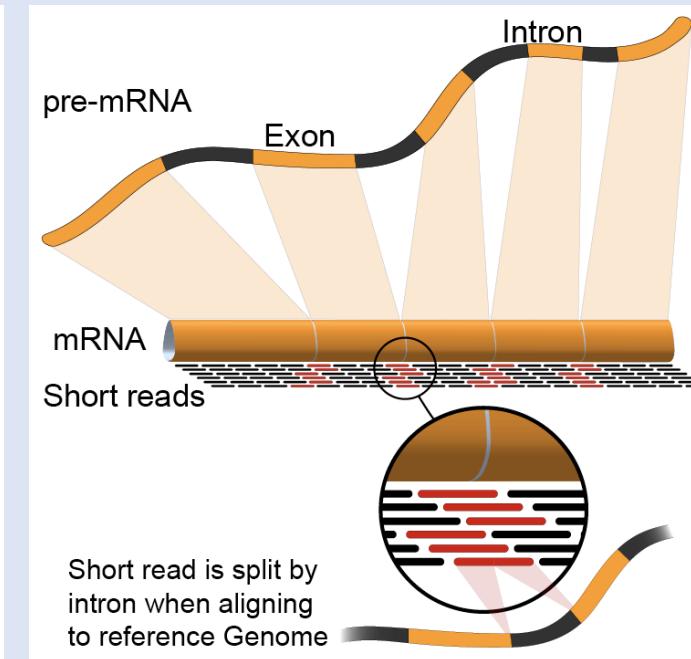
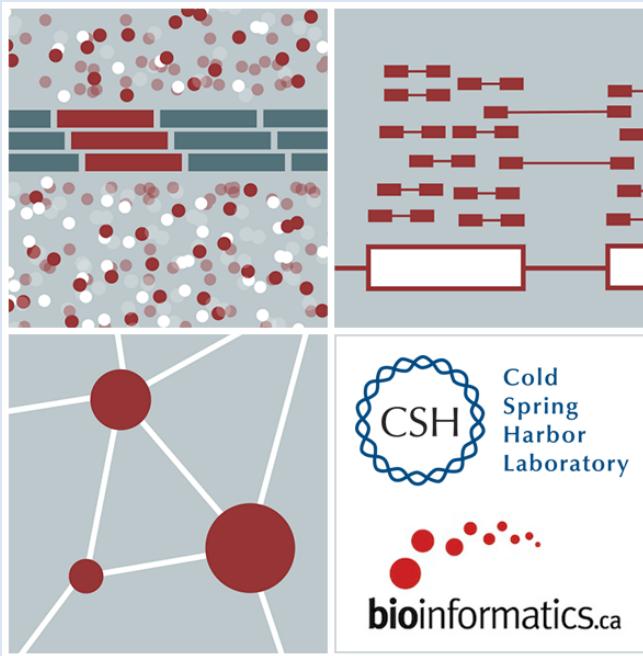




Cold  
Spring  
Harbor  
Laboratory



# RNA-Seq Week 2

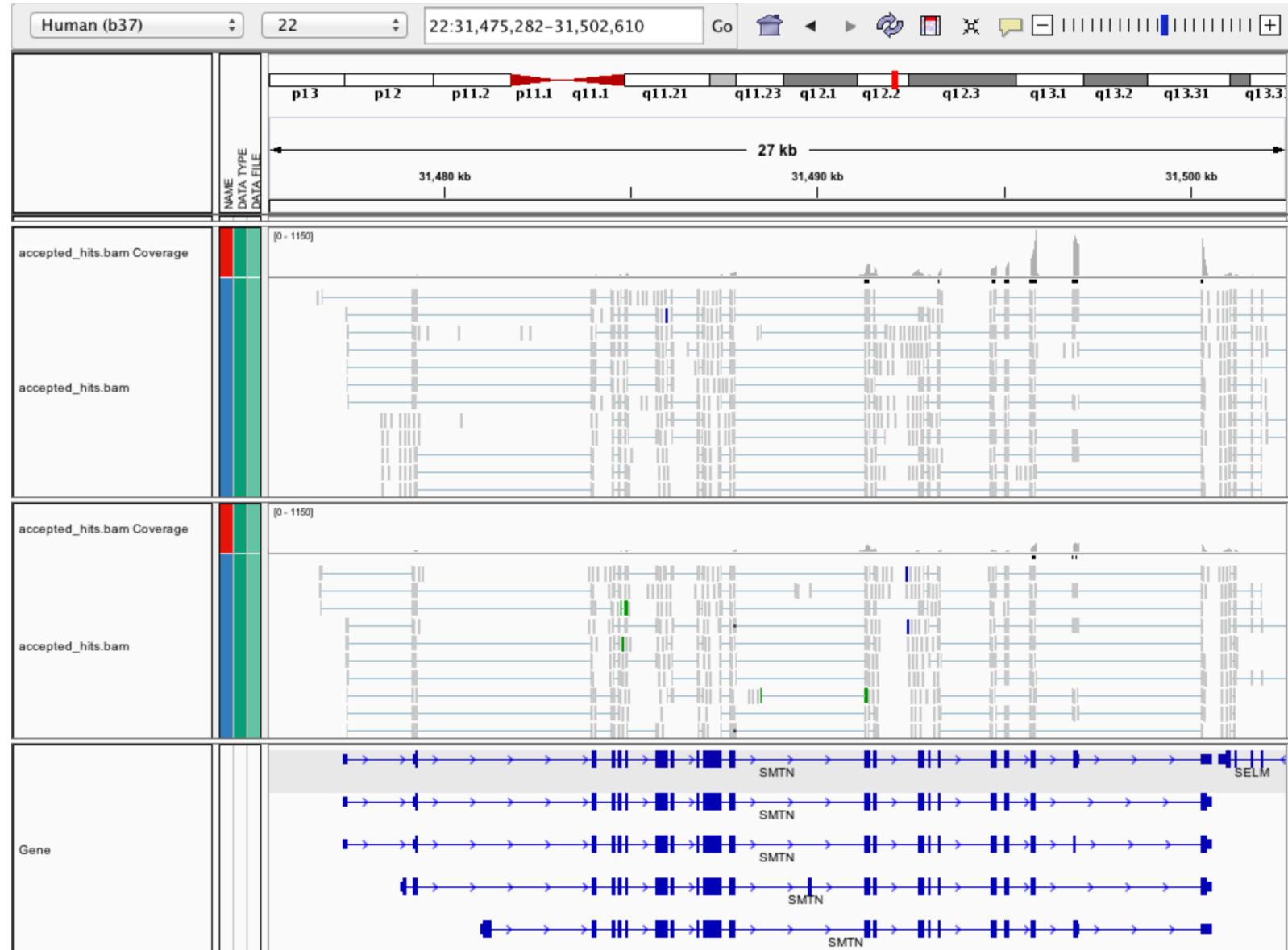
## Abundance Estimation, Differential Expression, and Alignment Free Expression Estimation (Kallisto)

Kelsy Cotto, Felicia Gomez, Obi Griffith, Malachi Griffith,  
Megan Richters, Huiming Xia

Advanced Sequencing Technologies & Bioinformatics Analysis November 16-20, 2020



# Expression estimation for known genes and transcripts



3' bias  
→

Down-regulated  
↓

# What is FPKM (RPKM)?

- RPKM: **Reads Per Kilobase of transcript per Million mapped reads.**
- FPKM: **Fragments Per Kilobase of transcript per Million mapped reads.**
- No essential difference - Just a terminology change to better describe paired-end reads!

# What is FPKM?

- Why not just count reads in my RNAseq data? → **Fragments**
- The relative expression of a transcript is proportional to the number of cDNA fragments that originate from it. However:
  - # fragments is biased towards larger genes → **Per Kilobase of transcript**
  - # fragments is related to total library depth → **per Million mapped reads.**

# What is FPKM?

- FPKM attempts to normalize for gene size and library depth
  - remember – RPKM is essentially the same!
- C = number of mappable fragments for a gene (transcript)
- N = total number of mappable fragments in the library
- L = number of base pairs in the gene (transcript)
  - $\text{FPKM} = (\text{C} / (\text{N} \times \text{L})) \times 1,000 \times 1,000,000$
  - $\text{FPKM} = (1,000,000,000 \times \text{C}) / (\text{N} \times \text{L})$
  - $\text{FPKM} = (\text{C} / (\text{N} / 1,000,000)) / (\text{L}/1000)$
- More reading:
  - <http://www.biostars.org/p/11378/>
  - <http://www.biostars.org/p/68126/>

# How do FPKM and TPM differ?

- TPM: Transcript per Kilobase Million
- The difference is in the order of operations:

## FPKM

- 1) Determine total fragment count, divide by 1,000,000 (per Million)
- 2) Divide each gene/transcript fragment count by #1 (Fragments Per Million)
- 3) Divide each FPM by length of each gene/transcript in kilobases (FPKM)

## TPM

- 1) Divide each gene/transcript fragment count by length of the transcript in kilobases (Fragments Per Kilobase)
- 2) Sum all FPK values for the sample and divide by 1,000,000 (per Million)
- 3) Divide #1 by #2 (TPM)

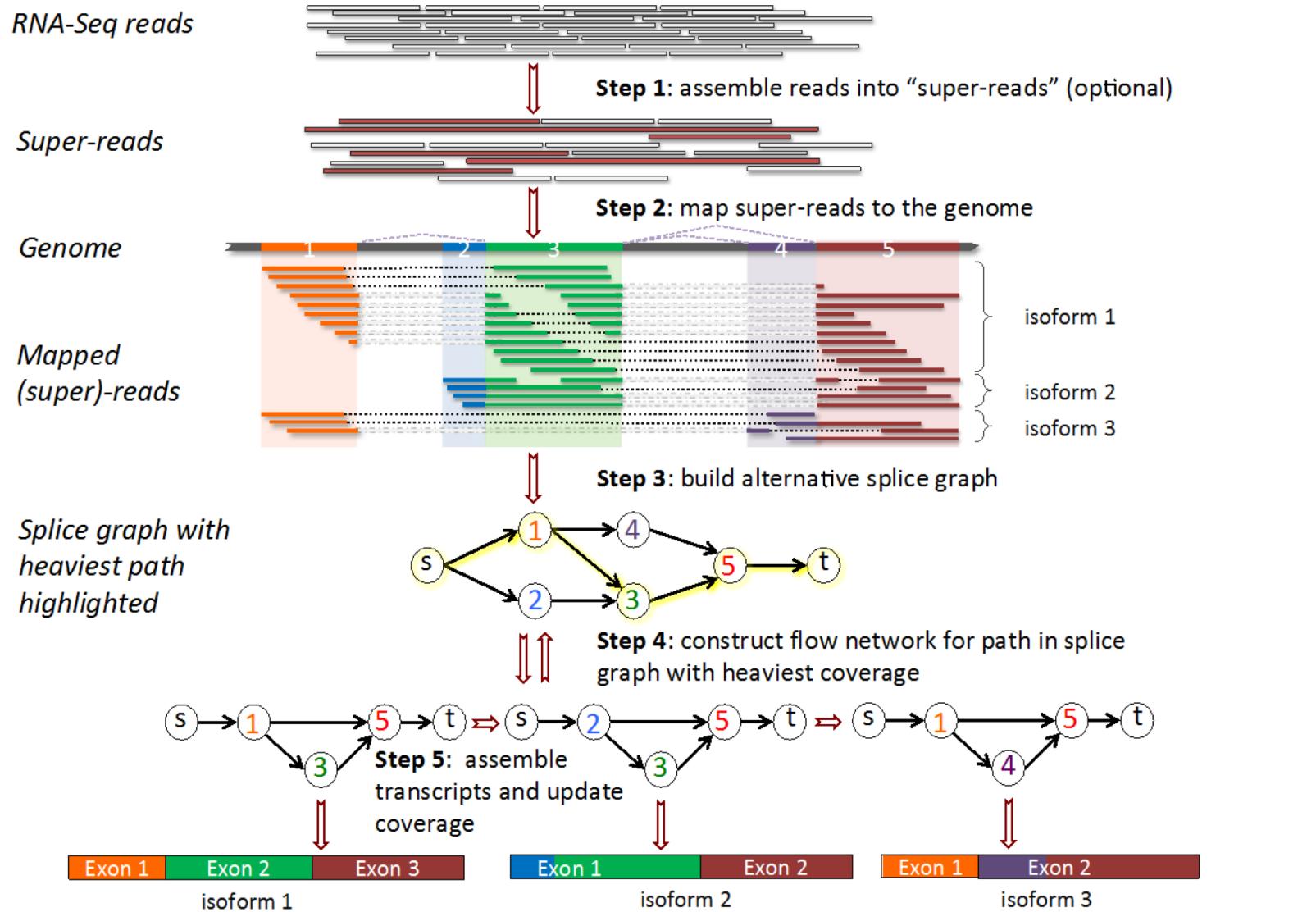
- The sum of all TPMs in each sample is the same. Easier to compare across samples!
- <http://www.rna-seqblog.com/rpkm-fpkm-and-tpm-clearly-explained/>
- <https://www.ncbi.nlm.nih.gov/pubmed/22872506>

# How does StringTie work?

Map reads to the genome

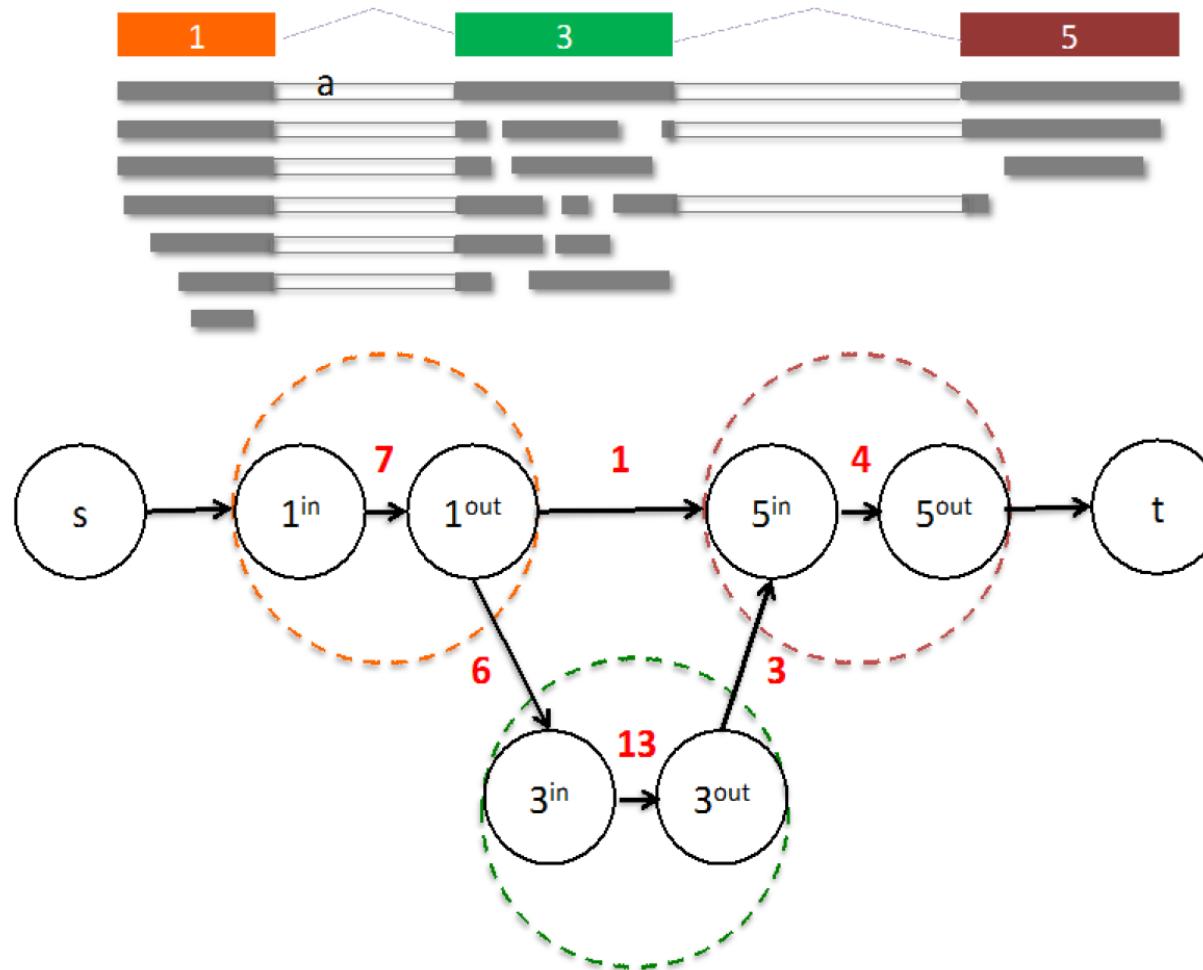
Infer isoforms:

- iteratively extract the heaviest path from a splice graph
- construct a flow network
- compute maximum flow to estimate abundance
- update the splice graph by removing reads that were assigned by the flow algorithm
- This process repeats until all reads have been assigned.



Perteau et al. Nature Biotechnology, 2015

# From flow network for each transcript, maximum flow is used to assemble transcript and estimate abundance



StringTie uses basic graph theory (splice graph), custom heuristics (heaviest path), more graph theory (flow network) and optimization theory (maximum flow). See StringTie paper for definitions and math.

# StringTie -merge

- Merge together all gene structures from all samples
  - Some samples may only partially represent a gene structure
- Incorporates known transcripts with assembled, potentially novel transcripts
- For de novo or reference guided mode, we will rerun StringTie with the merged transcript assembly.

Pertea et al. Nature Protocols, 2016

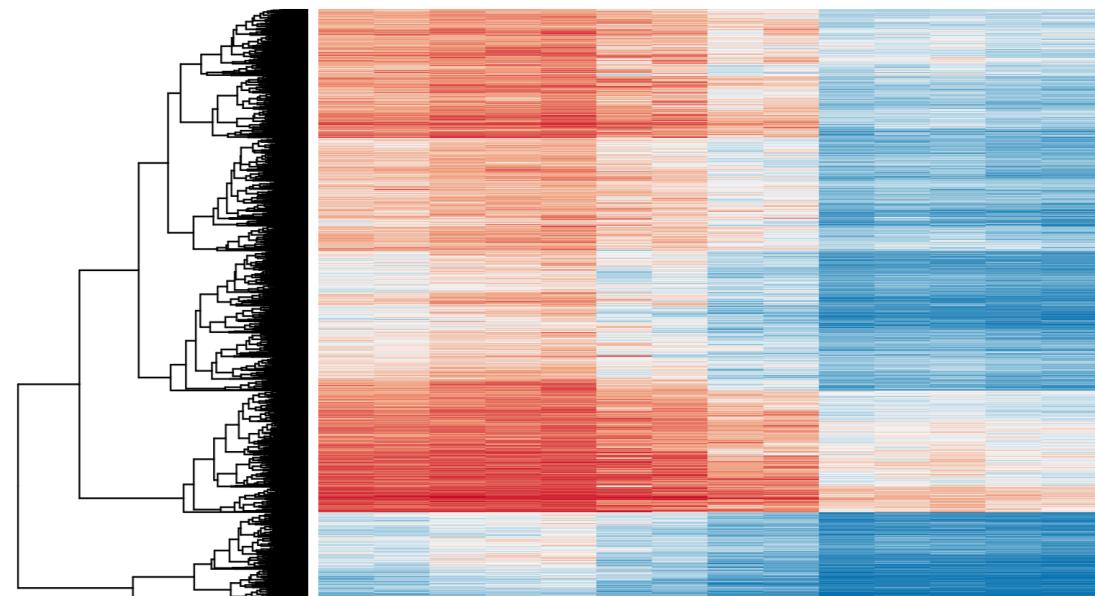
# gffcompare

- gffcompare will compare a merged transcript GTF with known annotation, also in GTF/GFF3 format
- <http://cole-trapnell-lab.github.io/cufflinks/cuffcompare/index.html#cuffcompare-output-files>

Priority	Code	Description
1	=	Complete match of intron chain
2	c	Contained
3	j	Potentially novel isoform (fragment): at least one splice junction is shared with a reference transcript
4	e	Single exon transfrag overlapping a reference exon and at least 10 bp of a reference intron, indicating a possible pre-mRNA fragment.
5	i	A transfrag falling entirely within a reference intron
6	o	Generic exonic overlap with a reference transcript
7	p	Possible polymerase run-on fragment (within 2Kbases of a reference transcript)
8	r	Repeat. Currently determined by looking at the soft-masked reference sequence and applied to transcripts where at least 50% of the bases are lower case
9	u	Unknown, intergenic transcript
10	x	Exonic overlap with reference on the opposite strand
11	s	An intron of the transfrag overlaps a reference intron on the opposite strand (likely due to read mapping errors)
12	.	(.tracking file only, indicates multiple classifications)

# Differential Expression

- Tying gene expression back to genotype/phenotype
- What genes/transcripts are being expressed at higher/lower levels in different groups of samples?
  - Are these differences 'significant', accounting for variance/noise?
- Examples (used in course):
  - UHR cells vs HBR brain
  - Tumor vs Normal tissue
  - Wild-type vs gene KO cells



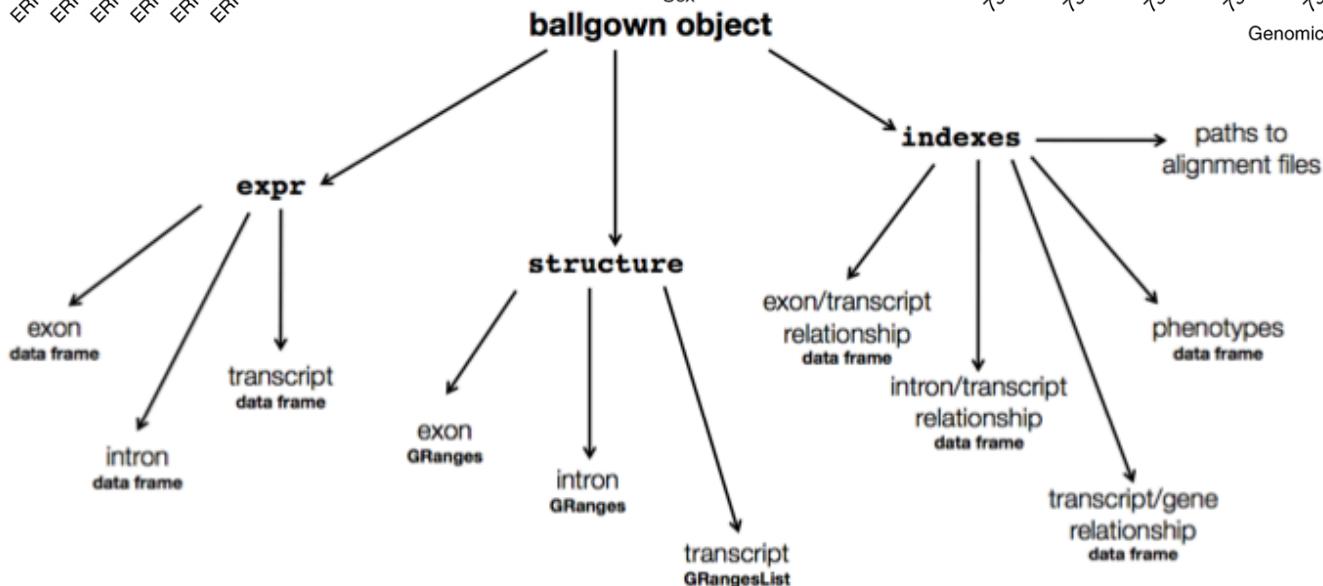
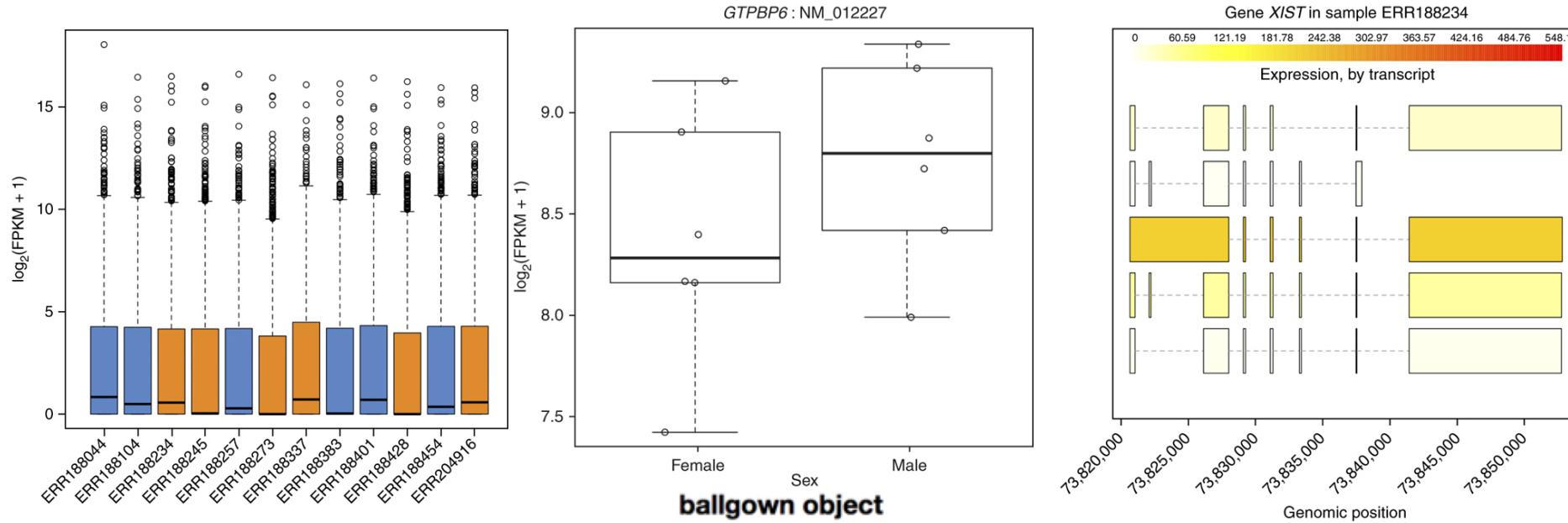
# Differential Expression with Ballgown

Parametric F-test comparing nested linear models

- Two models are fit to each feature, using expression as the outcome
  - one including the covariate of interest (e.g., case/control status or time) and one not including that covariate.
- An F statistic and p-value are calculated using the fits of the two models.
  - A significant p-value means the model including the covariate of interest fits significantly better than the model without that covariate, indicating differential expression.
- We adjust for multiple testing by reporting q-values:
  - $q < 0.05$  the false discovery rate should be controlled at  $\sim 5\%$ .

[Frazee et al. \(2014\)](#)

# Ballgown for Visualization with R



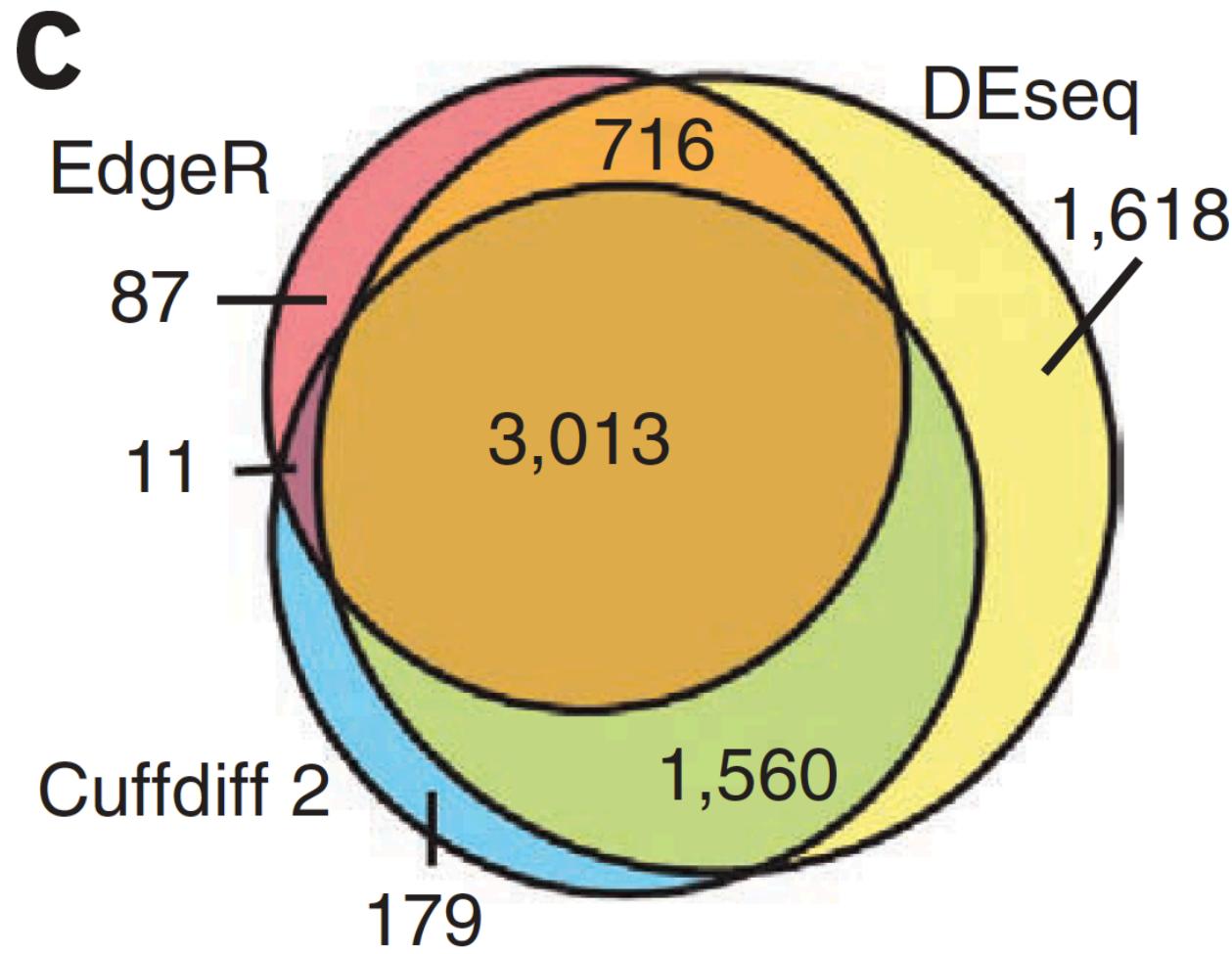
# Alternative differential expression methods

- Raw count approaches
  - DESeq2 - <http://www-huber.embl.de/users/anders/DESeq/>
  - edgeR - <http://www.bioconductor.org/packages/release/bioc/html/edgeR.html>
  - Others...

# **'FPKM/TPM' expression estimates vs. 'raw' counts**

- Which should I use?
  - Long running debate, but the general consensus:
- FPKM/TPM
  - When you want to leverage benefits of tuxedo suite
    - Isoform deconvolution
  - Good for visualization (e.g., heatmaps)
  - Calculating fold changes, etc.
- Counts
  - More robust statistical methods for differential expression
  - Accommodates more sophisticated experimental designs with appropriate statistical tests

# Multiple approaches advisable



# Lessons learned from microarray days

- Hansen et al. “Sequencing Technology Does Not Eliminate Biological Variability.” *Nature Biotechnology* 29, no. 7 (2011): 572–573.
- Power analysis for RNA-seq experiments
  - <http://scotty.genetics.utah.edu/>
- RNA-seq need for biological replicates
  - <http://www.biostars.org/p/1161/>
- RNA-seq study design
  - <http://www.biostars.org/p/68885/>

# Multiple testing correction

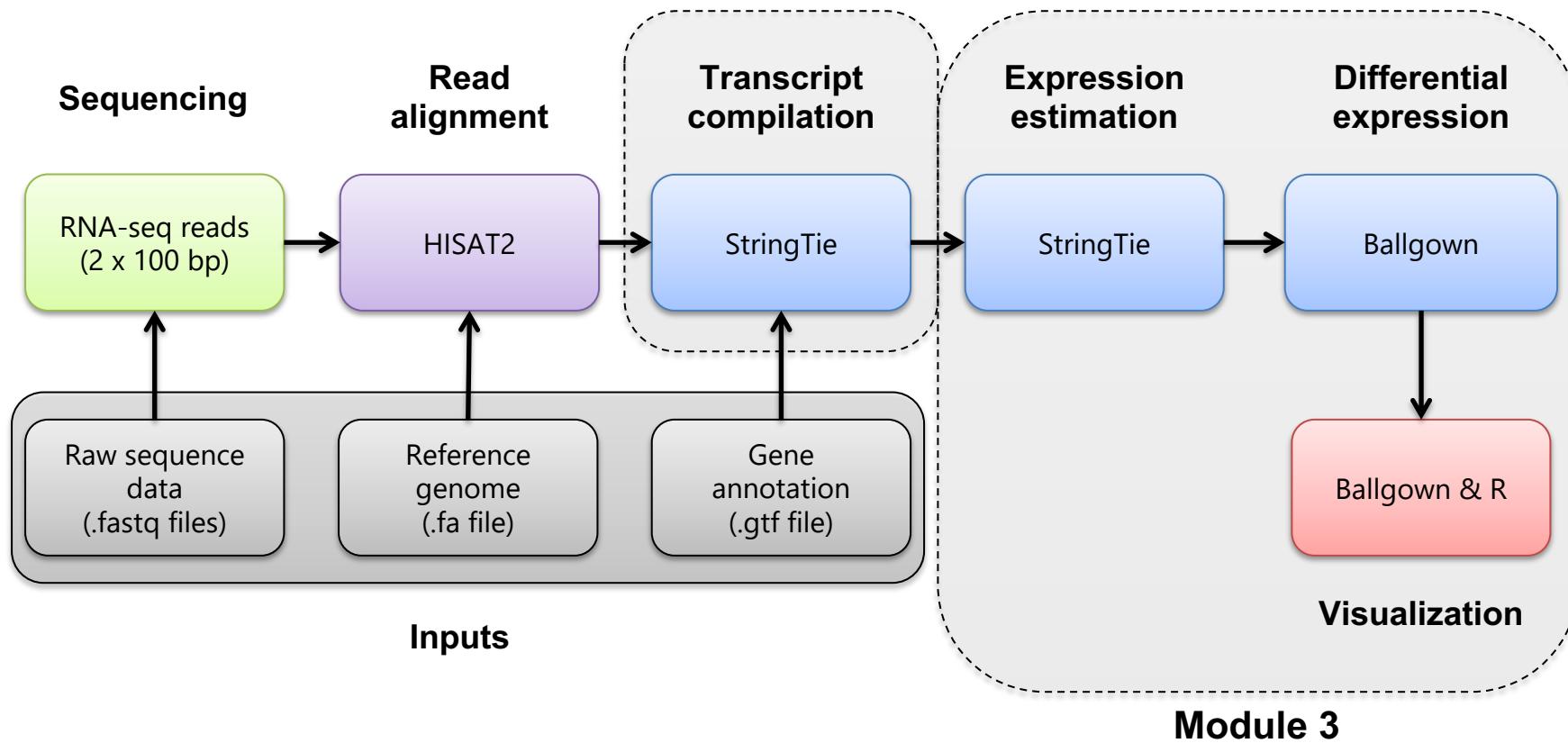
- As more attributes are compared, differences due solely to chance become more likely!
- Well known from array studies
  - 10,000s genes/transcripts
  - 100,000s exons
- With RNA-seq, more of a problem than ever
  - All the complexity of the transcriptome gives huge numbers of potential features
    - Genes, transcripts, exons, junctions, retained introns, microRNAs, lncRNAs, etc
- Bioconductor multtest
  - <http://www.bioconductor.org/packages/release/bioc/html/multtest.html>

# Downstream interpretation of expression analysis

- Topic for an entire course
- Expression estimates and differential expression lists from StringTie, Ballgown or other alternatives can be fed into many analysis pipelines
- See supplemental R tutorial for how to format expression data and start manipulating in R
- Clustering/Heatmaps
  - Provided by Ballgown
  - For more customized analysis various R packages exist:
    - hclust, heatmap.2, plotrix, ggplot2, etc.
- Classification
  - For RNA-seq data we still rarely have sufficient sample size and clinical details but this is changing
    - Weka is a good learning tool
    - RandomForests R package (biostar tutorial being developed)
- Pathway analysis
  - GSEA, IPA, Cytoscape, many R/BioConductor packages:  
<http://www.bioconductor.org/help/search/index.html?q=pathway>

[https://genviz.org/module-04-expression/0004/01/01/Expression\\_Profiling\\_and\\_Visualization/](https://genviz.org/module-04-expression/0004/01/01/Expression_Profiling_and_Visualization/)

# HISAT2/StringTie/Ballgown RNA-seq Pipeline



# What is a k-mer?

- A fixed sized ( $K$ ) sequence

1-mer

A
C
G
T

2-mer

AA	AC	AG	AT
CA	CC	CG	CT
GA	GC	GG	GT
TA	TC	TG	TT

- A string of length  $N$  contains  $N-K+1$  k-mers

ATTCGACAGTAGCCATGACTGG

...

- One can build  $K$ -mer index to represent a string

7-mer	iD	N
ATTCGAC	1	1
TTCGACA	2	1
TCGACAG	3	1
...		

Sailfish: Alignment-free Isoform Quantification from RNA-seq Reads using Lightweight Algorithms Rob Patro, Stephen M. Mount, and Carl Kingsford. *Manuscript Submitted* (2013) <http://www.cs.cmu.edu/~ckingsf/class/02714-f13/Lec05-sailfish.pdf>

<https://www.slideshare.net/duruofei/cmsc702-project-final-presentation>

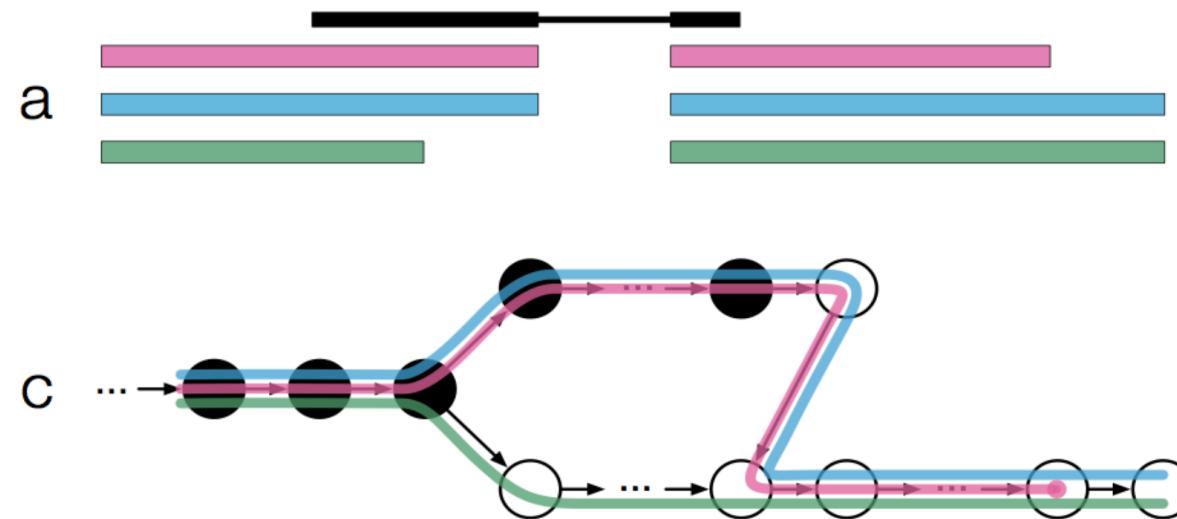
# Alignment free approaches for transcript abundance

1. Obtain reference transcript sequences
  - e.g. Ensembl, Refseq, or GENCODE
2. Build a **k-mer index** of all of the k-mers in each transcript sequence
  - Store each k-mer and its position within the transcript. “hashing”

# Alignment free approaches for transcript abundance

## 3. Count number of times each k-mer occurs within each RNAseq read

- Model relationship between RNA-seq read k-mers and the transcript k-mer index.
- What transcript is the most likely source for each read?
- Called “pseudoalignment”, “quasi-mapping”, etc.



Bray, 2016 doi:10.1038/nbt.3519

<https://tinyheero.github.io/2015/09/02/pseudoadalignments-kallisto.html>

## 4. Handle sequencing errors, isoforms, ambiguity, and determine abundance estimates

- Transcriptome de Bruijn graphs, likelihood function, expectation maximization, etc.

# Advantages/disadvantages of alignment free approaches

- Advantages
  - Very fast and efficient
    - Similar accuracy to alignment based approach but with much, much shorter run time.
  - Do not need a reference genome, only a reference transcriptome
- Disadvantages
  - You don't get a proper BAM file (though a pseudo-bam can be created)
  - Information in reads with sequence errors may be ignored
  - Limited potential for transcript discovery, variant calling, fusion detection, etc.

# Common alignment free tools

- Sailfish
  - “Sailfish enables alignment-free isoform quantification from RNA-seq reads using lightweight algorithms.” 2014
  - <https://www.ncbi.nlm.nih.gov/pubmed/24752080>
- RNA-Skim
  - “RNA-Skim: a rapid method for RNA-Seq quantification at transcript level.” 2014
  - <https://www.ncbi.nlm.nih.gov/pubmed/24931995>
- Kallisto
  - “Near-optimal probabilistic RNA-seq quantification.” 2016
  - <https://www.ncbi.nlm.nih.gov/pubmed/27043002>
- Salmon
  - “Salmon provides fast and bias-aware quantification of transcript expression.” 2017
  - <https://www.ncbi.nlm.nih.gov/pubmed/28263959>

# Which is best?

- Somewhat controversial ...
- <https://liorpachter.wordpress.com/2017/08/02/how-not-to-perform-a-differential-expression-analysis-or-science/>
- Various sources suggest that Salmon, Kallisto, and Sailfish results are quite comparable
- Usability, documentation, and supporting downstream tools could be used to decide

We are on a Coffee Break &  
Networking Session