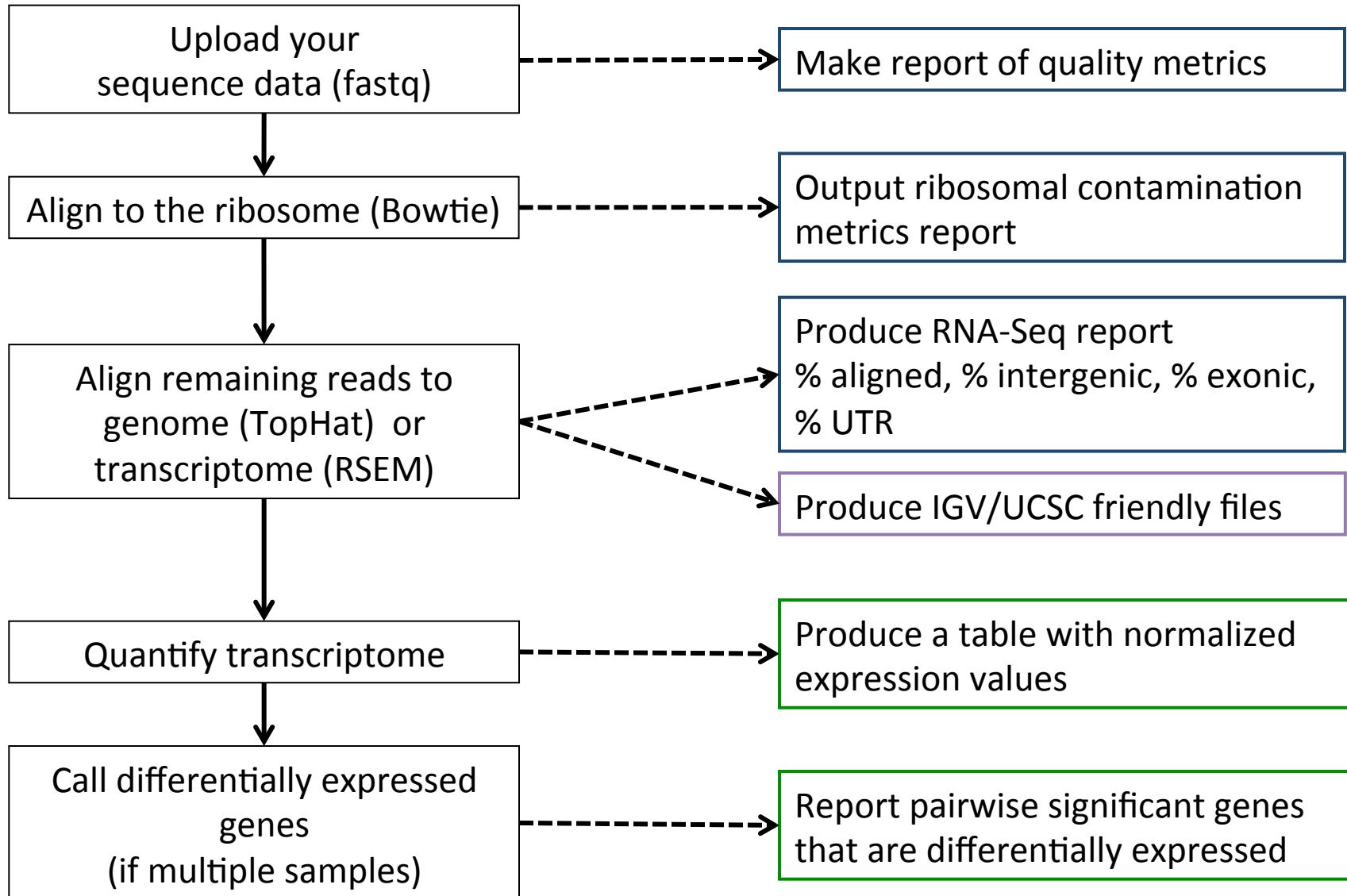


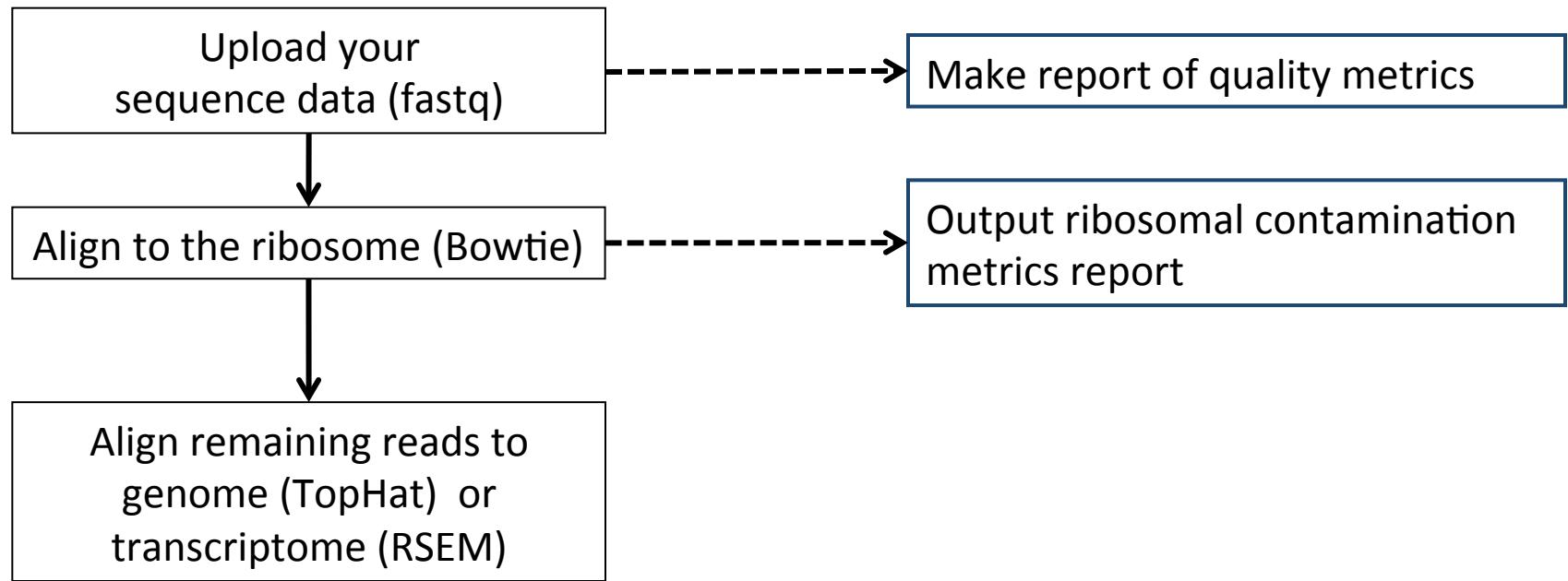


RNA-Seq primer

Our typical RNA quantification pipeline

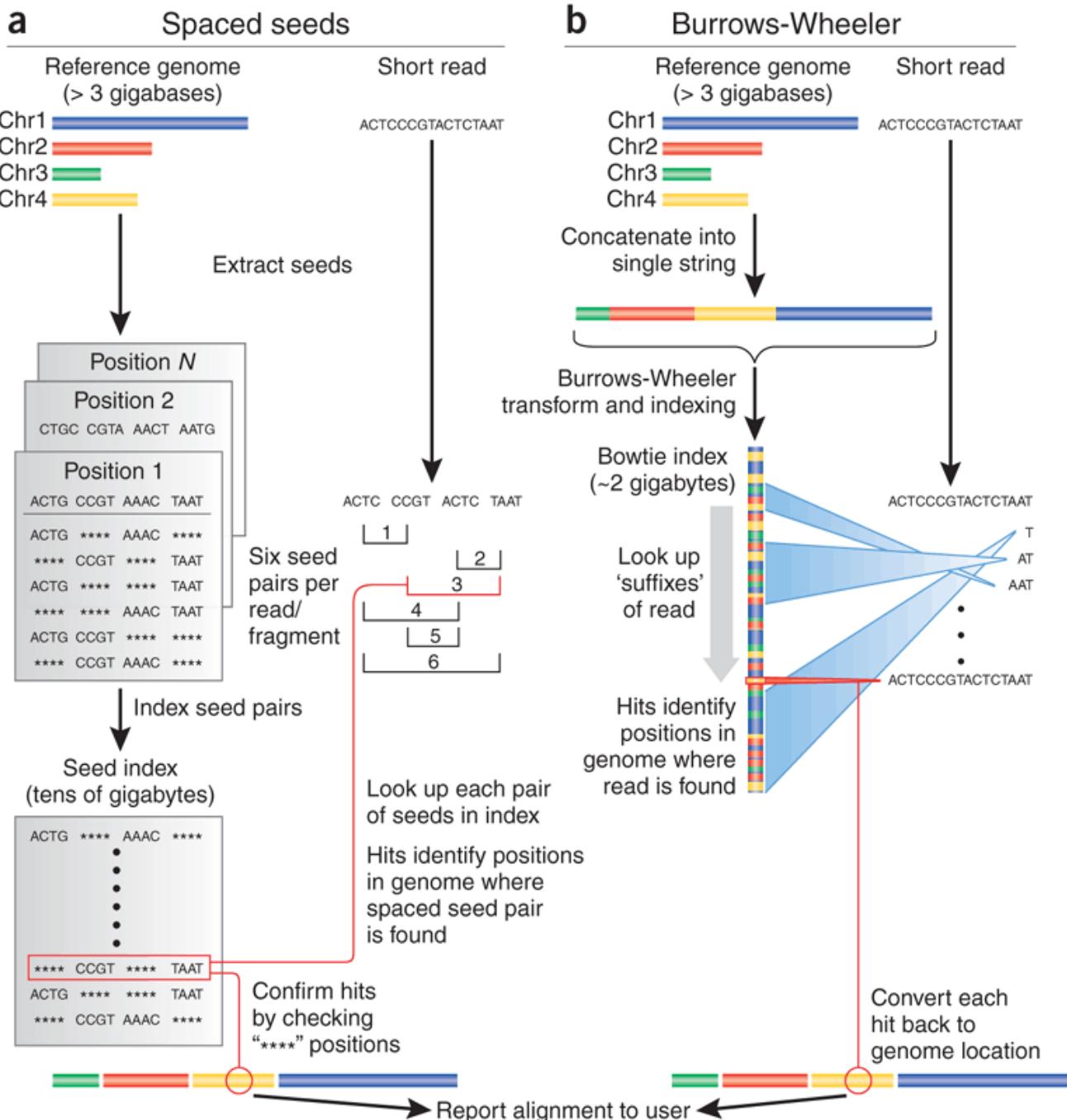


Alignment requires pre-processing



```
bowtie2-build -f mm10.fa mm10
```

```
rsem-prepare-reference \
--gtf ucsc.gtf --transcript-to-gene-map ucsc_into_genesymbol.rsem \
mm10.fa mm10.rsem
```



Mapping quality

What we want to estimate is $q_{MS} = -10 \log_{10} P(\text{read is wrongly mapped})$

That is, the posterior probability, the probability that the region starting at i was sequenced *given* that we observed the read q :

$$P(i|G, q) = \frac{P(q|G, i)P(i|G)}{P(q|G)} = \frac{P(q|G, i)P(i|G)}{\sum_j P(q|G, j)}$$

Fortunately, there are efficient ways to approximate this probability (see Li, H *genome Research* 2008, for example)

$$q_{MS} = -10 \log_{10} (1 - P(i|G, q))$$

Considerations

- Trade-off between sensitivity, speed and memory
 - Smaller seeds allow for greater mismatches at the cost of more tries
 - Smaller seeds result in a smaller tables (table size is at most 4^k), larger seeds increase speed (less tries, but more seeds)

Short read mapping software

Seed-extend

	Short indels	Use base qual
Maq	No	YES
RMAP	Yes	YES
SeqMap	Yes	NO
SHRiMP	Yes	NO

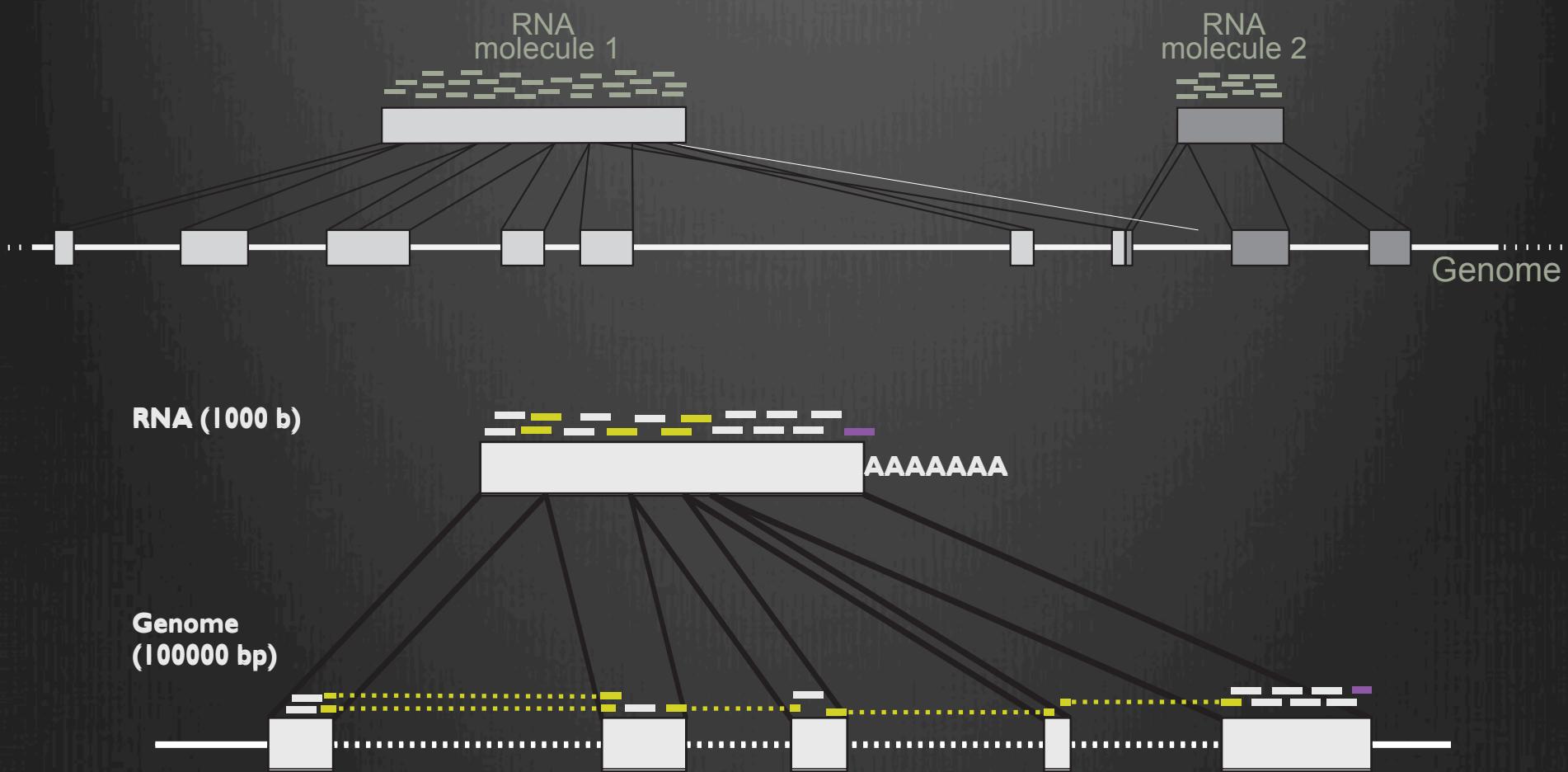
BWT

	Use Base qual
BWA	YES
Bowtie	NO
Stampy*	YES
Bowtie2*	(NO)

*Stampy is a hybrid approach which first uses BWA to map reads then uses seed-extend only to reads not mapped by BWA

*Bowtie2 breaks reads into smaller pieces and maps these “seeds” using a BWT genome.

RNA-Seq Read mapping



Mapping RNA-Seq reads: Exon-first spliced alignment (e.g. TopHat)



Mapping RNA-Seq reads: Seed-extend spliced alignment (e.g. GSNAP)



Short read mapping software for RNA-Seq

Seed-extend

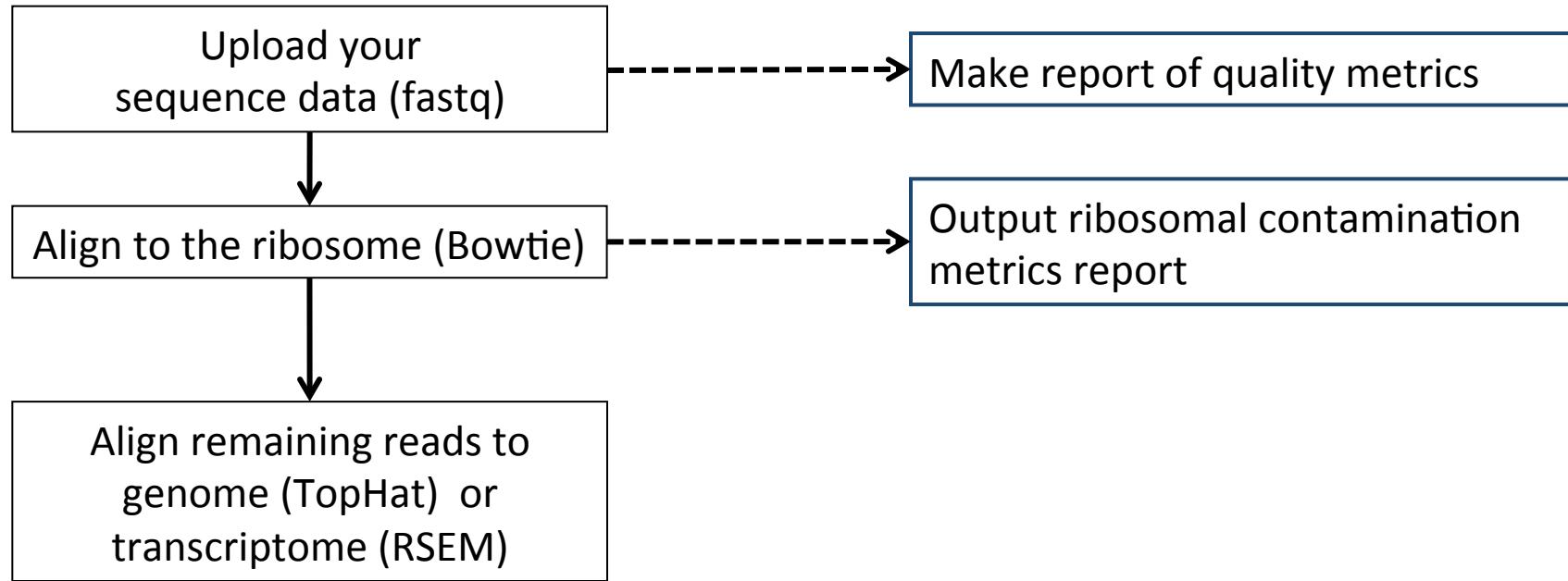
	Short indels	Use base qual
GSNAP	Yes	?
QPALMA	Yes	NO
BLAT	Yes	NO

Exon-first

	Use base qual
STAR	NO
TopHat	NO

Exon-first alignments will map contiguous first at the expense of spliced hits

Short read alignment



```
tophat2 --library-type fr-firststrand --segment-length 20 \  
-G  genome.quantification/ucsc.gtf -o tophat/th.quant.ctrl1 \  
genome.quantification/mm10 fastq.quantification/control_rep1.1.fq \  
fastq.quantification/control_rep1.2.fq
```

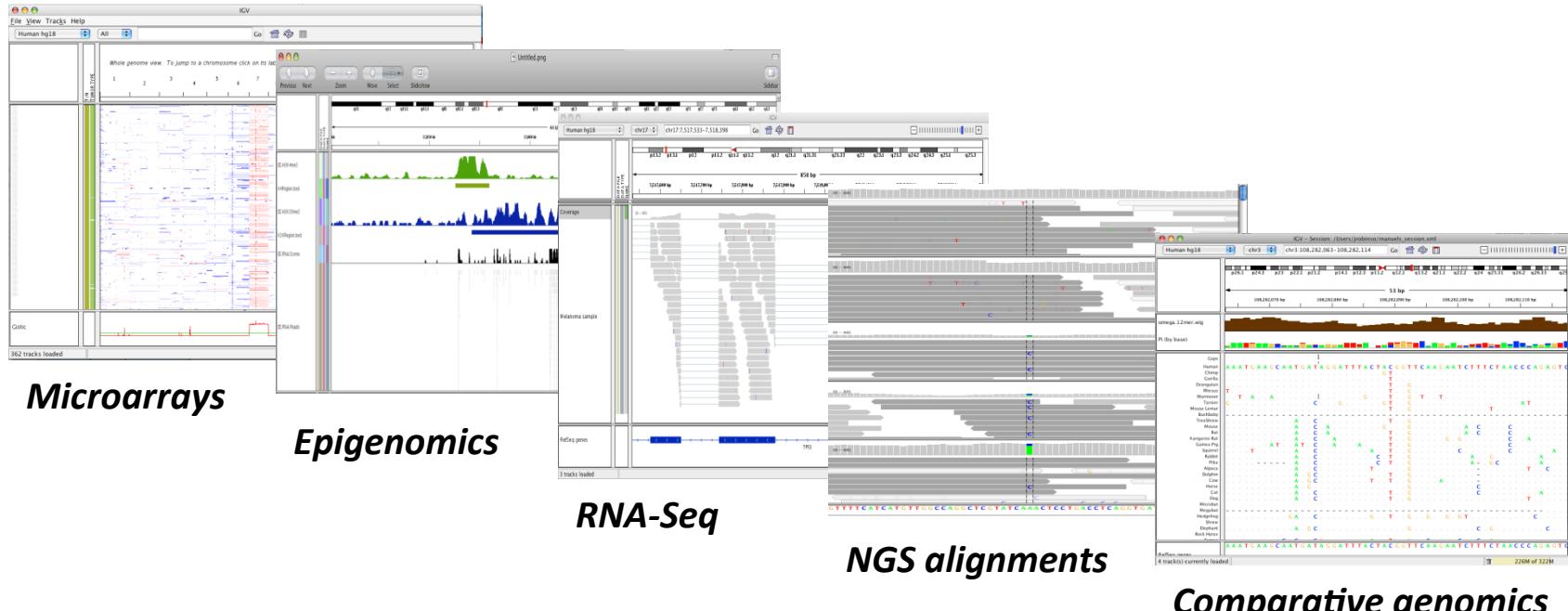
```
/project/umw_biocore/bin/igvtools.sh count -w 5 tophat/th.quant.ctrl1.bam \  
tophat/th.quant.ctrl1.bam.tdf genome.quantification/mm10.fa
```

IGV: Integrative Genomics Viewer

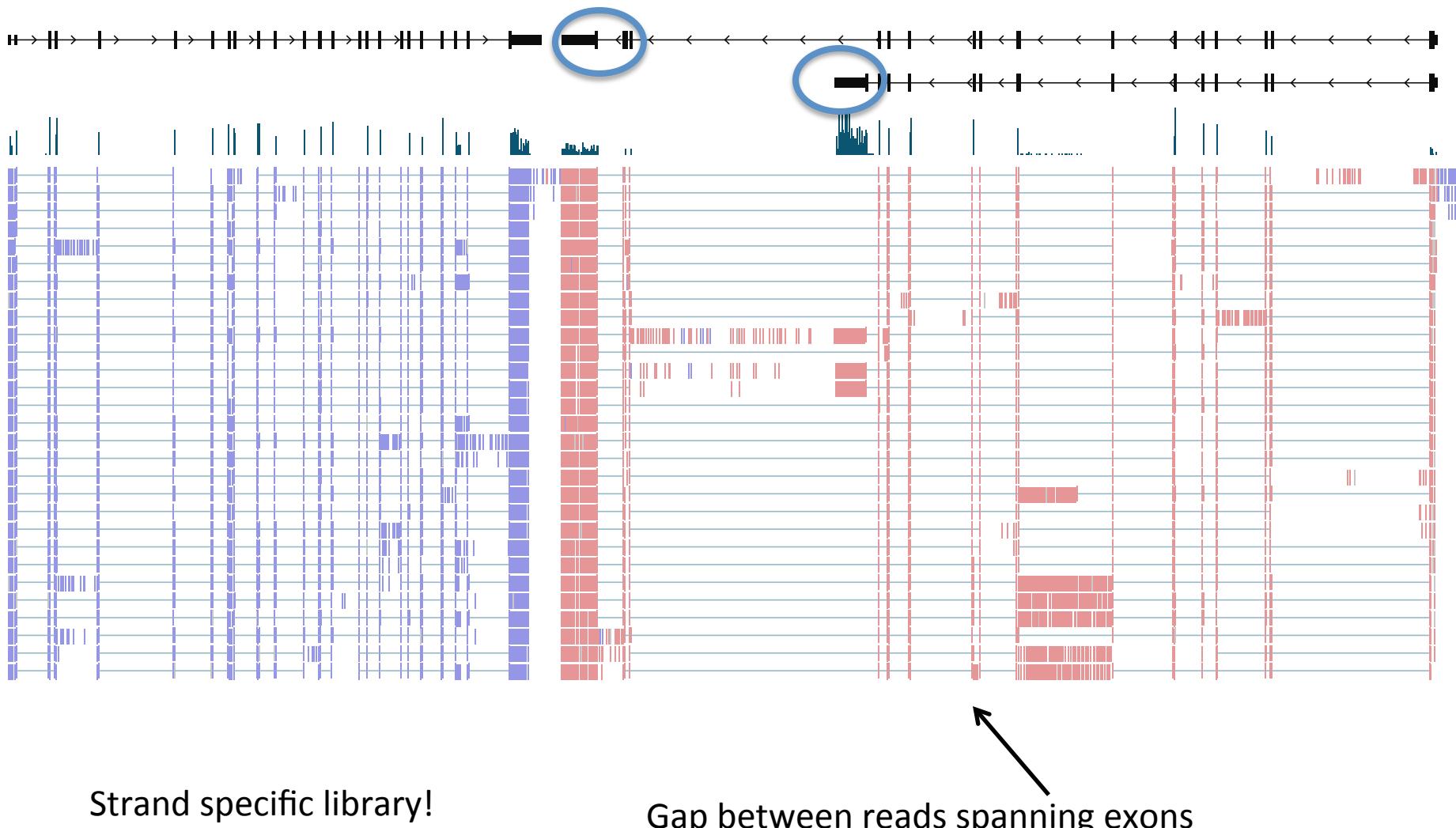


A desktop application

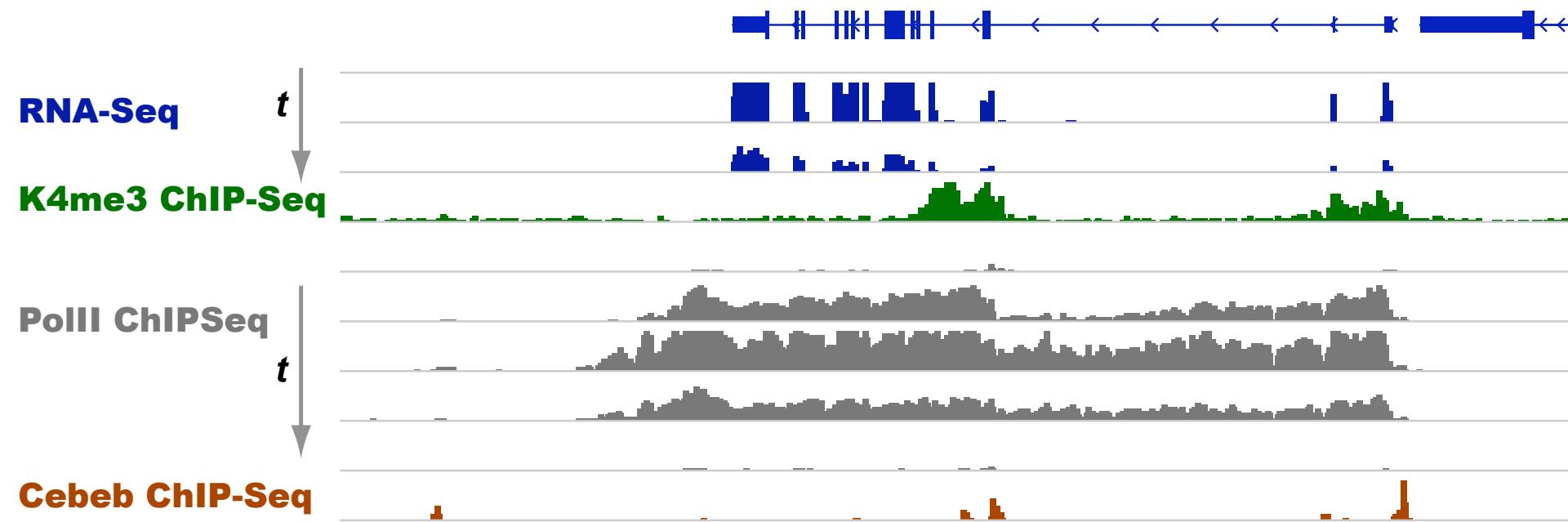
for the visualization and interactive exploration
of genomic data



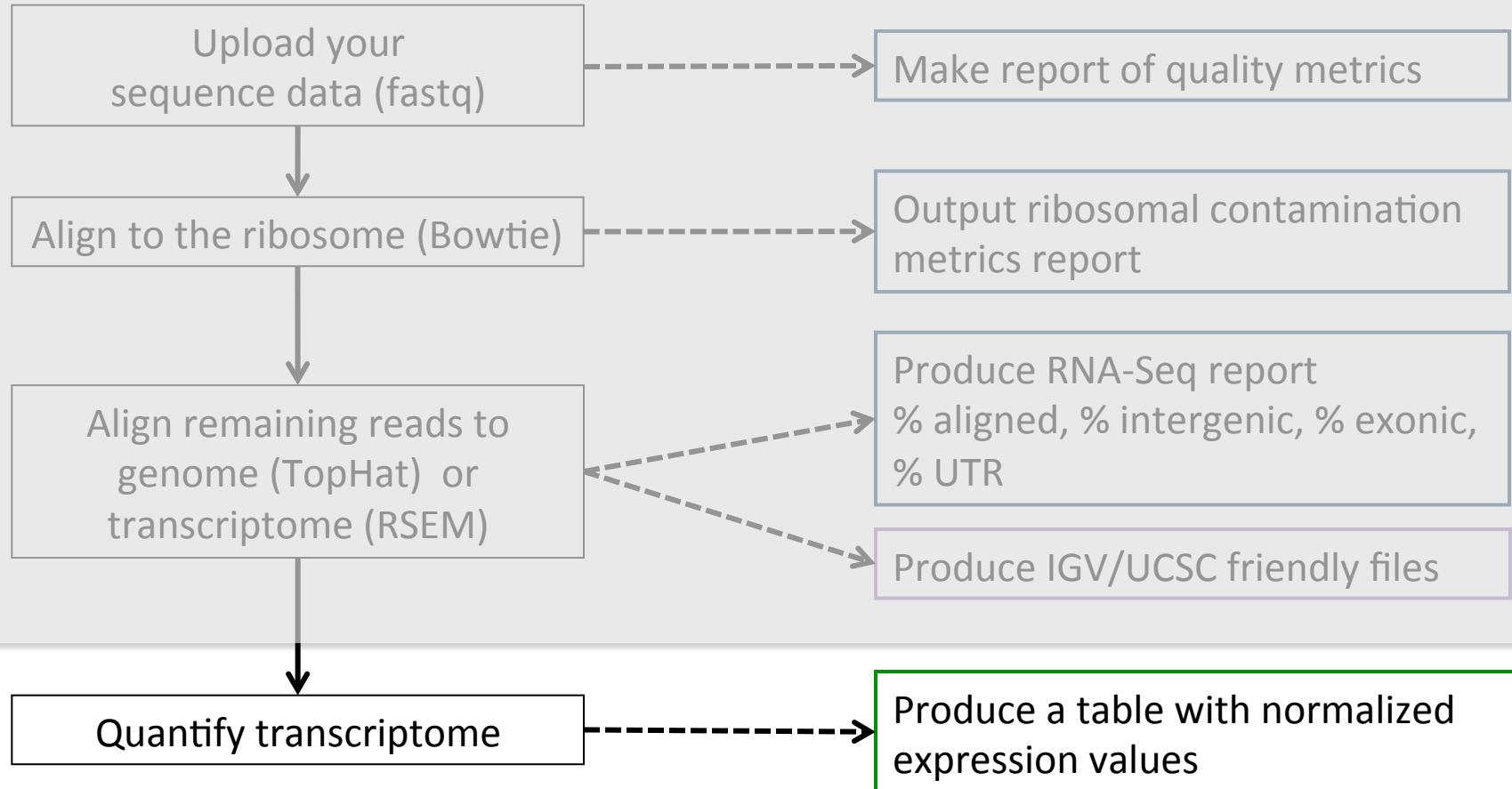
Visualizing read alignments with IGV — RNASeq



Visualizing read alignments with IGV — zooming out



Computing gene expression

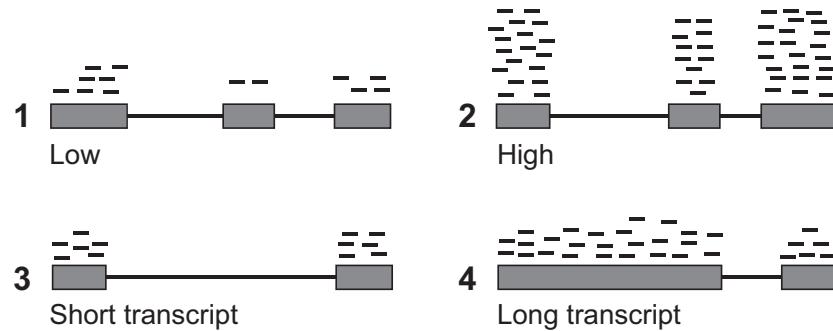


```
rsem-calculate-expression --paired-end --strand-specific -p 2 \  
--output-genome-bam fastq.quantification/control_rep1.1.fq \  
fastq.quantification/control_rep1.2.fq genome.quantification/mm10.rsem \  
rsem/ctrl1.rsem
```

RNA-Seq quantification

- Is a given gene (or isoform) expressed?
- Is expression gene A > gene B?
- Is expression of gene A isoform a_1 > gene A isoform a_2 ?
- Given two samples is expression of gene A in sample 1 > gene A in sample 2?

Quantification: only one isoform



$$RPKM = 10^9 \frac{\#reads}{length \times Total\,Reads}$$

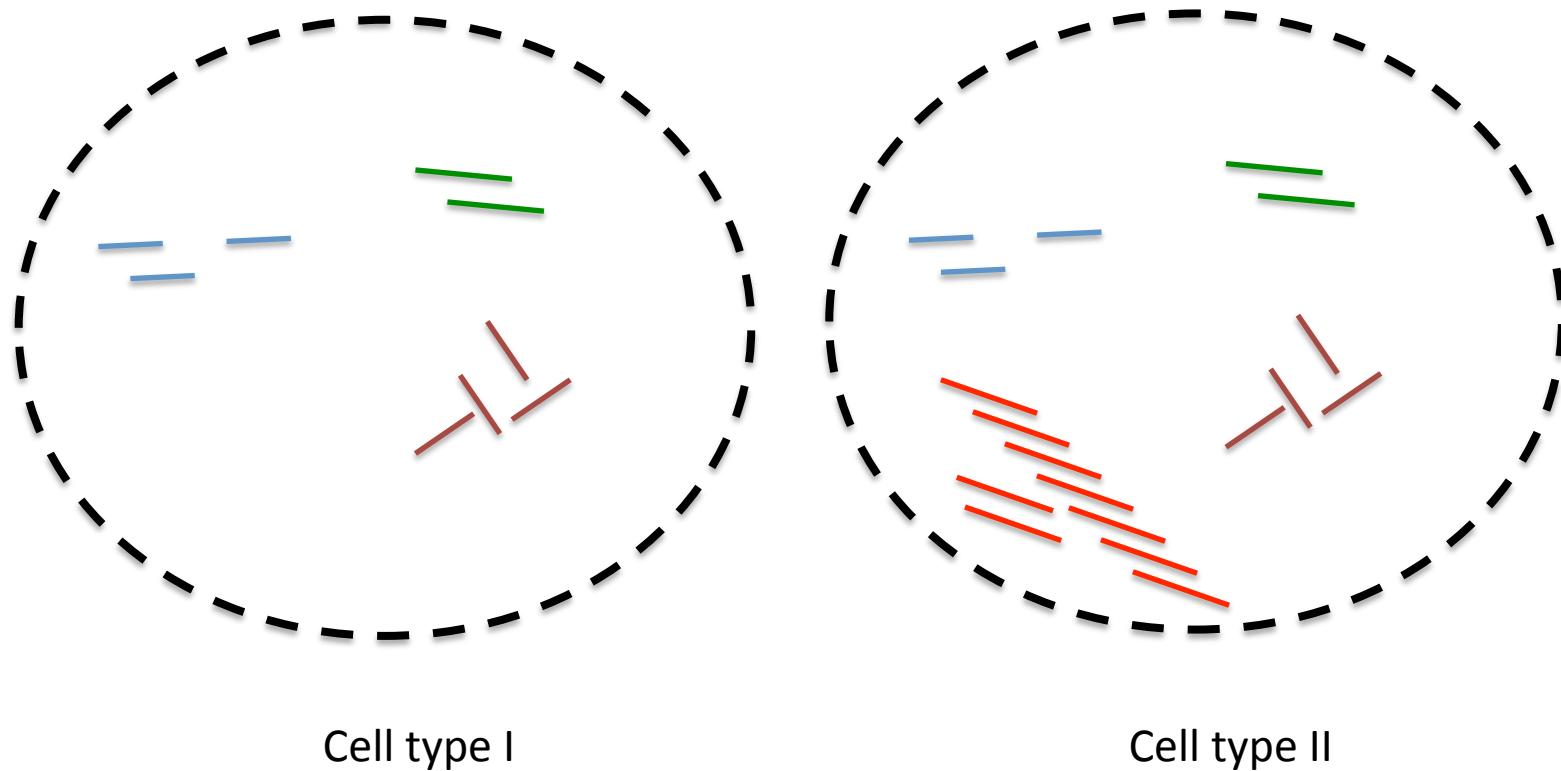
Reads per kilobase of exonic sequence per million mapped reads
(Mortazavi et al Nature methods 2008)

- Fragmentation of transcripts results in length bias: longer transcripts have higher counts
- Different experiments have different yields. Normalization is key for cross lane comparisons

Normalization for comparing two different genes

- To compare within a sequence run (lane), RPKM accounts for length bias.
- RPKM is not optimal for cross experiment comparisons.
 - Different samples may have different compositions.
- FPKM superseded RPKM
- And later TPM = $10^6 \times$ Fraction of transcript

Normalization for comparing a gene across samples



Normalizing by total reads does not work well for samples with very different RNA composition

Step2: More robust normalization

$$s_j = \text{median}_i \frac{k_{ij}}{\left(\prod_{v=1}^m k_{iv} \right)^{1/m}}.$$

Counts for gene i in experiment j

Geometric mean for that gene
over ALL experiments

The diagram consists of two text boxes with arrows pointing to specific parts of the equation. The top box contains the text "Counts for gene i in experiment j" and has an arrow pointing to the term k_{ij} . The bottom box contains the text "Geometric mean for that gene over ALL experiments" and has an arrow pointing to the term $\left(\prod_{v=1}^m k_{iv} \right)^{1/m}$.

i runs through all n genes

j through all m samples

k_{ij} is the observed counts for gene i in sample j

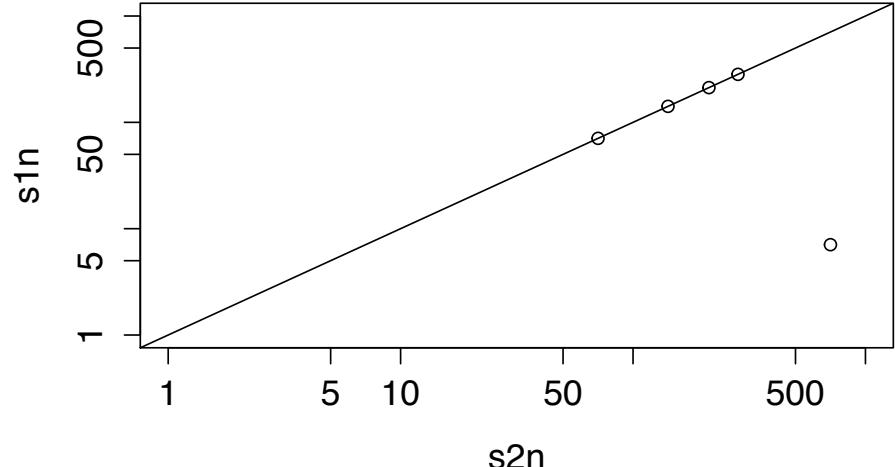
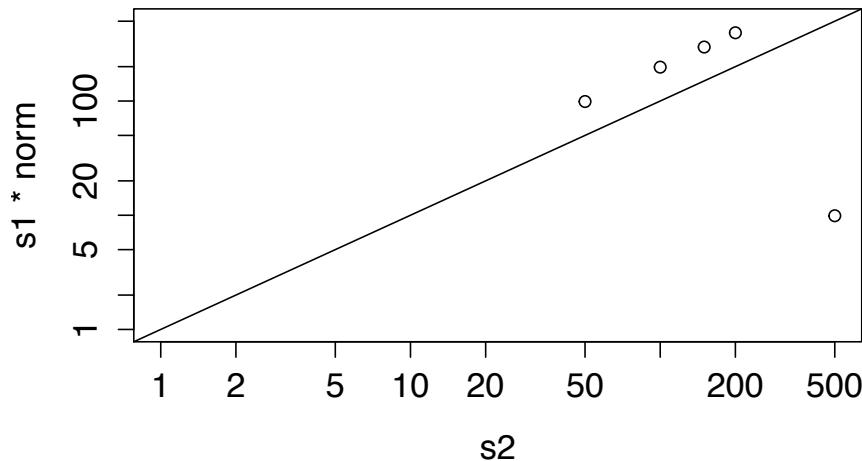
s_j is the normalization constant

Lets do an experiment (and do a short R practice)

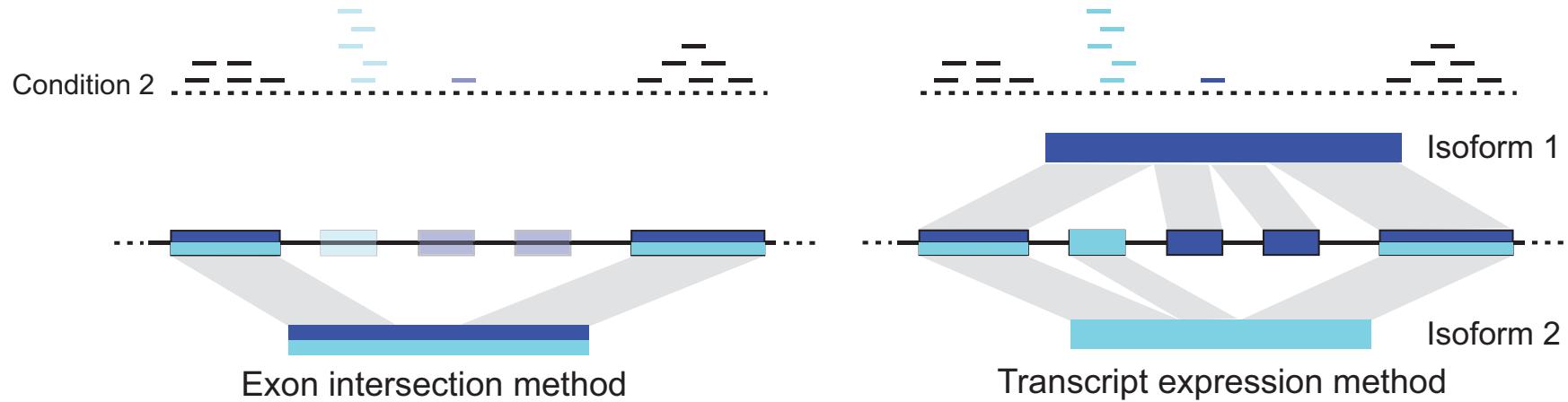
```
> s1 = c(100, 200, 300, 400, 10)  
> s2 = c(50, 100, 150, 200, 500)  
  
> norm=sum(s2)/sum(s1)  
> plot(s2, s1*norm,log="xy")  
> abline(a = 0, b = 1)  
  
> g = sqrt(s1 * s2t)  
> s1n = s1/median(s1/g); s2n = s2/median(s2/g)  
> plot(s2n, s1n,log="xy")  
> abline(a = 0, b = 1)
```

Similar read number,
one transcript many fold changed

Size normalization results in 2-fold
changes in *all* transcripts



But, how to compute counts for complex gene structures?



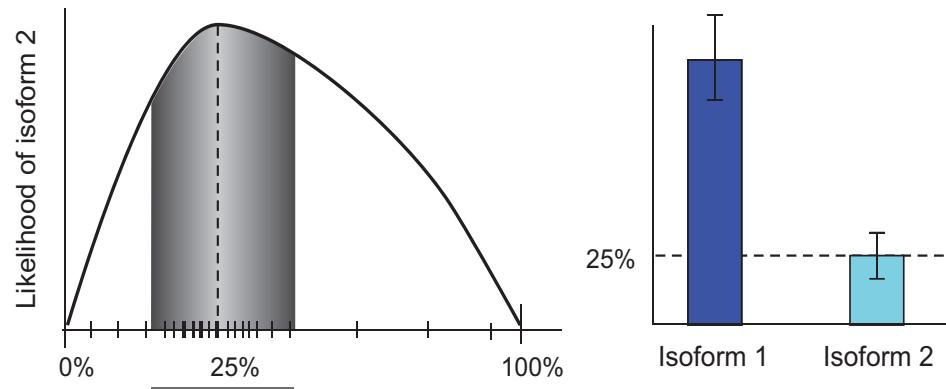
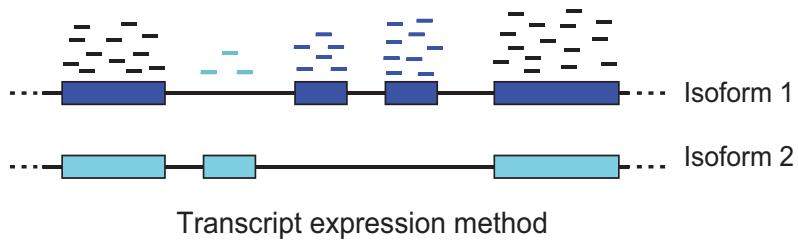
Three popular options:

Exon *intersection* model: Score constituent exons

Exon *union* model: Score the the “merged” transcript

Transcript expression model: Assign reads uniquely to different isoforms. *Not a trivial problem!*

Quantification: read assignment method

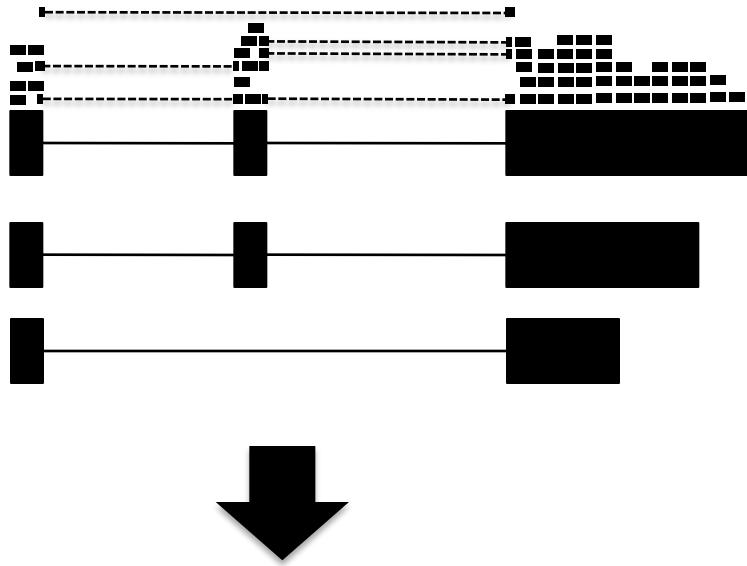


Quantification with multiple isoforms



How do we define the gene expression?
How do we compute the expression of each isoform?

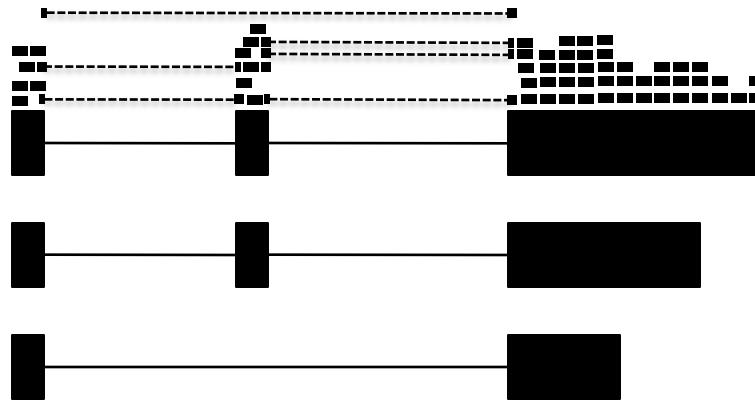
Computing gene expression



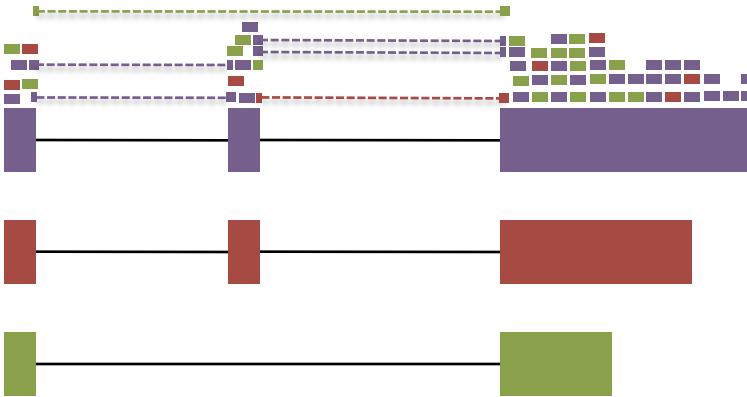
Idea1: RPKM of the
constitutive reads
(Neuma, Alexa-Seq,
Scripture)



Computing gene expression — isoform deconvolution



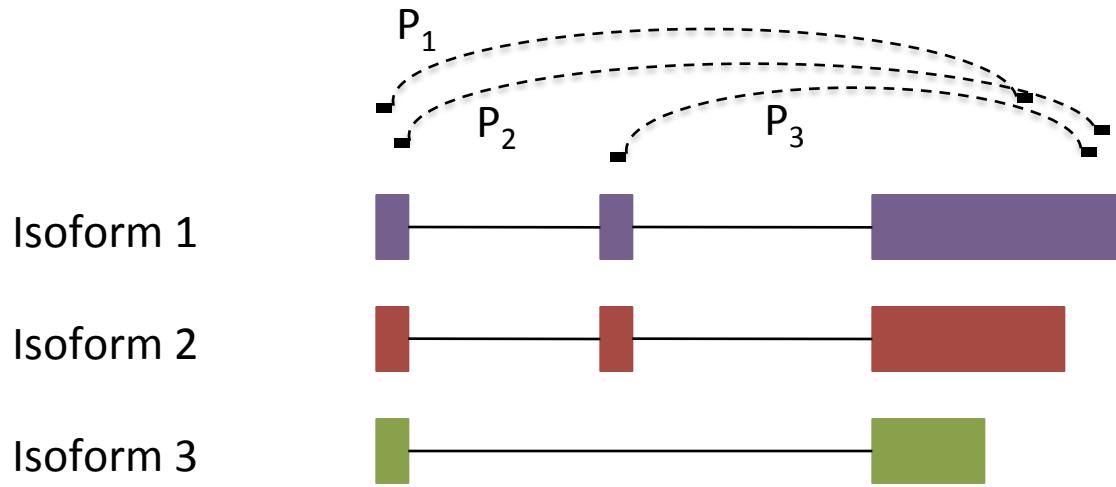
Computing gene expression — isoform deconvolution



If we knew the origin of the reads we could compute each isoform's expression. The gene's expression would be the sum of the expression of all its isoforms.

$$E = \text{RPKM}_1 + \text{RPKM}_2 + \text{RPKM}_3$$

Paired-end reads are easier to associate to isoforms

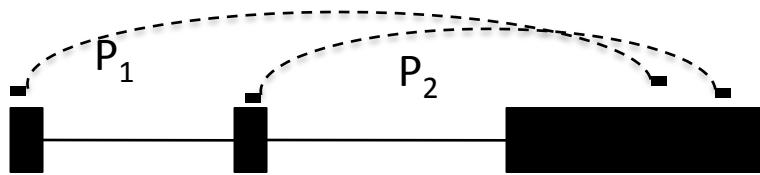


Paired ends increase isoform deconvolution confidence

- P₁ originates from isoform 1 or 2 but not 3.
- P₂ and P₃ originate from isoform 1

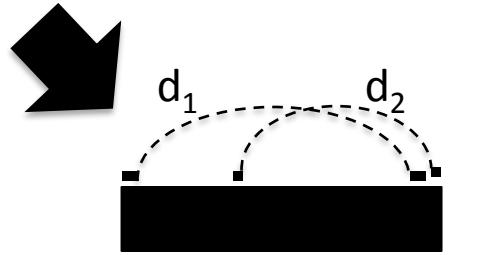
Do paired-end reads also help identifying reads originating in isoform 3?

We can estimate the insert size distribution

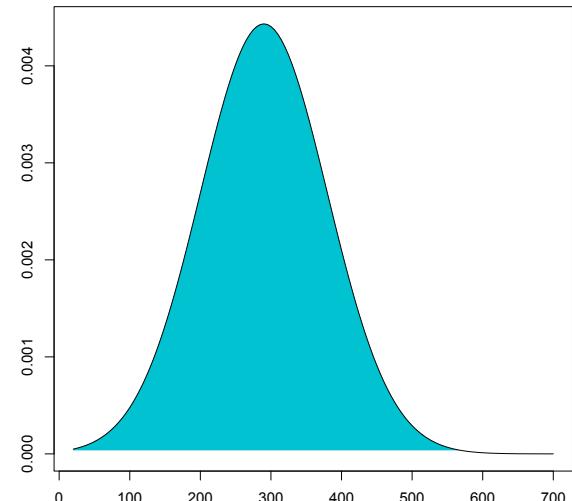


Get all single isoform reconstructions

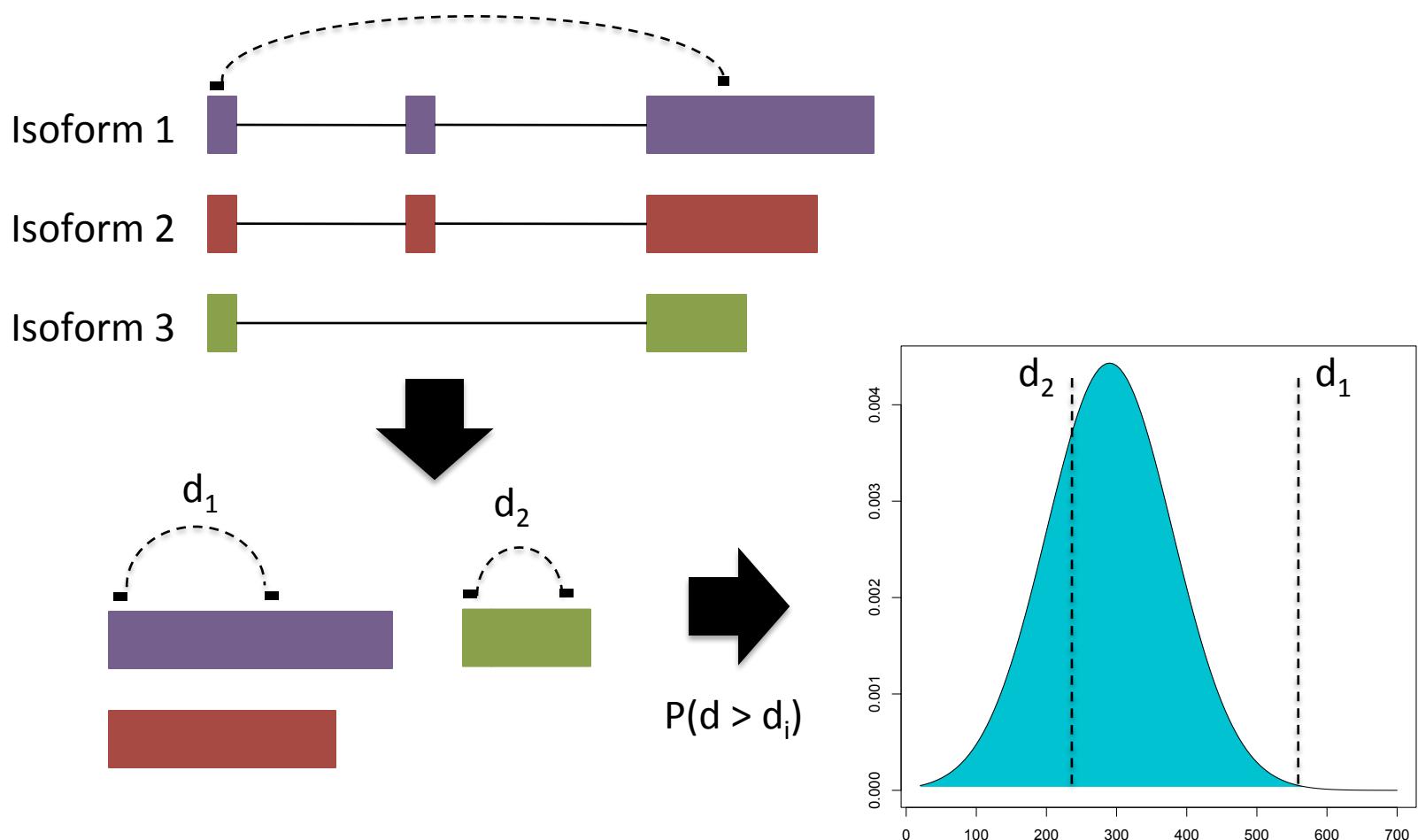
Splice and compute insert distance



Estimate insert size empirical distribution



... and use it for probabilistic read assignment

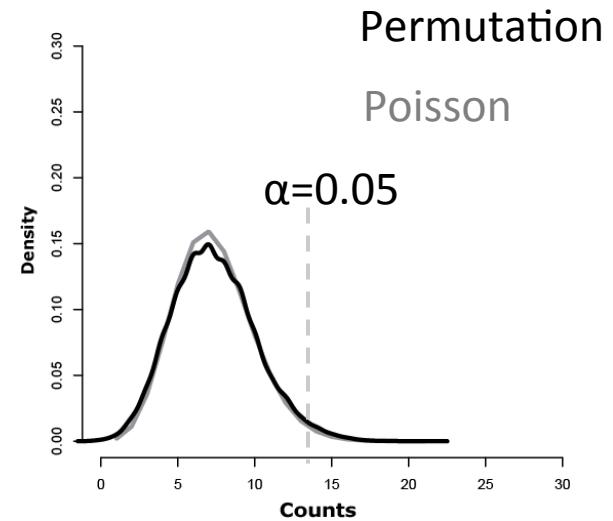
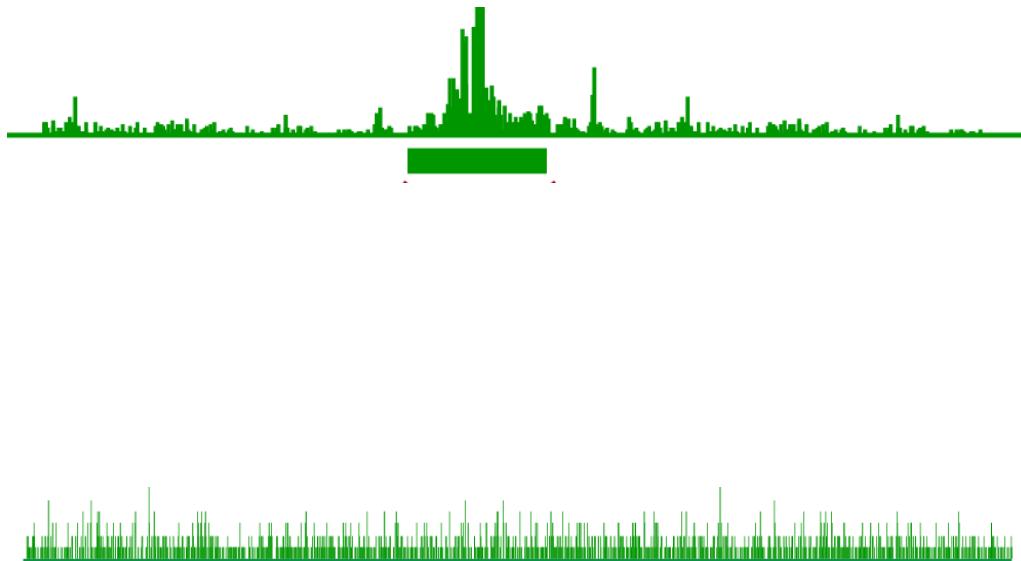


For methods such as MISO, Cufflinks and RSEM, it is critical to have paired-end data

RNA-Seq quantification summary

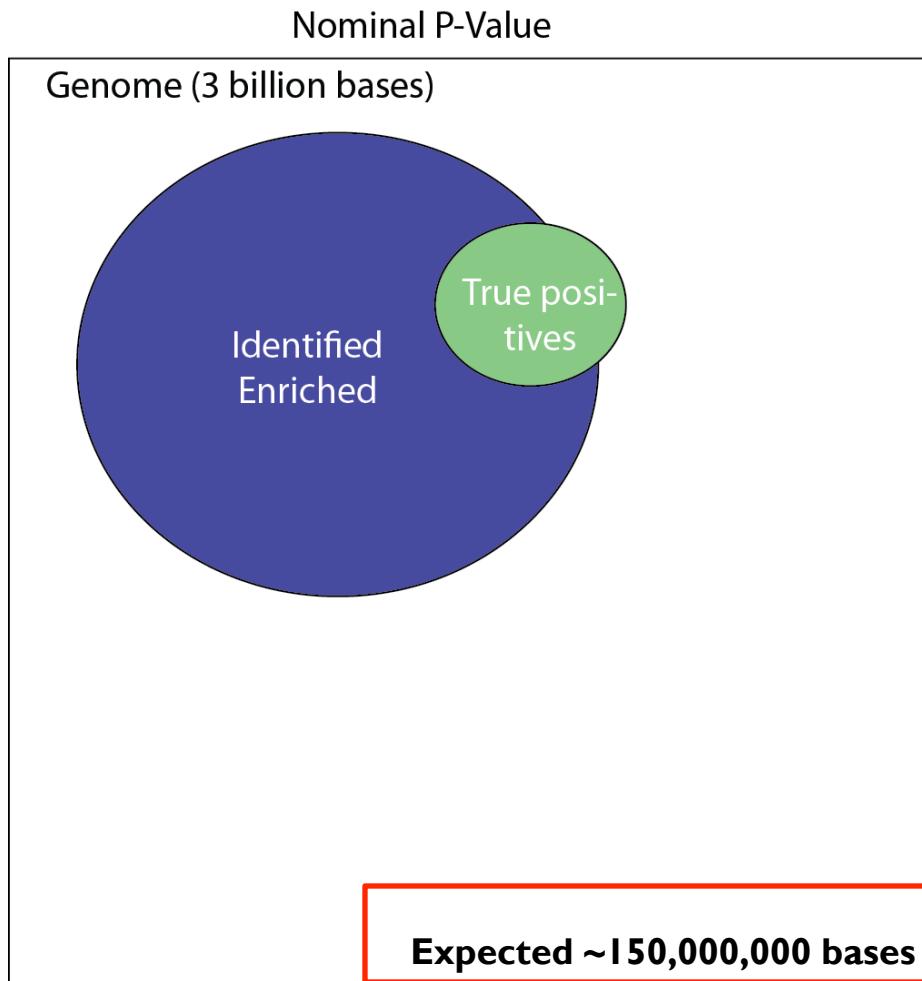
- Counts must be estimated from ambiguous read/transcript assignment.
 - Using simplified gene models (intersection)
 - Probabilistic read assignment
- Counts must be normalized
 - RPKM is sufficient for intra-library comparisons
 - More sophisticated normalizations to account for differences in library composition for inter-library comparisons.

Our approach



We have an efficient way to compute read count p-values ...

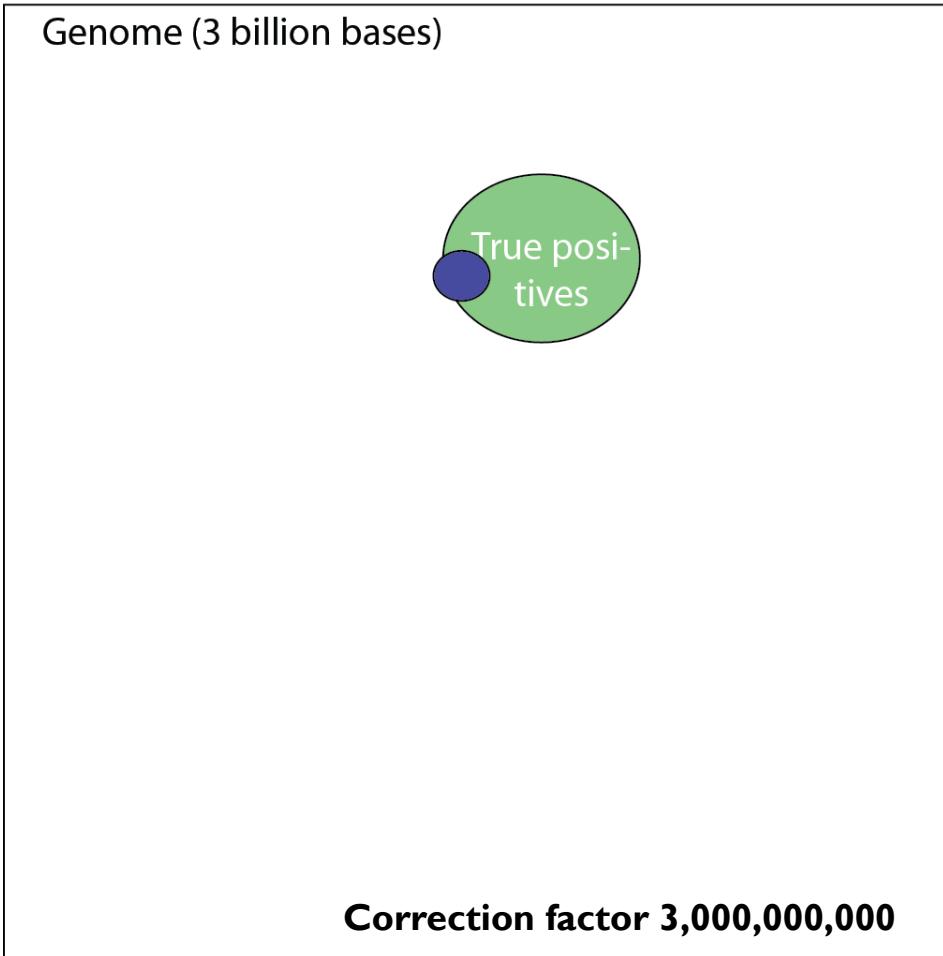
The genome is large, many things happen by chance



We need to correct for multiple hypothesis testing

Bonferroni correction is way to conservative

FWER-Bonferroni

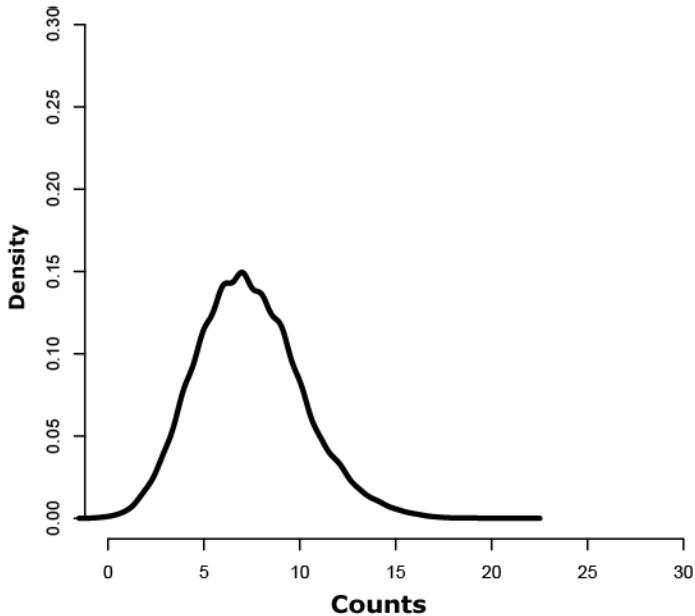


Bonferroni corrects the number of hits but misses many true hits because its too conservative – How do we get more power?

Controlling FWER

Max Count distribution

$$\alpha=0.05 \quad \alpha_{FWER}=0.05$$



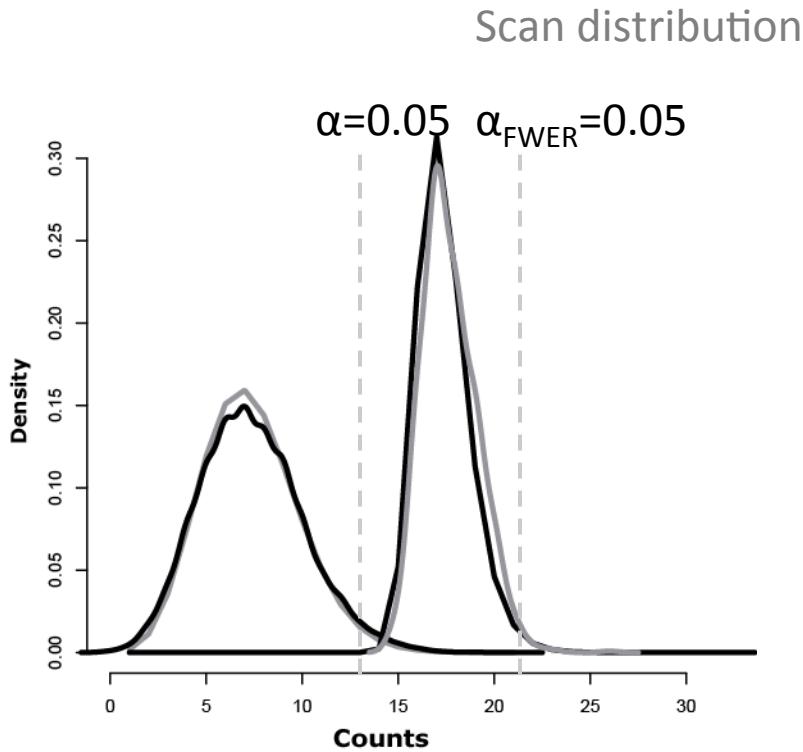
Count distribution (Poisson)

Given a region of size w and an observed read count n . What is the probability that one or more of the 3×10^9 regions of size w has read count $\geq n$ under the null distribution?

We could go back to our permutations and compute an FWER: **max of the genome-wide distributions of same sized region) →** but really really slow!!!

Scan distribution, an old problem

- Is the observed number of read counts over our region of interest high?
- Given a set of Geiger counts across a region find clusters of high radioactivity
- Are there time intervals where assembly line errors are high?



Poisson distribution

Thankfully, the **Scan Distribution** computes a closed form for this distribution.

ACCOUNTS for dependency of overlapping windows thus more powerful!

Scan distribution for a Poisson process

The probability of observing k reads on a window of size w in a genome of size L given a total of N reads can be approximated by (Alm 1983):

$$P(k|\lambda w, N, L) \approx 1 - F_p(k-1|\lambda w) e^{-\frac{k-w\lambda}{k}\lambda(T-w)} P(k-1|\lambda w)$$

where

$P(k-1|\lambda w)$ is the Poisson probability of observing $k-1$ counts given an expected count of λw

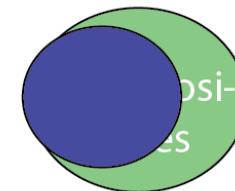
and

$F_p(k-1|\lambda w)$ is the Poisson probability of observing $k-1$ or fewer counts given an expectation of λw reads

The scan distribution gives a computationally very efficient way to estimate the FWER

FWER-Scan Statistics

Genome (3 billion bases)



By utilizing the dependency of overlapping windows we have greater power, while still controlling the same genome-wide false positive rate.