



Canadian Bioinformatics Workshops

www.bioinformatics.ca

This page is available in the following languages:

Afrikaans Български Català Dansk Deutsch Ελληνικά English English (CA) English (GB) English (US) Esperanto
Castellano Castellano (AR) Español (CL) Castellano (CO) Español (Ecuador) Castellano (MX) Castellano (PE)
Euskara Suomeksi français français (CA) Galego ລາວ hrvatski Magyar Italiano 日本語 한국어 Macedonian Melayu
Nederlands Norsk Sesotho sa Leboa polski Português română slovenski jezik srpski (latinica) Sotho svenska
中文 華語 (台灣) isiZulu



Attribution-Share Alike 2.5 Canada

You are free:



to Share — to copy, distribute and transmit the work

to Remix — to adapt the work



Under the following conditions:



Attribution. You must attribute the work in the manner specified by the author or licensor (but not in any way that suggests that they endorse you or your use of the work).



Share Alike. If you alter, transform, or build upon this work, you may distribute the resulting work only under the same or similar licence to this one.

- For any reuse or distribution, you must make clear to others the licence terms of this work.
- Any of the above conditions can be waived if you get permission from the copyright holder.
- The author's moral rights are retained in this licence.

[Disclaimer](#)

Your fair dealing and other rights are in no way affected by the above.

This is a human-readable summary of the Legal Code (the full licence) available in the following languages:

[English](#) [French](#)

Module 6

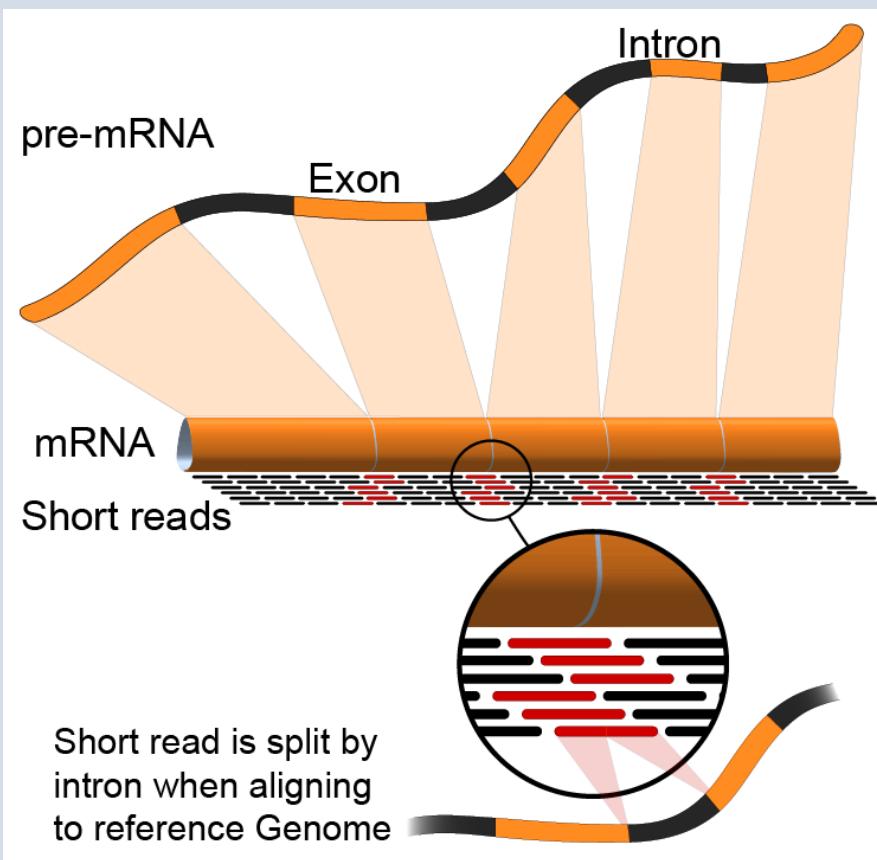
Gene Expression (RNA-Seq)

Bioinformatics for Cancer Genomics

May 31, 2017

Fouad Yousif

Adapted from: Obi Griffith and Malachi Griffith



Learning Objectives of The Module

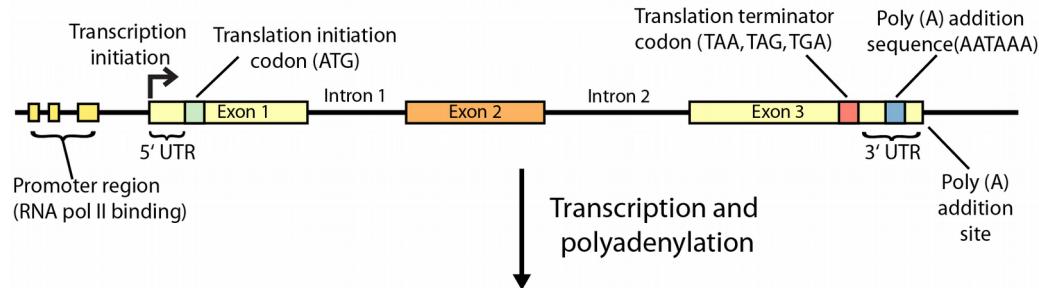
- **Part I: Introduction to RNA sequencing**
- Part II: RNA-seq alignment and visualization
- Part III: Expression and Differential Expression

Learning Objectives of Part I

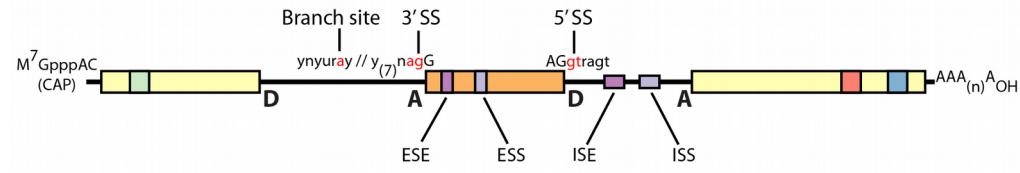
- Introduction to the theory and practice of RNA sequencing (RNA-seq) analysis
 - Rationale for sequencing RNA
 - Challenges specific to RNA-seq
 - General goals and themes of RNA-seq analysis work flows
 - Common technical questions related to RNA-seq analysis

Gene expression

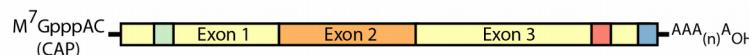
Double-stranded genomic DNA template



Single-stranded pre-mRNA (nuclear RNA)



Mature mRNA

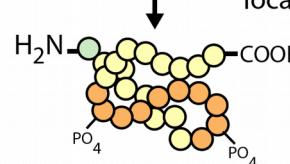


RNA processing

Protein (amino acid sequence)

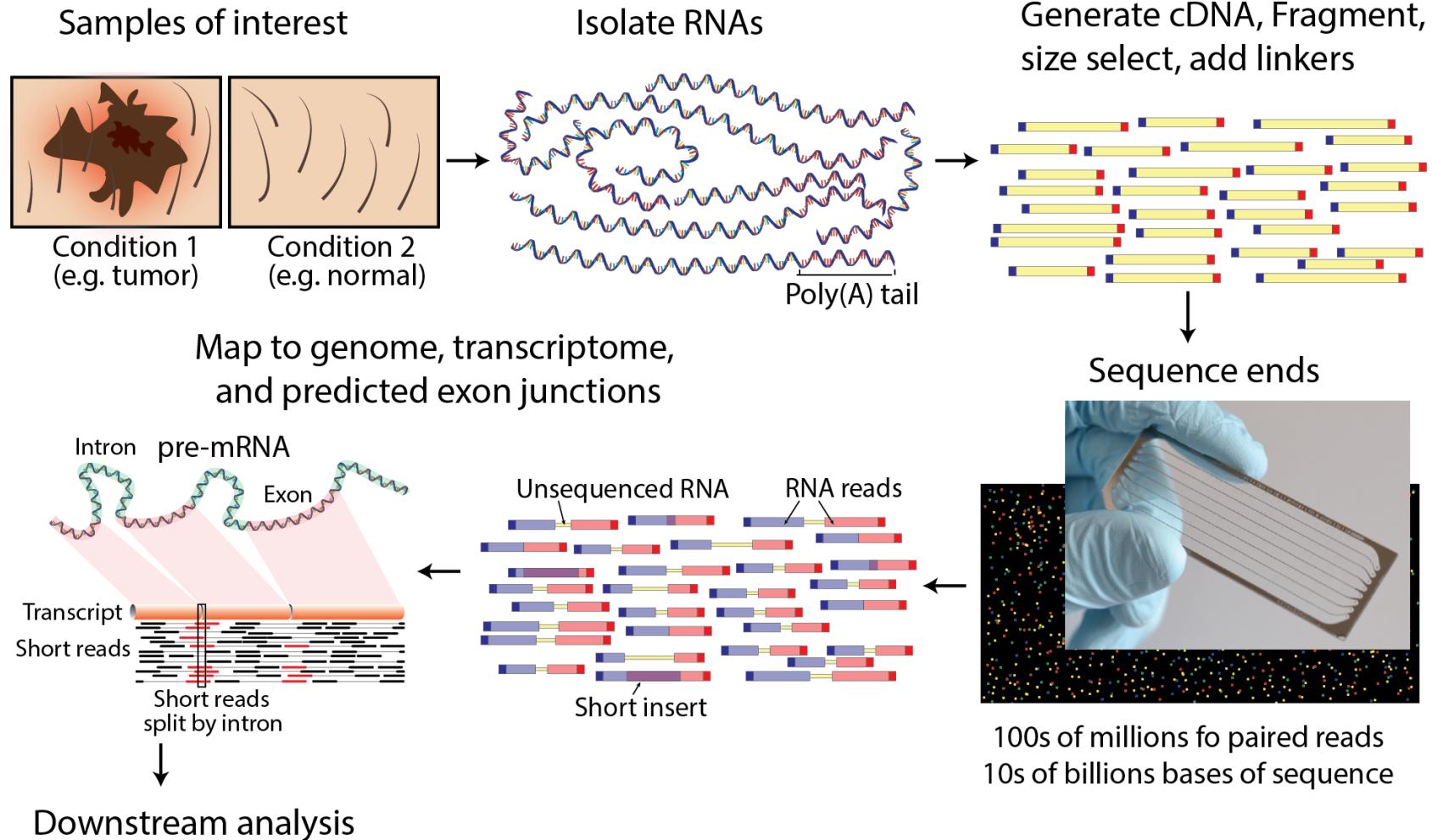


Export to cytoplasm and translation



Folding, posttranslational modification, subcellular localization, etc.

RNA sequencing



Why sequence RNA (versus DNA)?

- Functional studies
 - Genome may be constant but an experimental condition has a pronounced effect on gene expression
 - e.g. Drug treated vs. untreated cell line
 - e.g. Wild type versus knock out mice
- Predicting transcript sequence from genome sequence is difficult
 - Gene annotation is revolutionized by RNA-seq
- Some molecular features can only be observed at the RNA level
 - Alternative isoforms, fusion transcripts, RNA editing

Why sequence RNA (versus DNA)?

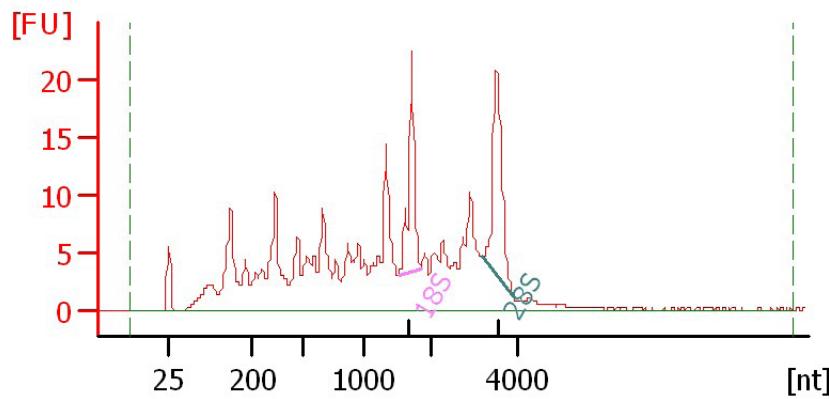
- Interpreting mutations that do not have an obvious effect on protein sequence
 - ‘Regulatory’ mutations that affect what mRNA isoform is expressed and how much
- Prioritizing protein coding somatic mutations (often heterozygous)
 - If the gene is not expressed, a mutation in that gene would be less interesting
 - If the gene is expressed but only from the wild type allele, this might suggest loss-of-function (haploinsufficiency)
 - If the mutant allele itself is expressed, this might suggest a candidate drug target

Challenges

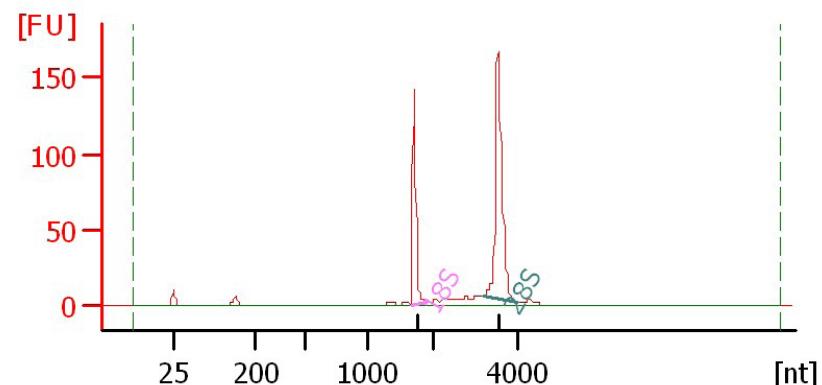
- Sample
 - Purity?, quantity?, quality?
- RNAs consist of small exons that may be separated by large introns
 - Mapping reads to genome is challenging
- The relative abundance of RNAs vary wildly
 - 10^5 – 10^7 orders of magnitude
 - Since RNA sequencing works by random sampling, a small fraction of highly expressed genes may consume the majority of reads
 - Ribosomal and mitochondrial genes
- RNAs come in a wide range of sizes
 - Small RNAs must be captured separately
 - PolyA selection of large RNAs may result in 3' end bias
- RNA is fragile compared to DNA (easily degraded)

Agilent example / interpretation

- https://github.com/griffithlab/rnaseq_tutorial/wiki/Resources/Agilent_Trace_Examples.pdf
- ‘RIN’ = RNA integrity number
 - 0 (bad) to 10 (good)

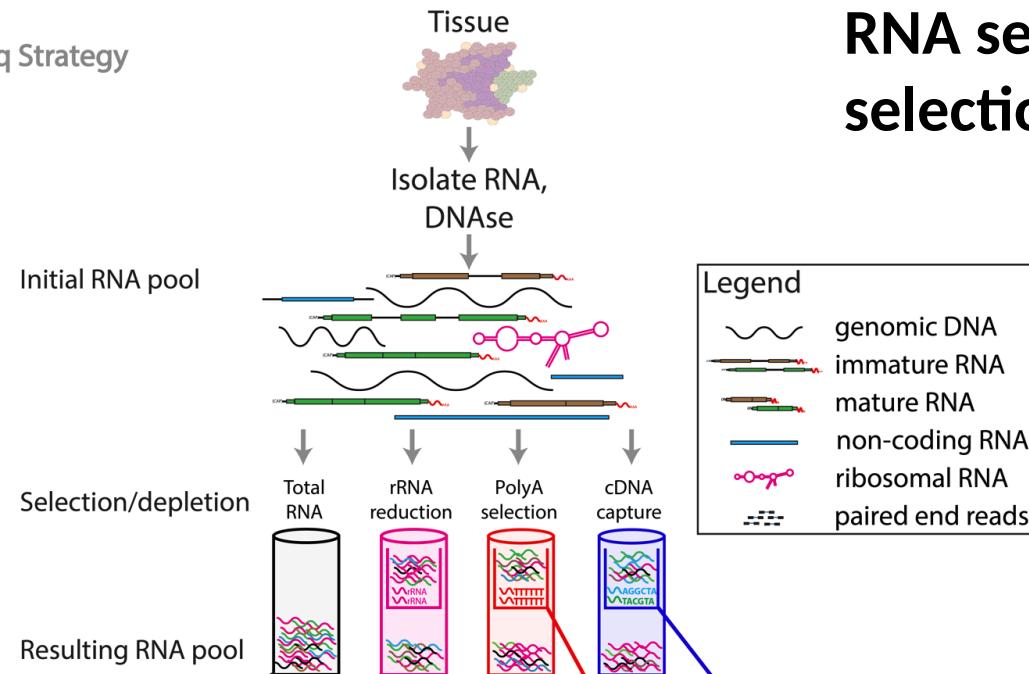


RIN = 6.0



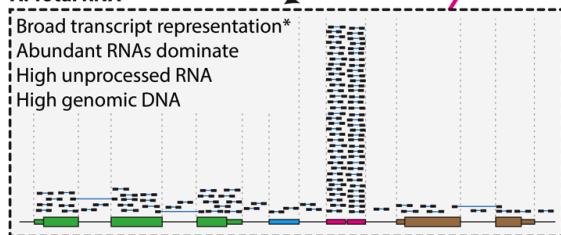
RIN = 10

RNA-seq Strategy

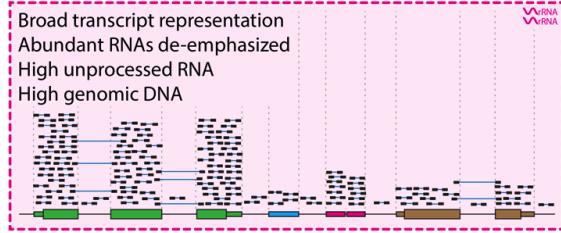


RNA sequence selection/depletion

A. Total RNA

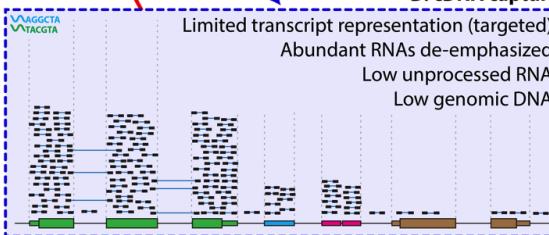


B. rRNA reduction

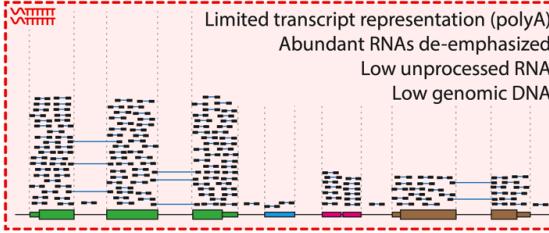


Expected Alignments

D. cDNA capture



C. PolyA selection

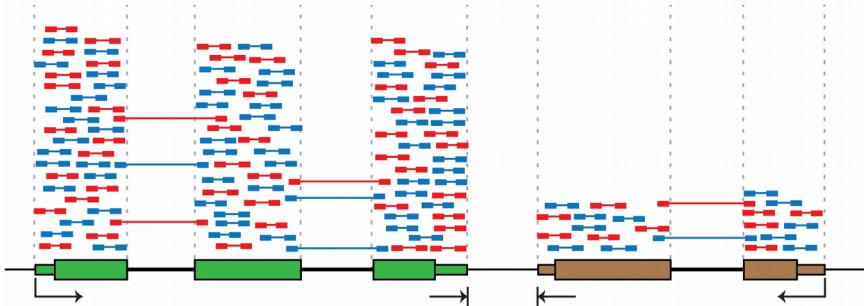


Stranded vs. unstranded

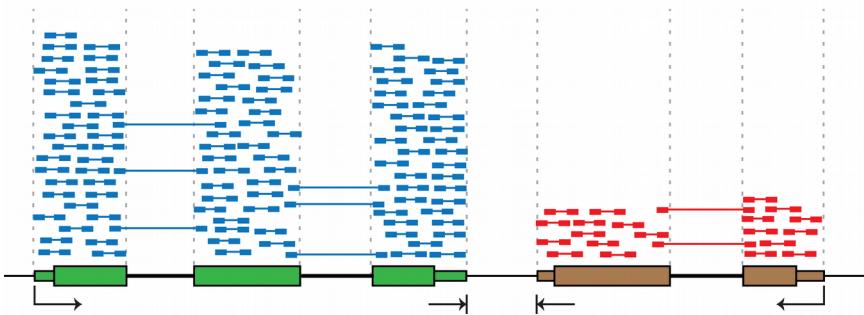
A. Depiction of cDNA fragments from an unstranded library

Legend

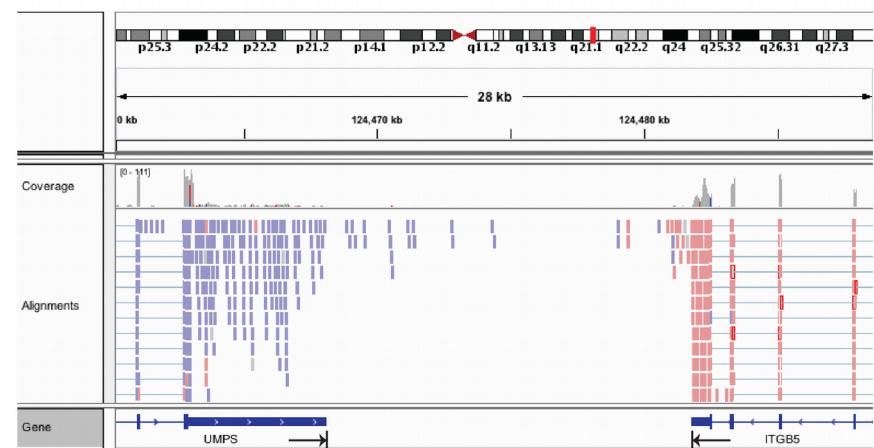
- Transcription start site and direction
- ← PolyA site (transcription end)
- Read sequenced from positive strand (forward)
- Read sequenced from negative strand (reverse)



B. Depiction of cDNA fragments from an stranded library

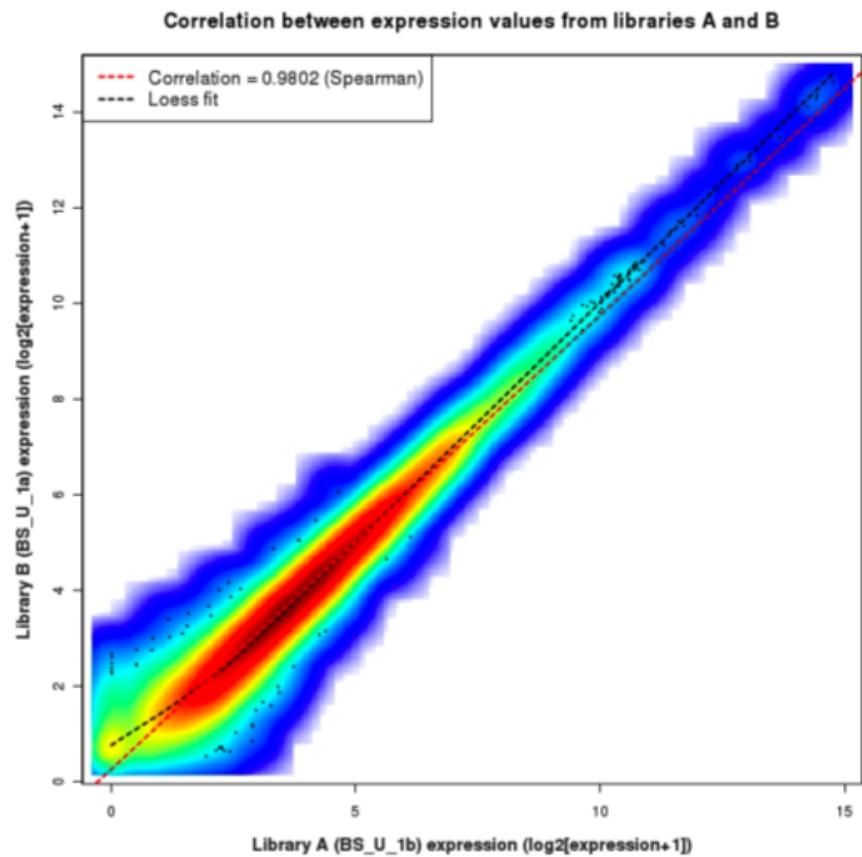


C. Viewing strand of aligned reads in IGV



Replicates

- Technical Replicate
 - Multiple instances of sequence generation
 - Flow Cells, Lanes, Indexes
- Biological Replicate
 - Multiple isolations of cells showing the same phenotype, stage or other experimental condition
 - Some example concerns/challenges:
 - Environmental Factors, Growth Conditions, Time
 - Correlation Coefficient 0.92-0.98



Common analysis goals of RNA-Seq analysis (what can you ask of the data?)

- Gene expression and differential expression
- Alternative expression analysis
- Transcript discovery and annotation
- Allele specific expression
 - Relating to SNPs or mutations
- Mutation discovery
- Fusion detection
- RNA editing

General themes of RNA-seq workflows

- Each type of RNA-seq analysis has distinct requirements and challenges but also a common theme:
 1. Obtain raw data (convert format)
 2. Align/assemble reads
 3. Process alignment with a tool specific to the goal
 - e.g. 'cufflinks' for expression analysis, 'defuse' for fusion detection, etc.
 4. Post process
 - Import into downstream software (R, Matlab, Cytoscape, Ingenuity, etc.)
 5. Summarize and visualize
 - Create gene lists, prioritize candidates for validation, etc.

Common questions: Should I remove duplicates for RNA-seq?

- Maybe... more complicated question than for DNA
- Concern.
 - Duplicates may correspond to biased PCR amplification of particular fragments
 - For highly expressed, short genes, duplicates are expected even if there is no amplification bias
 - Removing them may reduce the dynamic range of expression estimates
- If you do remove them, assess duplicates at the level of paired-end reads (fragments) not single end reads

Common questions: How much library depth is needed for RNA-seq?

- Depends on a number of factors:
 - Question being asked of the data. Gene expression? Alternative expression? Mutation calling?
 - Tissue type, RNA preparation, quality of input RNA, library construction method, etc.
 - Sequencing type: read length, paired vs. unpaired, etc.
 - Computational approach and resources
- Identify publications with similar goals
- Pilot experiment
- Good news: 1-2 lanes of recent Illumina HiSeq data should be enough for most purposes

Common questions: What mapping strategy should I use for RNA-seq?

- Depends on read length
- < 50 bp reads
 - Use aligner like BWA and a genome + junction database
 - Junction database needs to be tailored to read length
 - Or you can use a standard junction database for all read lengths and an aligner that allows substring alignments for the junctions only (e.g. BLAST ... slow).
 - Assembly strategy may also work (e.g. Trans-ABySS)
- > 50 bp reads
 - Spliced aligner such as Bowtie/TopHat/HISAT

Common questions: What if I don't have a reference genome for my species?

- Have you considered sequencing the genome of your species?
- If that is not practical or you simply prefer a transcript discovery approach that does not rely on prior knowledge of the genome or transcriptome there are some tools available ...
 - Unfortunately de novo transcriptome assembly is beyond the scope of this workshop
 - The good news is that the skills you learn here will help you figure out how to install and run those tools yourself

Learning Objectives of The Module

- Part I: Introduction to RNA sequencing
- **Part II: RNA-seq alignment and visualization**
- Part III: Expression and Differential Expression

Learning Objectives of Part II

- RNA-seq alignment challenges and common questions
- Alignment strategies
- Bowtie/TopHat/HISAT2
- Introduction to the BAM and BED formats
- Basic manipulation of BAMs
- Visualization of RNA-seq alignments in IGV
- Alignment QC Assessment
- BAM read counting and determination of variant allele expression status

RNA-seq alignment challenges

- Computational cost
 - 100's of millions of reads
- Introns!
 - Spliced vs. unspliced alignments
- Can I just align my data once using one approach and be done with it?
 - Unfortunately probably not
- Is HISAT2 the only mapper to consider for RNA-seq data?
 - <http://www.biostars.org/p/60478/>

Three RNA-seq mapping strategies

De novo assembly

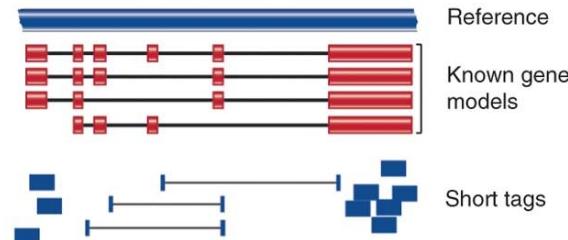


Assemble transcripts from overlapping tags



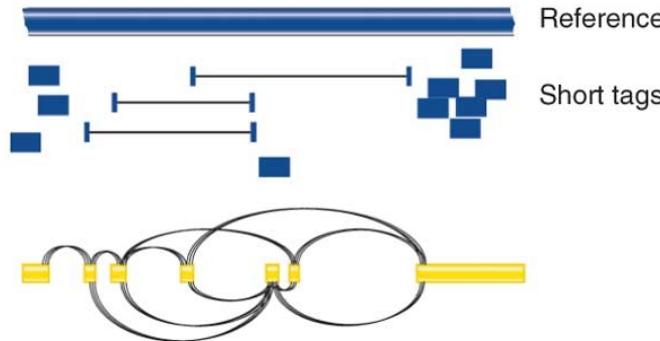
Optional: align to genome to get exon structure

Align to transcriptome



Use known and/or predicted gene models to examine individual features

Align to reference genome



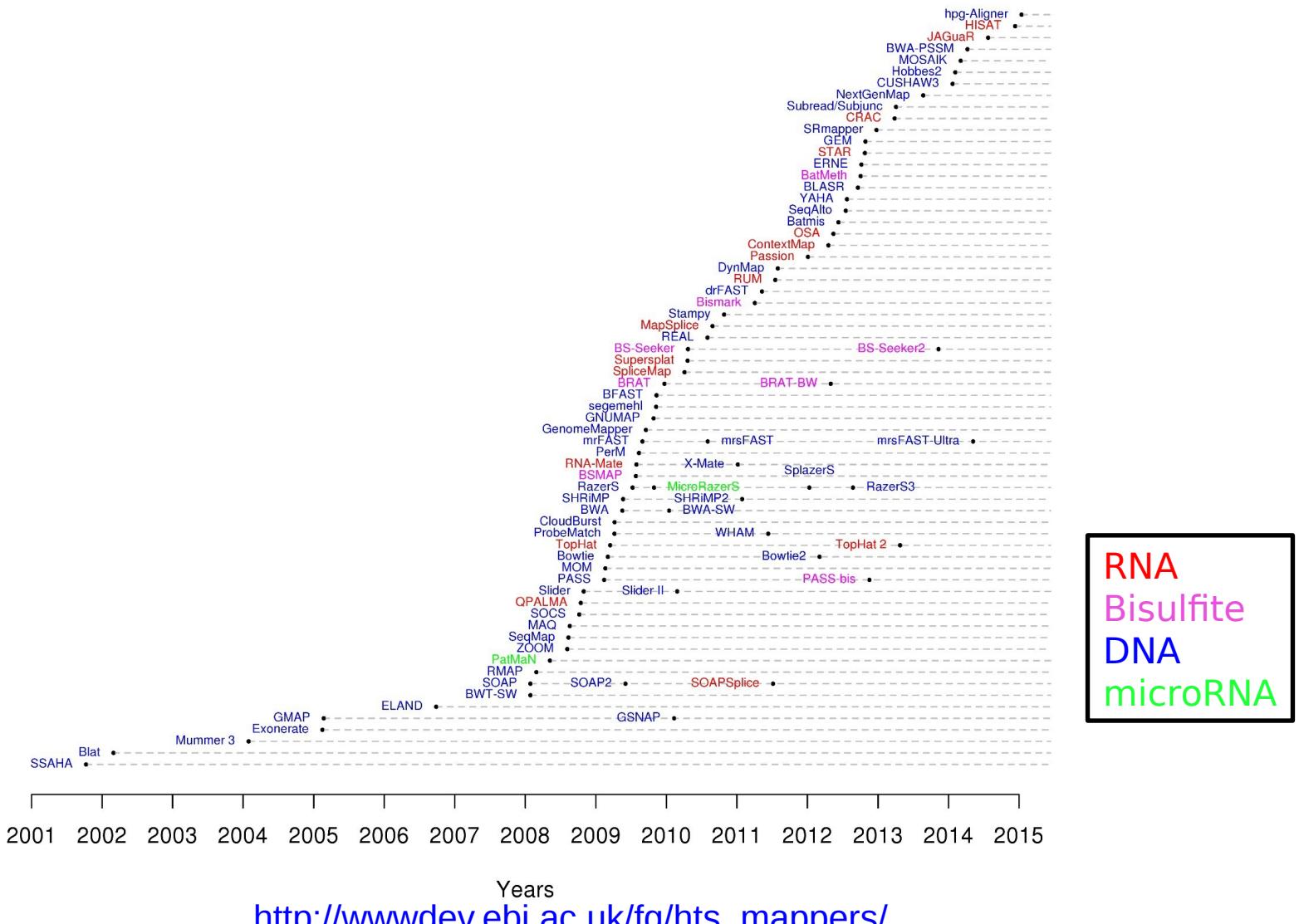
Infer possible transcripts and abundance

Diagrams from Cloonan & Grimmond, Nature Methods 2010

Which alignment strategy is best?

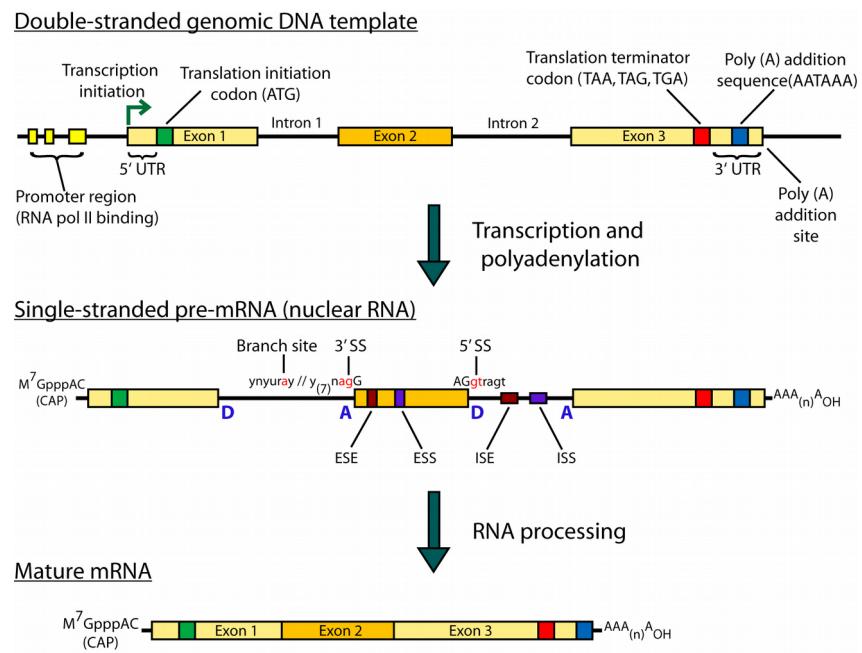
- De novo assembly
 - If a reference genome does not exist for the species being studied
 - If complex polymorphisms/mutations/haplotypes might be missed by comparing to the reference genome
- Align to transcriptome
 - If you have short reads (< 50bp)
- Align to reference genome
 - All other cases
- Each strategy involves different alignment/assembly tools

Which read aligner should I use?



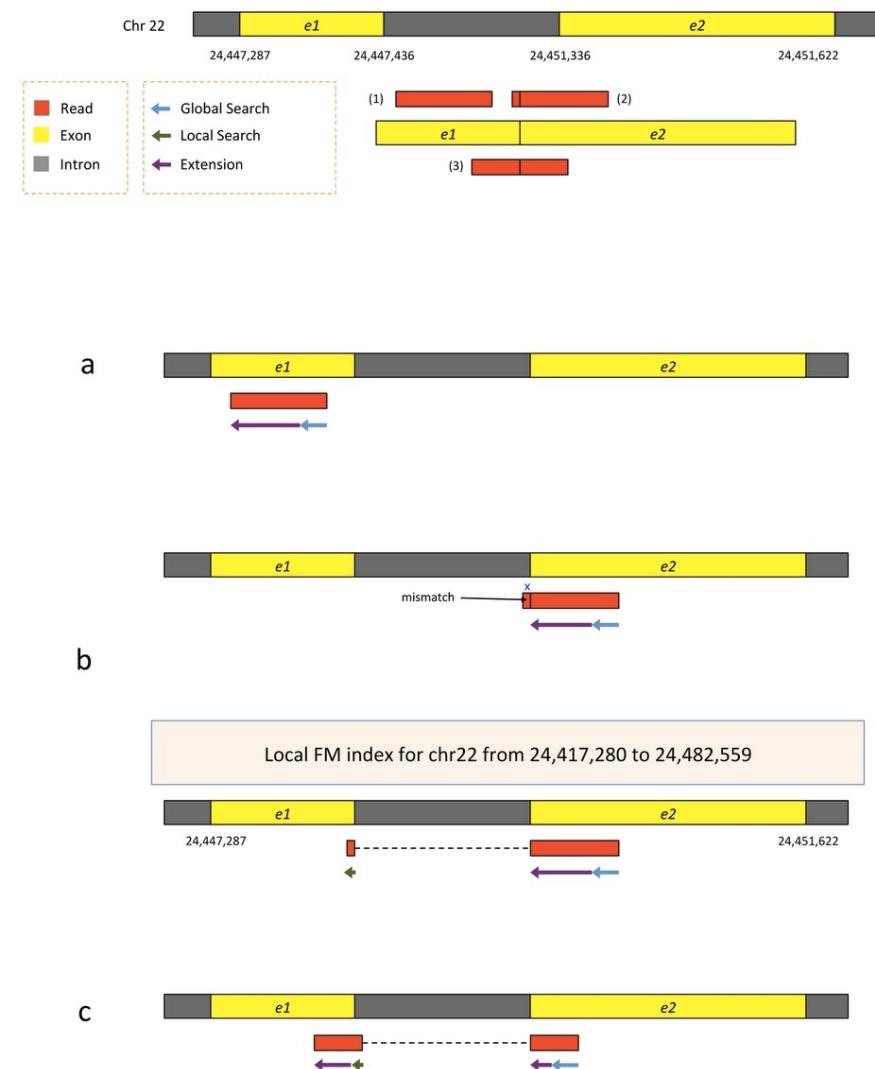
Should I use a splice-aware or unspliced mapper

- RNA-seq reads may span large introns
- The fragments being sequenced in RNA-seq represent mRNA and therefore the introns are removed
- But we are usually aligning these reads back to the reference genome
- Unless your reads are short (<50bp) you should use a splice-aware aligner
 - HISAT2, TopHat, STAR, MapSplice, etc.



HISAT/HISAT2

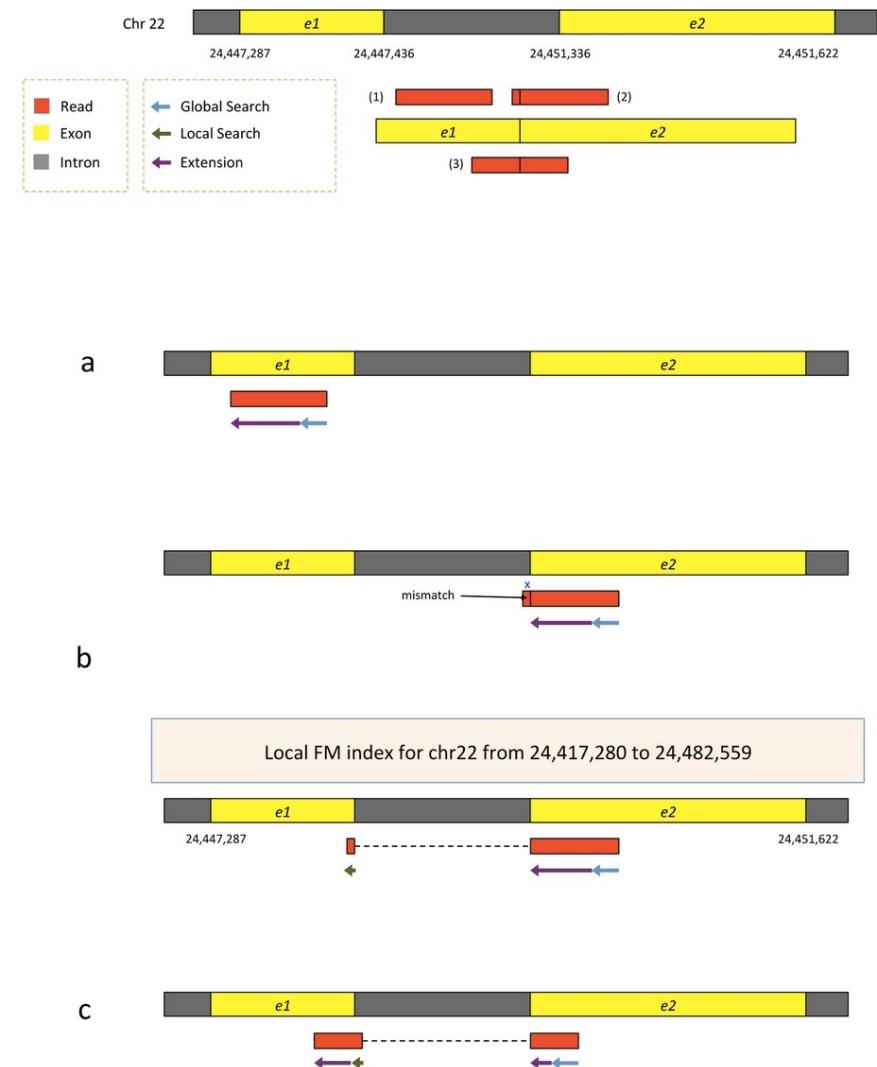
- HISAT is a ‘splice-aware’ RNA-seq read aligner
- Requires a reference genome
- Very fast
- Uses an indexing scheme based on the Burrows-Wheeler transform and the Ferragina-Manzini (FM) index
- Multiple types of indexes for alignment
 - a whole-genome FM index to anchor each alignment
 - numerous local FM indexes for very rapid extensions of these alignments.
 - Whole-genome indices with SNPs and known transcript structures accounted for



Kim et al. 2015. Nat Methods 12:357–360

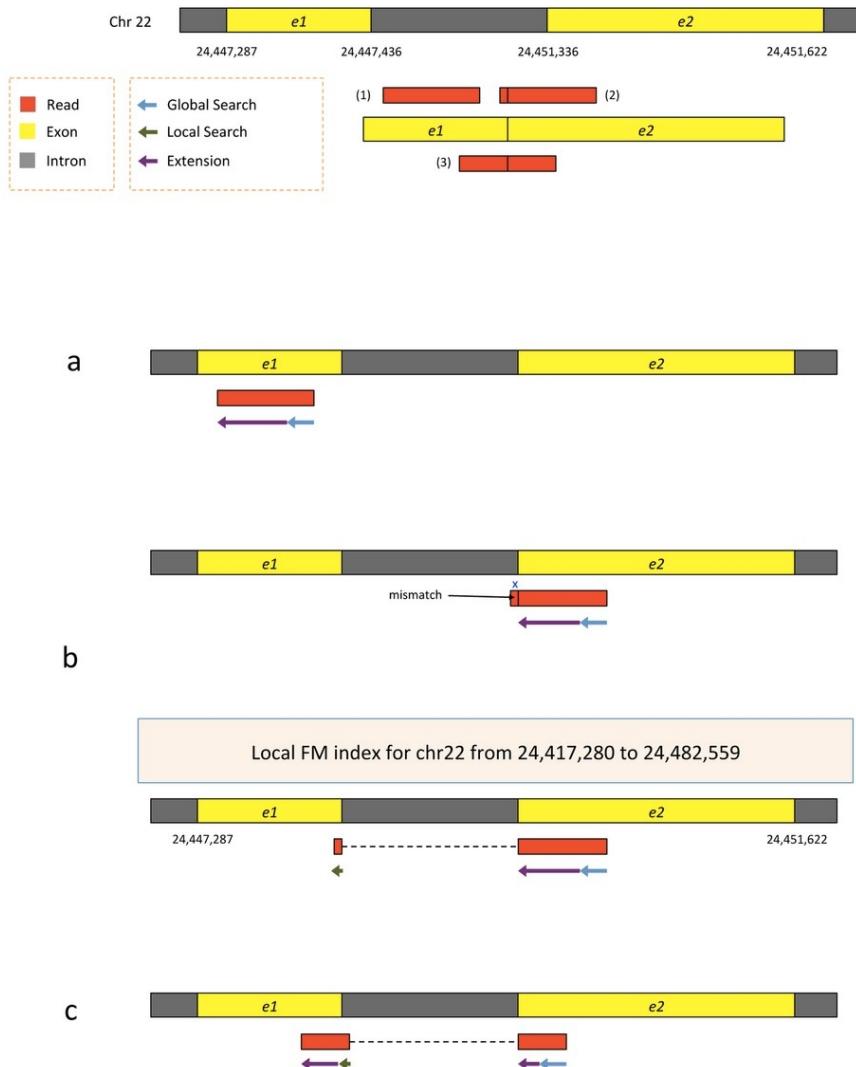
HISAT/HISAT2

- Uses hierarchical indexing algorithm and several adaptive strategies, based on the position of a read with respect to splice sites
- First tries to find candidate locations across the target genome from which the read may have originated by mapping part of each read using the global FM index, which in most cases identifies one or a small number of candidates
- Then selects one of ~48,000 local indexes for each candidate and uses it to align the remainder of the read
- For paired reads, each mate is separately aligned and the alignments of both mates are combined
 - If a read fails to align, then the alignments of its mate are used as anchors to map the unaligned mate



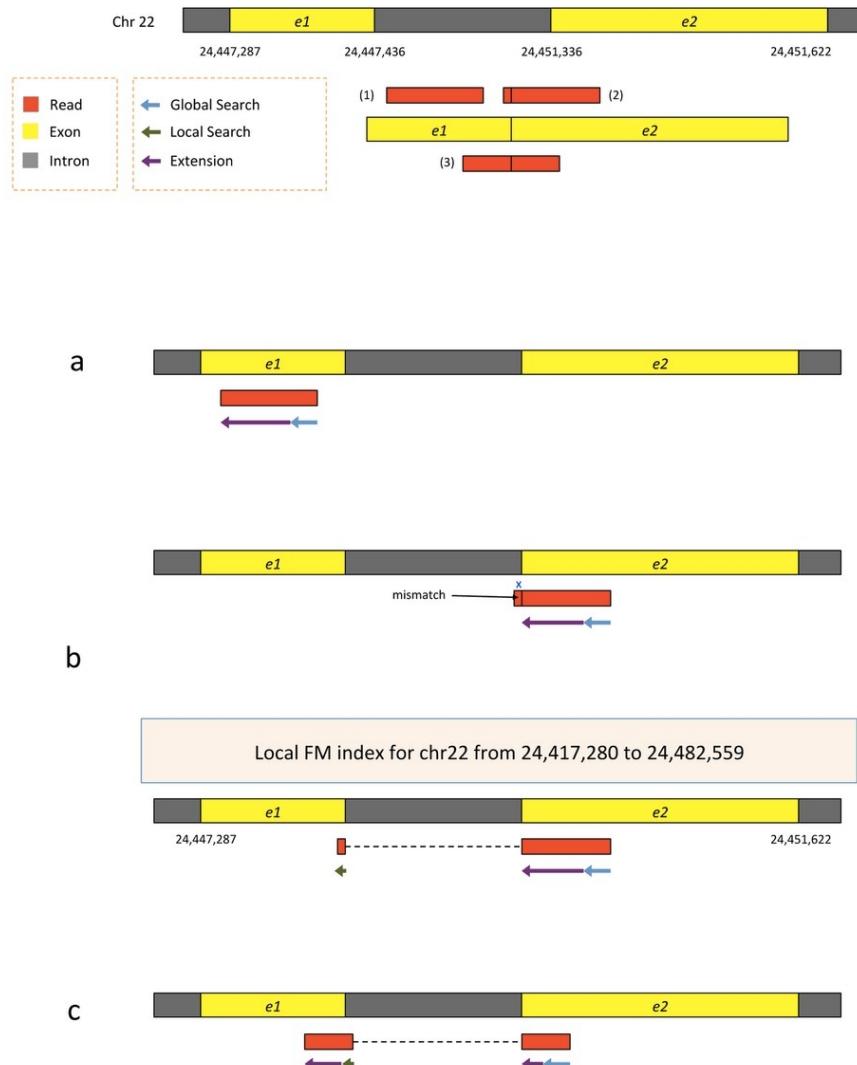
Kim et al. 2015. Nat Methods 12:357–360

HISAT/HISAT2



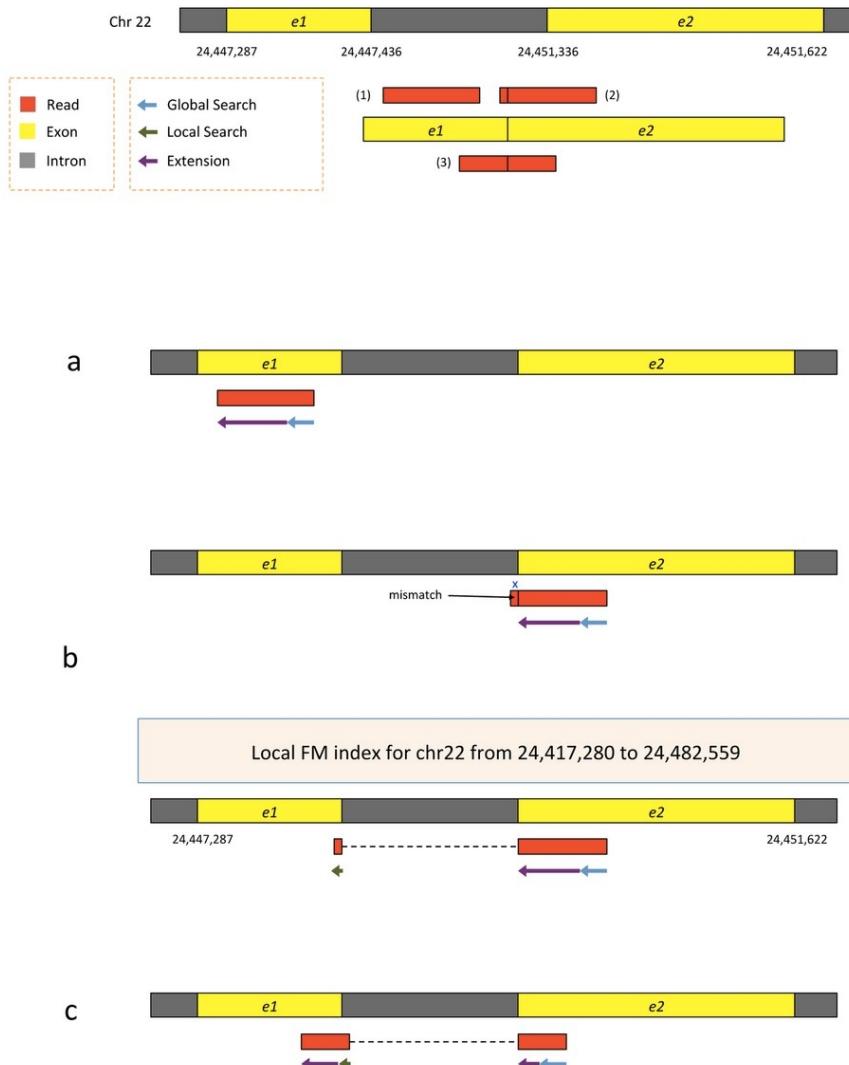
- First align read with global index (slower)
- Once at least 28bp and exactly one location switch to extension mode against reference genome (faster)

HISAT/HISAT2



- Again use global search until exactly one match of at least 28bp (slower)
- Extend as before until mismatch at 93bp (faster)
- Switch to local FM index to align remaining 8bp
 - Because the index covers only a small region, in this case we find just one match for the 8bp segment.
- Check for compatibility and combine into single spliced alignment

HISAT/HISAT2



- Again use global search until exactly one match of at least 28bp (slower)
- Extend as before until mismatch at 51bp (faster)
- Switch to local FM index to align first 8bp of remaining read
 - If too many matches increase prefix size
- Extend again
- Check for compatibility and combine into single spliced alignment

Should I allow ‘multi-mapped’ reads?

- Depends on the application
- In *DNA* analysis it is common to use a mapper to randomly select alignments from a series of equally good alignments
- In *RNA* analysis this is less common
 - Perhaps disallow multi-mapped reads if you are variant calling
 - Definitely should allow multi-mapped reads for expression analysis with Cufflinks (StringTie)
 - Definitely should allow multi-mapped reads for gene fusion discovery

What is the output of bowtie/tophat/HISAT2?

- A SAM/BAM file
 - SAM stands for Sequence Alignment/Map format
 - BAM is the binary version of a SAM file
- Remember, compressed files require special handling compared to plain text files
- How can I convert BAM to SAM?
 - <http://www.biostars.org/p/1701/>

Example of SAM/BAM file format

Example SAM/BAM header section (abbreviated)

```
mrg riffit@linus270 ~> samtools view -H /gscmnt/gc13001/info/model_data/2891632684/build136494552	alignments/136080019.bam | grep -P "SN:22|HD|RG|PG"
@HD VN:1.4 SO:coordinate
@SQ SN:22 LN:51304566 UR:ftp://ftp.ncbi.nih.gov/genbank/genomes/Eukaryotes/vertebrates_mammals/Homo_sapiens/GRCh37/special_requests/GRCh37-lite.fa.gz AS:GRCh37-lite M5:a718aca6135fdca8357d5bfe9
4211dd SP:Homo sapiens
@RG ID:2888721359 PL:illumina PU:D1BA4ACXX.3 LB:H_KA-452198-0817007-cDNA-3-lib1 PI:365 DS:paired end DT:2012-10-03T19:00:00-0500 SM:H_KA-452198-0817007 CN:WUGSC
@PG ID:2888721359 VN:2.0.8 CL:tophat --library-type fr-secondstrand --bowtie-version=2.1.0
@PG ID:MarkDuplicates PN:MarkDuplicates PP:2888721359 VN:15.0(exported) CL:net.sf.picard.sam.MarkDuplicates INPUT=/gscmnt/gc13001/info/build_merged_alignments/merged-alignment-blade10-2-5.gsc.wustl.edu-jwalker-15434-136080019/scratch-Ilg6Y/H_KA-452198-0817007-cDNA-3-lib1-2888360300.bam OUTPUT=/gscmnt/gc13001/info/build_merged_alignments/merged-alignment-blade10-2-5.gsc.wustl.edu-jwalker-15434-136080019/scratch-Ilg6Y/H_KA-452198-0817007-cDNA-3-lib1-2888360300-post_dup.bam METRICS_FILE=/gscmnt/gc13001/info/build_merged_alignments/merged-alignment-blade10-2-5.gsc.wustl.edu-jwalker-15434-136080019/staging-liu5/H_KA-452198-0817007-cDNA-3-lib1-2888360300.metrics REMOVE_DUPLICATES=false ASSUME_SORTED=true MAX_FILE_HANDLES_FOR_READ_ENDS_MAP=95000 TMP_DIR=/gscmnt/gc13001/info/build_merged_alignments/merged-alignment-blade10-2-5.gsc.wustl.edu-jwalker-15434-136080019/scratch-Ilg6Y] VALIDATION_STRINGENCY=SILENT MAX_RECORDS_IN_RAM=500000 PROGRAM_RECORD_ID=MarkDuplicates PROGRAM_GROUP_NAME=MarkDuplicates MAX_SEQUENCES_FOR_DISK_READ_ENDS_MAP=50000 SORTING_COLLECTION_SIZE_RATIO=0.25 READ_NAME_REGEX=[a-zA-Z0-9]+:[0-9]+:[0-9]+:[0-9]+.* OPTICAL_DUPLICATE_PIXEL_DISTANCE=100 VERBOSITY=INFO QUIET=false COMPRESSION_LEVEL=5 CREATE_INDEX=false CREATE_MD5_FILE=false
mrg riffit@linus270 ~>
```

Example SAM/BAM alignment section (only 10 alignments shown)

Introduction to the BED format

- When working with BAM files, it is very common to want to examine a focused subset of the reference genome
 - e.g. the exons of a gene
- These subsets are commonly specified in ‘BED’ files
 - <https://genome.ucsc.edu/FAQ/FAQformat.html#format1>
- Many BAM manipulation tools accept regions of interest in BED format
- Basic BED format (tab separated):
 - Chromosome name, start position, end position
 - Coordinates in BED format are 0 based

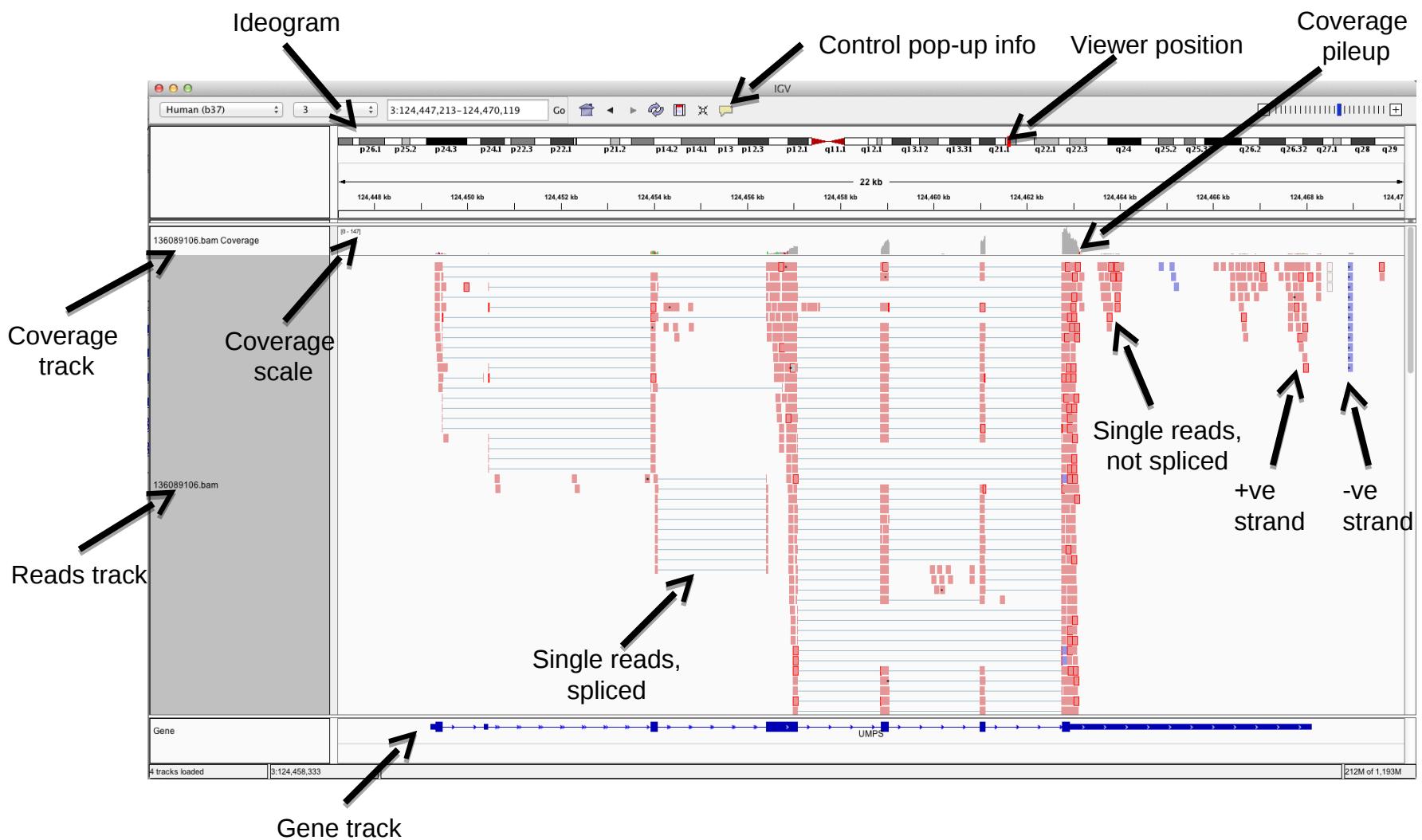
Manipulation of SAM/BAM and BED files

- Several tools are used ubiquitously in sequence analysis to manipulate these files
- SAM/BAM files
 - samtools
 - bamtools
 - picard
- BED files
 - bedtools
 - bedops

How should I sort my SAM/BAM file?

- Generally BAM files are sorted by position
 - This is for performance reasons
 - When sorted and indexed, arbitrary positions in a massive BAM file can be accessed rapidly
- Certain tools require a BAM sorted by read name
 - Usually this is when we need to easily identify both reads of a pair
 - The insert size between two reads may be large
 - In fusion detection we are interested in read pairs that map to different chromosomes...

Visualization of RNA-seq alignments in IGV browser



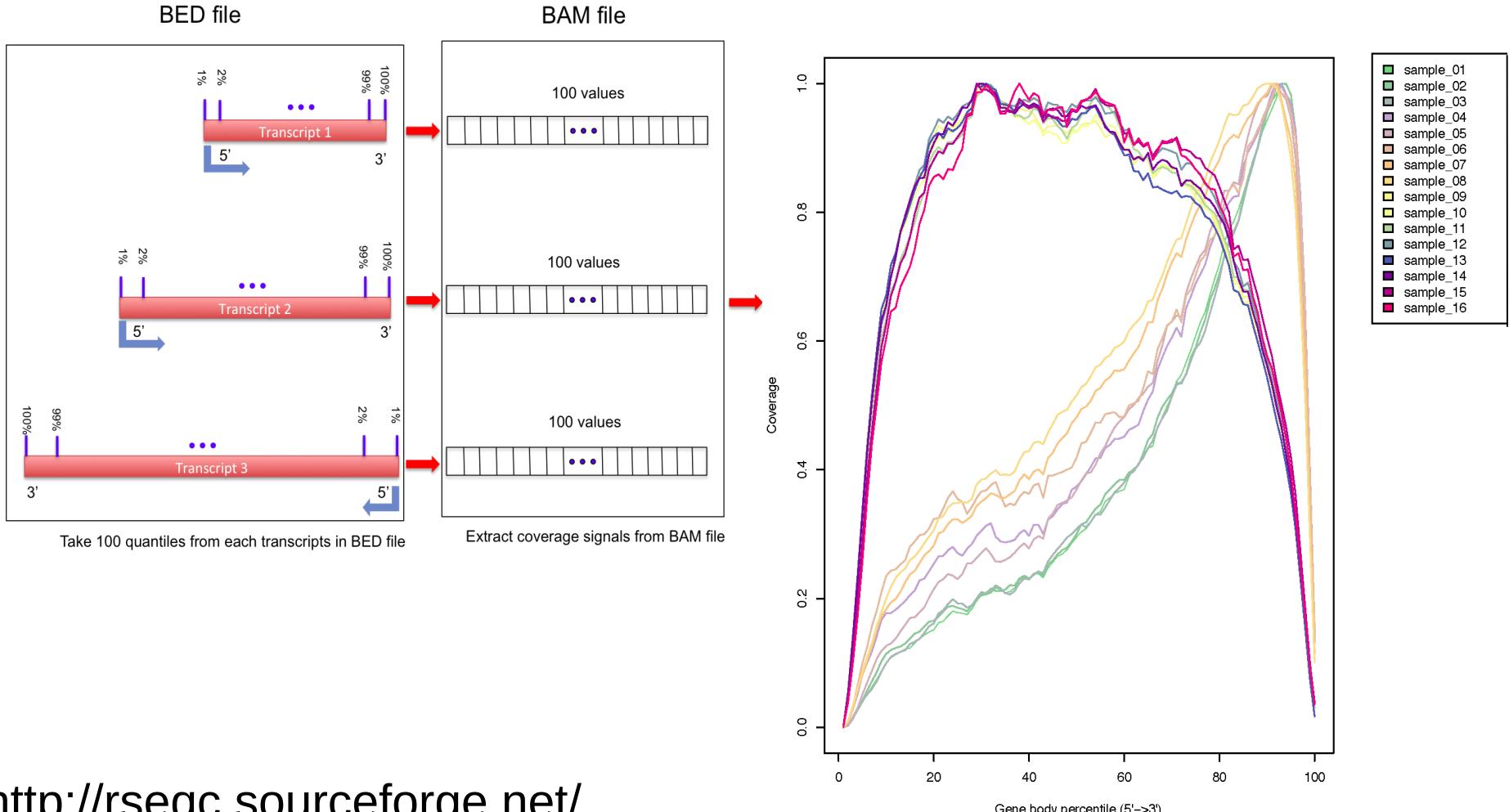
Alternative viewers to IGV

- Alternative viewers to IGV
 - <http://www.biostars.org/p/12752/>
 - <http://www.biostars.org/p/71300/>
- Artemis, BamView, Chipster, gbrowse2, GenoViewer, MagicViewer, **Savant**, Tablet, tview

Alignment QC Assessment

- 3' and 5' Bias
- Nucleotide Content
- Base/Read Quality
- PCR Artifact
- Sequencing Depth
- Base Distribution
- Insert Size Distribution

Alignment QC: 3' & 5' Bias



<http://rseqc.sourceforge.net/>

Alignment QC: Nucleotide Content

- **Random primers** are used to reverse transcribe RNA fragments into double-stranded complementary DNA (dsDNA)
- Causes certain patterns to be over represented at the beginning (5'end) of reads
- Deviation from expected A% = C% = G% = T% = 25%

Journal List > Nucleic Acids Res > v.38(12); 2010 Jul > PMC2896536

Nucleic Acids Research

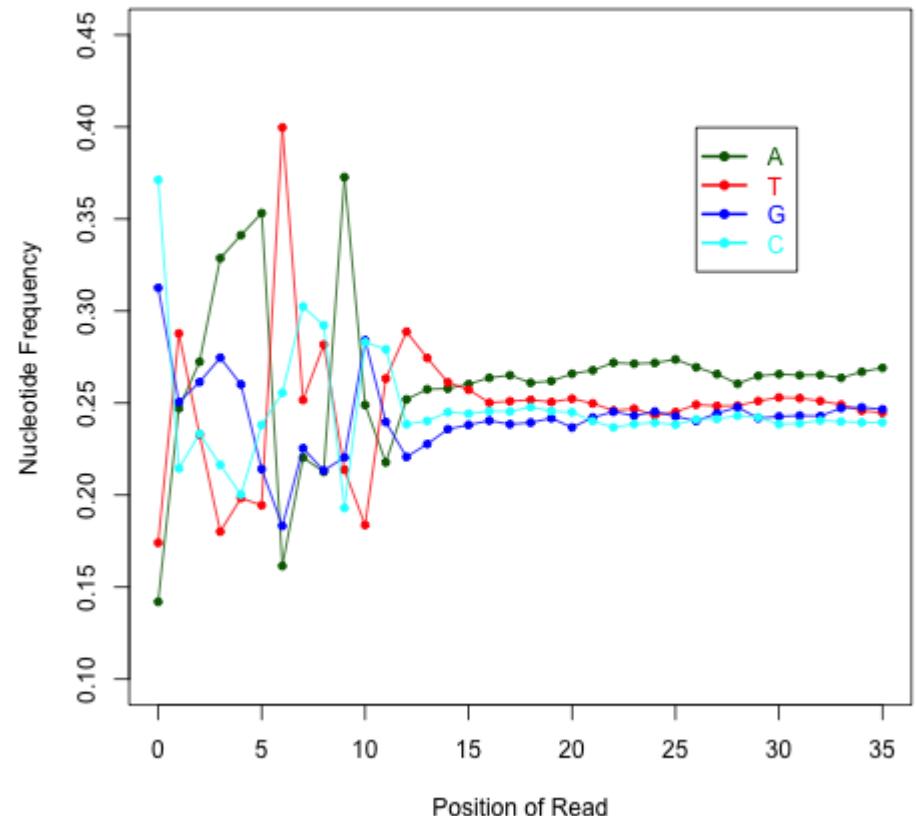
Nucleic Acids Res. 2010 Jul; 38(12): e131.
Published online 2010 Apr 14. doi: [10.1093/nar/gkq224](https://doi.org/10.1093/nar/gkq224)

Biases in Illumina transcriptome sequencing caused by random hexamer priming

Kasper D. Hansen,^{1,*} Steven E. Brenner,² and Sandrine Dudoit^{1,3}

[Author information ▾](#) [Article notes ▾](#) [Copyright and License information ▾](#)

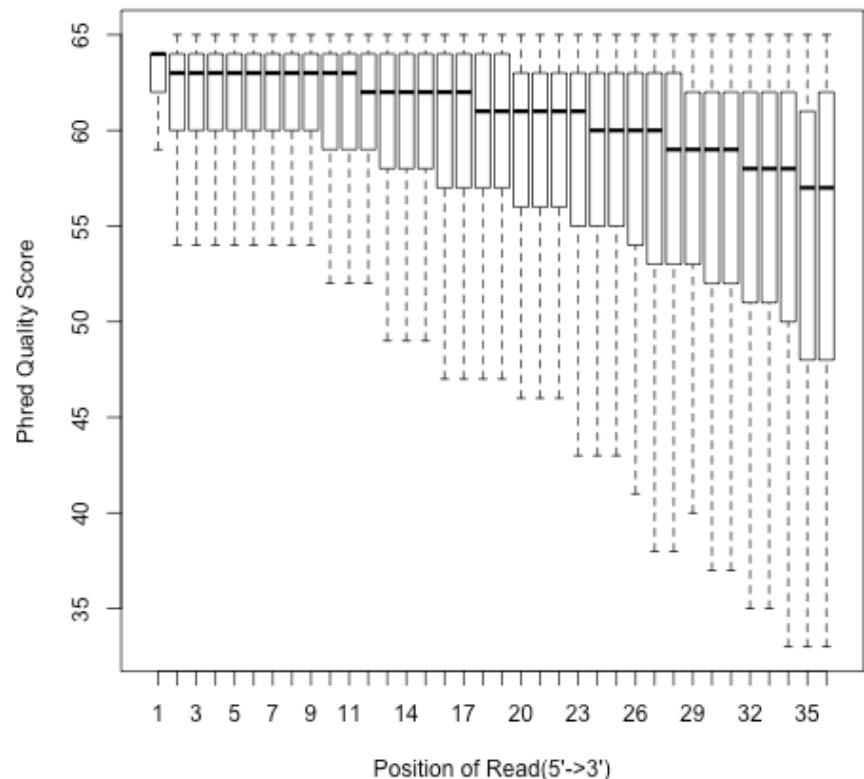
This article has been [cited by](#) other articles in PMC.



<http://rseqc.sourceforge.net/>

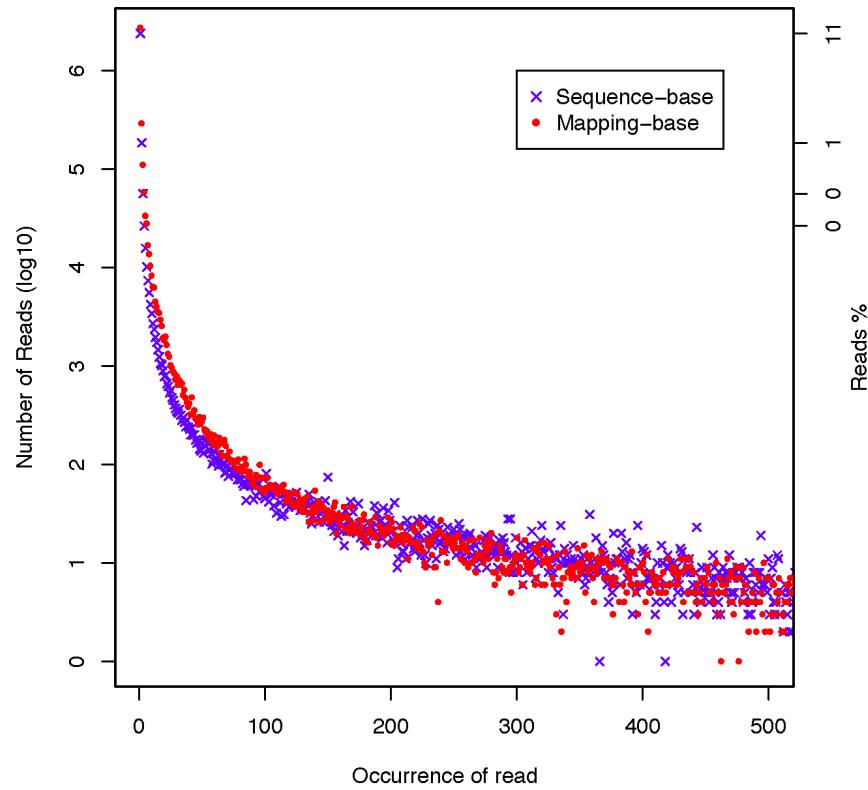
Alignment QC: Quality Distribution

- Phred quality score is widely used to characterize the quality of base-calling
- Phred quality score = $-10 \times \log(10)P$, here P is probability that base-calling is wrong
- Phred score of 30 means there is 1/1000 chance that the base-calling is wrong
- The quality of the bases tend to drop at the end of the read, a pattern observed in sequencing by synthesis techniques



Alignment QC: PCR Duplication

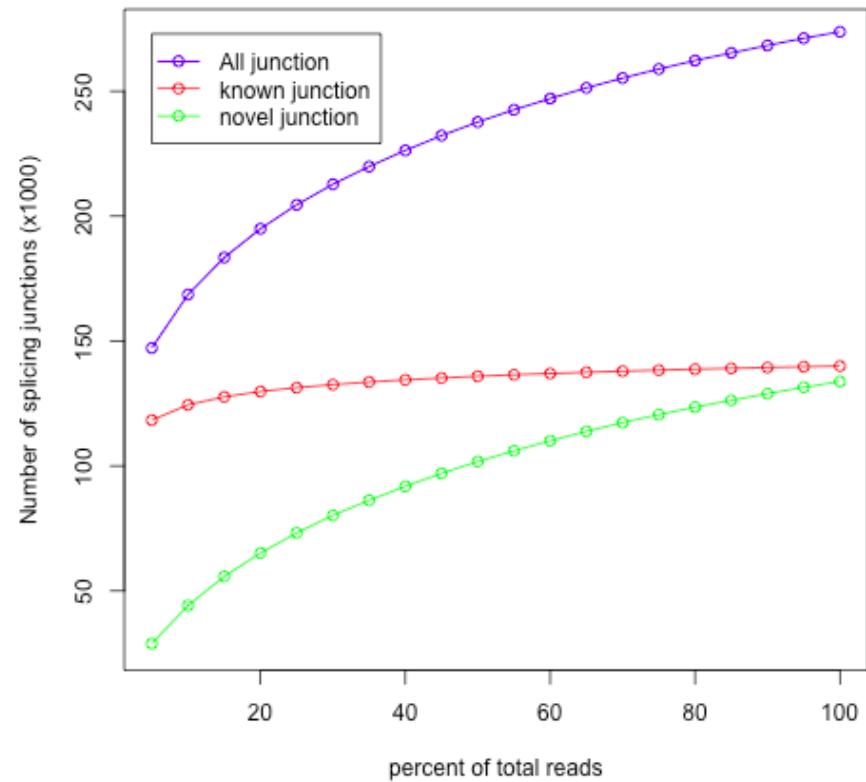
- Duplicate reads are reads that have the same start/end positions and same exact sequence
- In DNA-seq, reads/start point is used as a metric to assess PCR duplication rate
- In DNA-seq, duplicate reads are collapsed using tools such as picard
- How is RNA-seq different from DNA-seq?



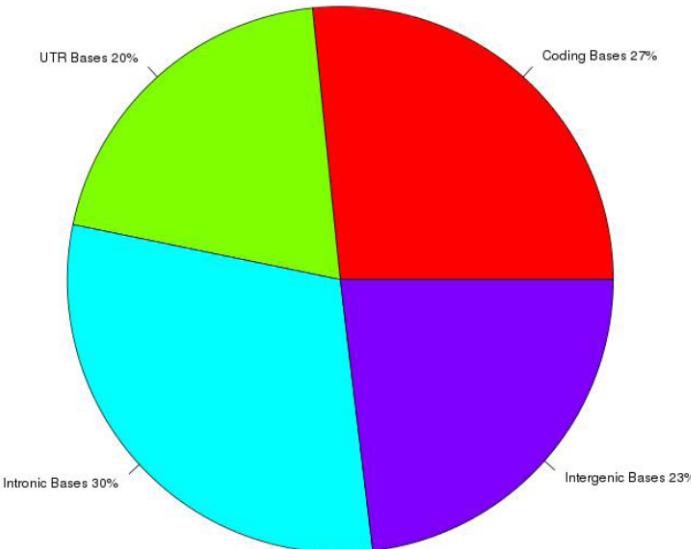
<http://rseqc.sourceforge.net>

Alignment QC: Sequencing Depth

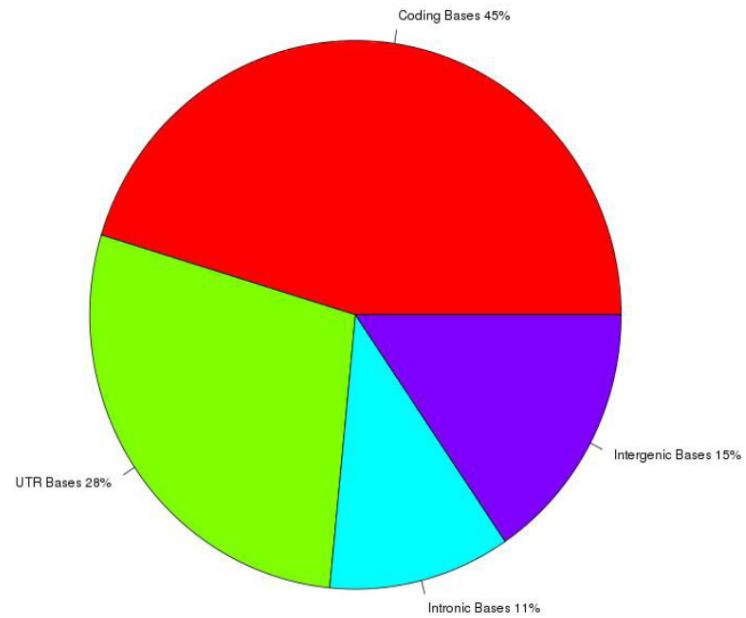
- Have we sequenced deep enough?
- In DNA-seq, we can determine this by looking at the average coverage over the sequenced region. Is it above a certain threshold?
- In RNA-seq, this is a challenge due to the variability in gene abundance
- Use splice junctions detection rate as a way to identify desired sequencing depth
- Check for saturation by resampling 5%, 10%, 15%, ..., 95% of total alignments from aligned file, and then detect splice junctions from each subset and compare to reference gene model.
- This method ensures that you have sufficient coverage to perform alternative splicing analyses



Alignment QC: Base Distribution



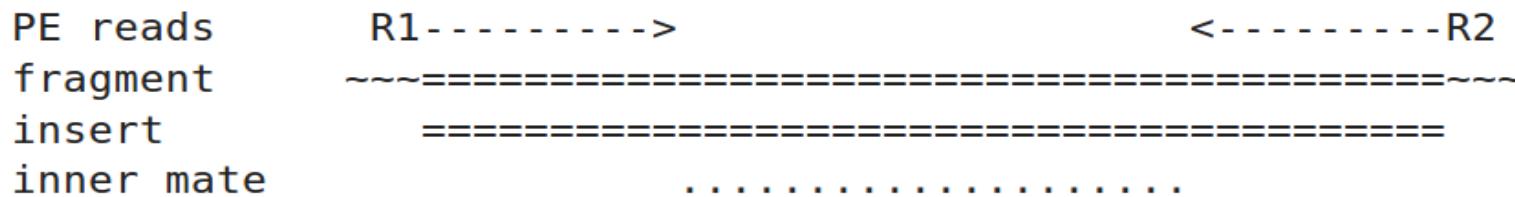
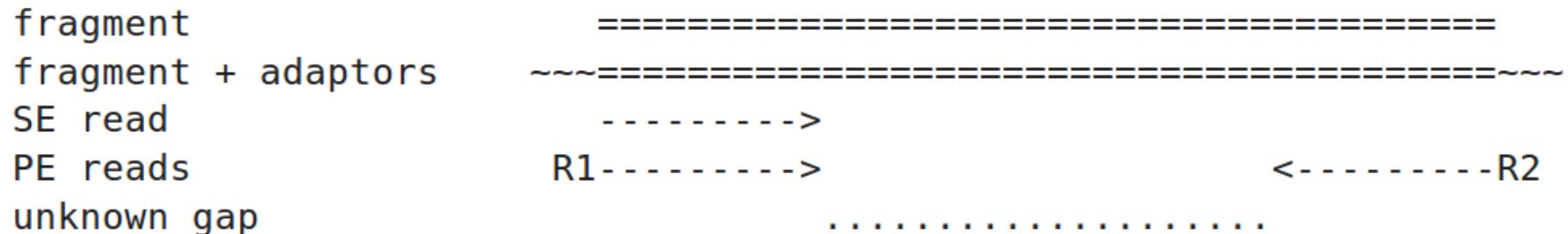
Whole Transcriptome Library



PolyA mRNA library

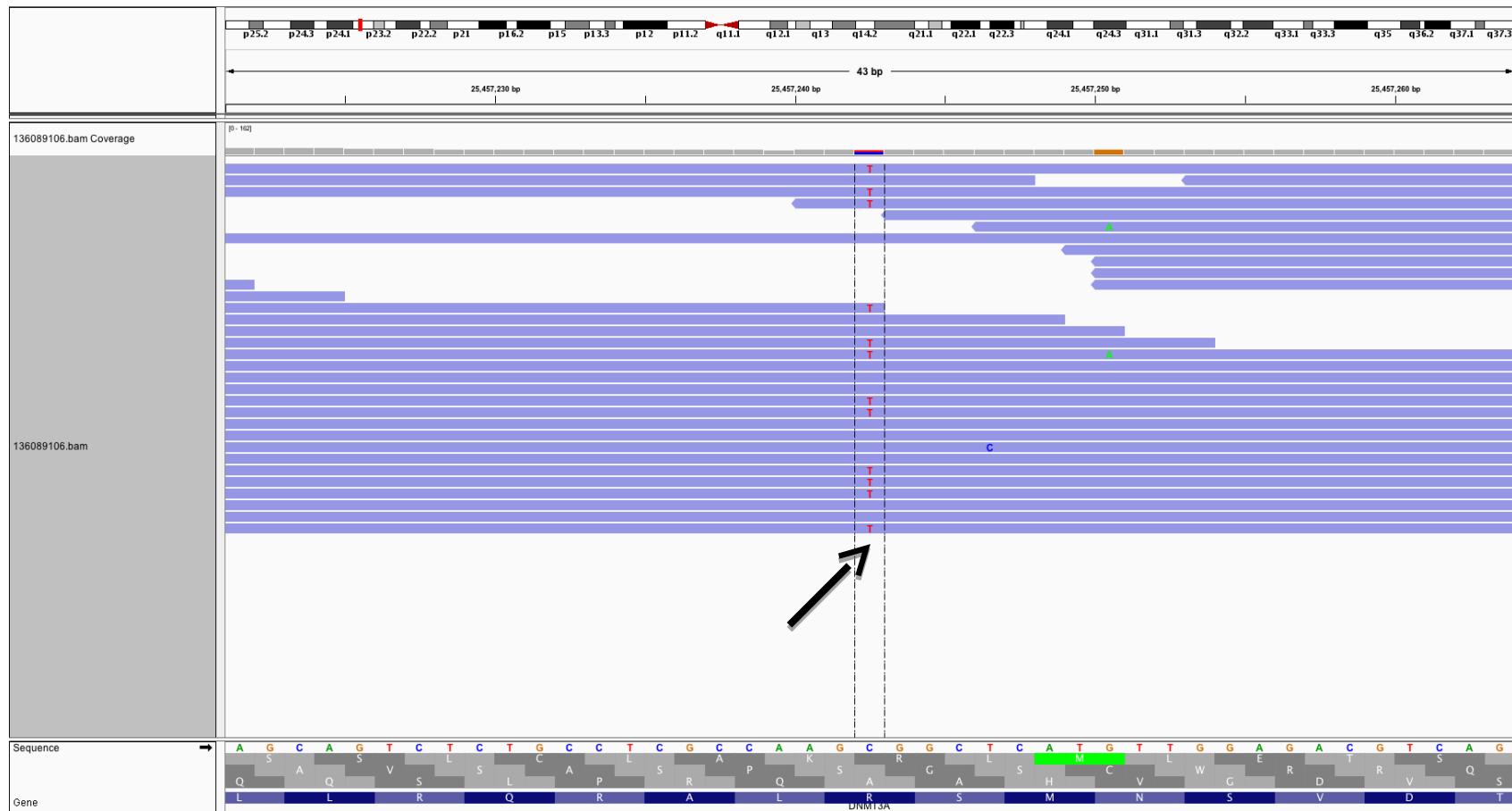
- Your sequenced bases distribution will depend on the library preparation protocol selected

Alignment QC: Insert Size



<http://thegenomefactory.blogspot.ca/2013/08/paired-end-read-confusion-library.html>

BAM read counting and variant allele expression status



- A variant C->T is observed in 12 of 25 reads covering this position. Variant allele frequency (VAF) $12/25 = 48\%$.
- Both alleles appear to be expressed equally (not always the case) -> heterozygous, no allele specific expression
- How can we determine variant read counts, depth of coverage, and VAF without manually viewing in IGV?

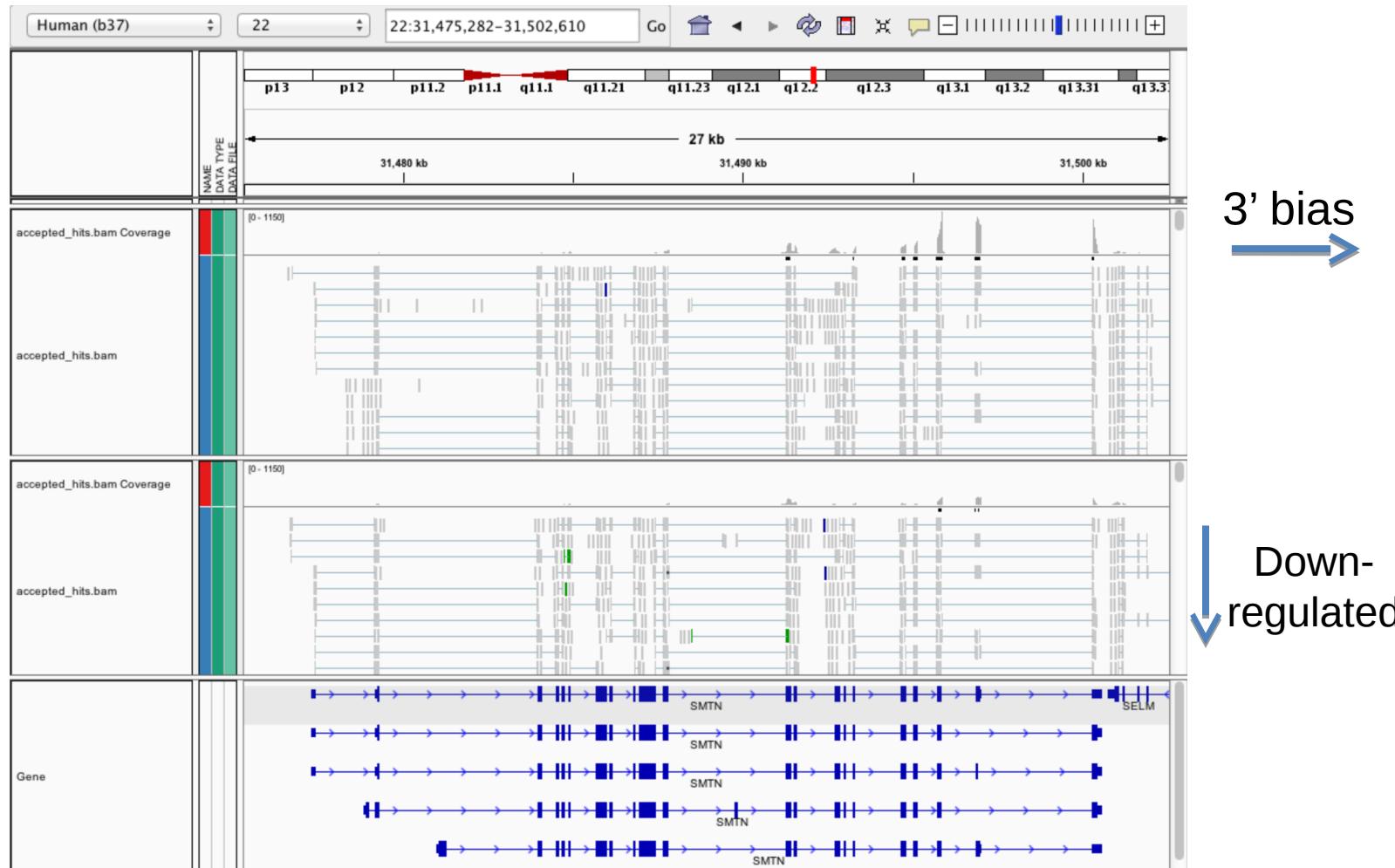
Learning Objectives of The Module

- Part I: Introduction to RNA sequencing
- Part II: RNA-seq alignment and visualization
- **Part III: Expression and Differential Expression**

Learning Objectives of Part III

- Expression estimation for known genes and transcripts
- ‘FPKM’ expression estimates vs. ‘raw’ counts
- Differential expression methods
- Downstream interpretation of expression and differential estimates
 - multiple testing, clustering, heatmaps, classification, pathway analysis, etc.

Expression estimation for known genes and transcripts



What is FPKM (RPKM)

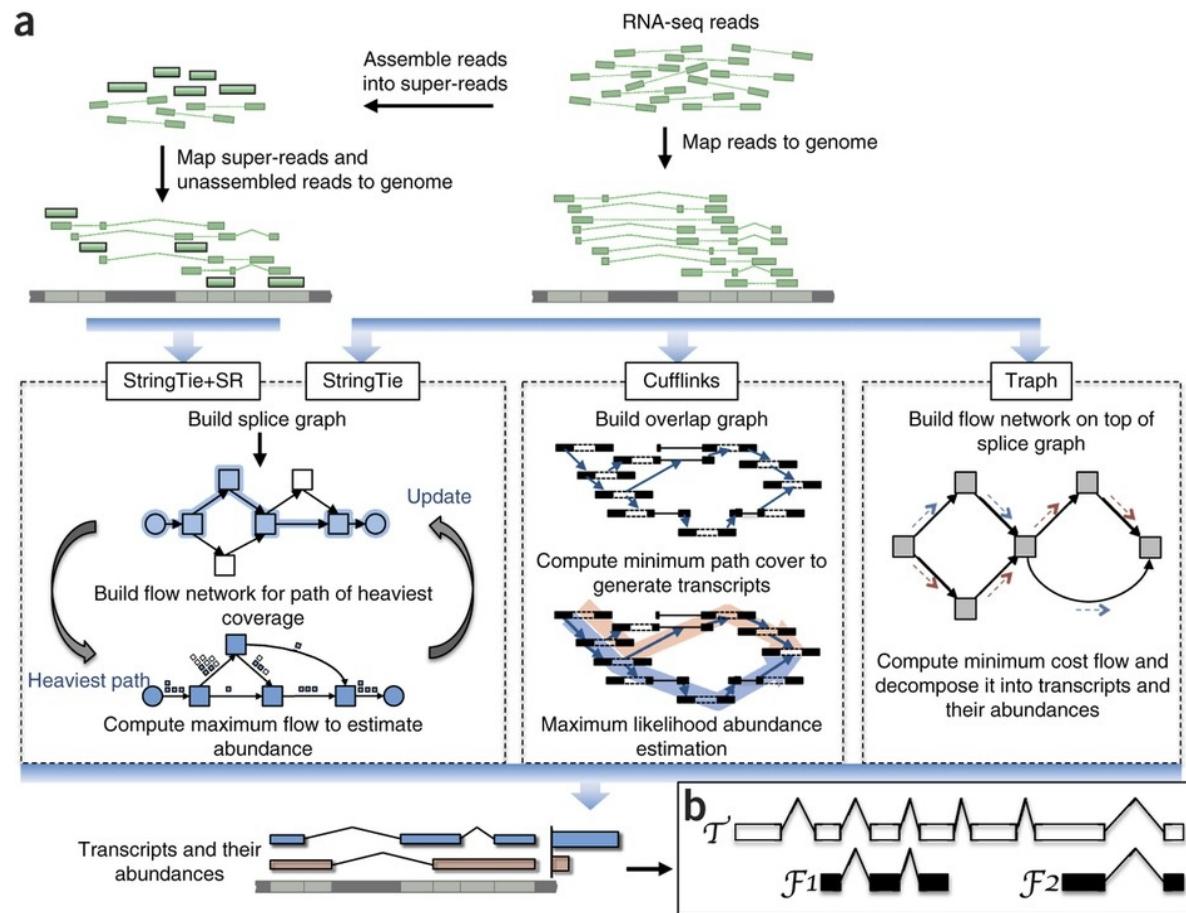
- RPKM: Reads Per Kilobase of transcript per Million mapped reads.
- FPKM: Fragments Per Kilobase of transcript per Million mapped reads.
- In RNA-Seq, the relative expression of a transcript is proportional to the number of cDNA fragments that originate from it. However:
 - The number of fragments is also biased towards larger genes
 - The total number of fragments is related to total library depth
- FPKM (or RPKM) attempt to normalize for gene size and library depth
- $\text{RPKM} \text{ (or FPKM)} = (10^9 * C) / (N * L)$
 - C = number of mappable reads/fragments for a gene/transcript/exon/etc
 - N = total number of mappable reads/fragments in the library
 - L = number of base pairs in the gene/transcript/exon/etc
- <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) Sum sample/library fragments per million
 - 2) Divide gene/transcript fragment count by #1
 - fragments per million, FPM
 - 3) Divide FPM by length of gene in kilobases (FPKM)
 - TPM
 - 1) Divide fragment count by length of transcript
 - fragments per kilobase, FPK
 - 2) Sum all FPK for sample/library per million
 - 3) Divide #1 by #3 (TPM)
- <http://www.rna-seqblog.com/rpkf-fpkf-and-tpm-clearly-explained/>

How does StringTie work?

- StringTie iteratively extracts the heaviest path from a splice graph, constructs a flow network, computes maximum flow to estimate abundance, and then updates the splice graph by removing reads that were assigned by the flow algorithm. This process repeats until all reads have been assigned.
- Annotated transcript T for which read data covers only the fragments F1 and F2.



Pertea et al. Nature Biotechnology, 2015

StringTie -merge

- Merge together all gene structures from all samples
 - Some samples may only partially represent a gene structure
- Allows for the incorporation of 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>

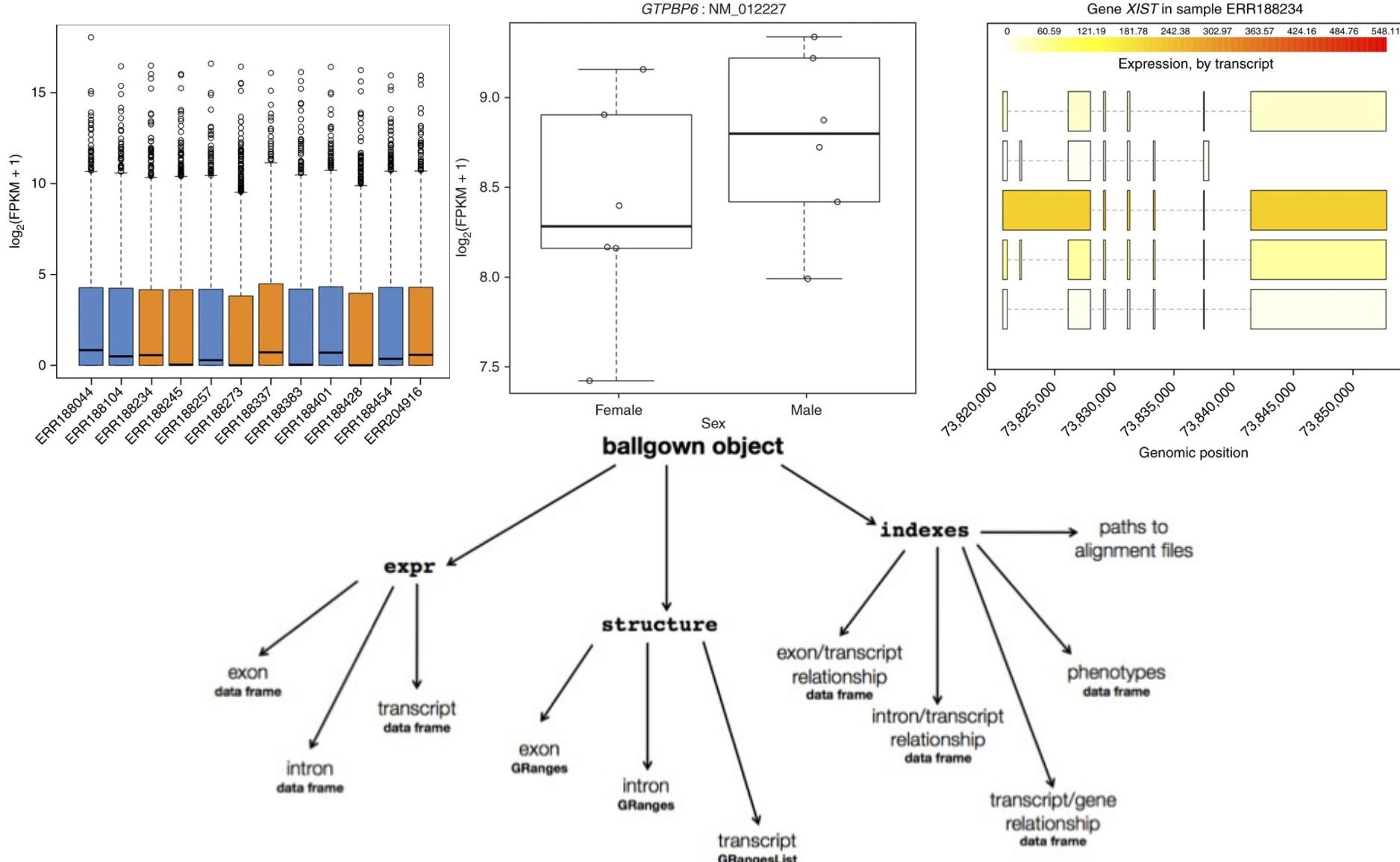
Priority	Code	Description
1	=	Complete match of intron chain
2	c	Contained
3	j	Potentially novel isoform (fragment): at least one splice junction is shared with a reference transcript
4	e	Single exon transfrag overlapping a reference exon and at least 10 bp of a reference intron, indicating a possible 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)

Ballgown for Differential Expression

- 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 ~5%.

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

Ballgown for Visualization with R



Alternatives to FPKM

- Raw read counts as an alternate for differential expression analysis
 - Instead of calculating FPKM, simply assign reads/fragments to a defined set of genes/transcripts and determine “raw counts”
 - Transcript structures could still be defined by something like cufflinks
- HTSeq (htseq-count)
 - <http://www-huber.embl.de/users/anders/HTSeq/doc/count.html>
 - htseq-count --mode intersection-strict --stranded no --minaqual 1 --type exon --idattr transcript_id accepted_hits.sam chr22.gff > transcript_read_counts_table.tsv
 - Important caveat of ‘transcript’ analysis by htseq-count:
 - <http://seqanswers.com/forums/showthread.php?t=18068>

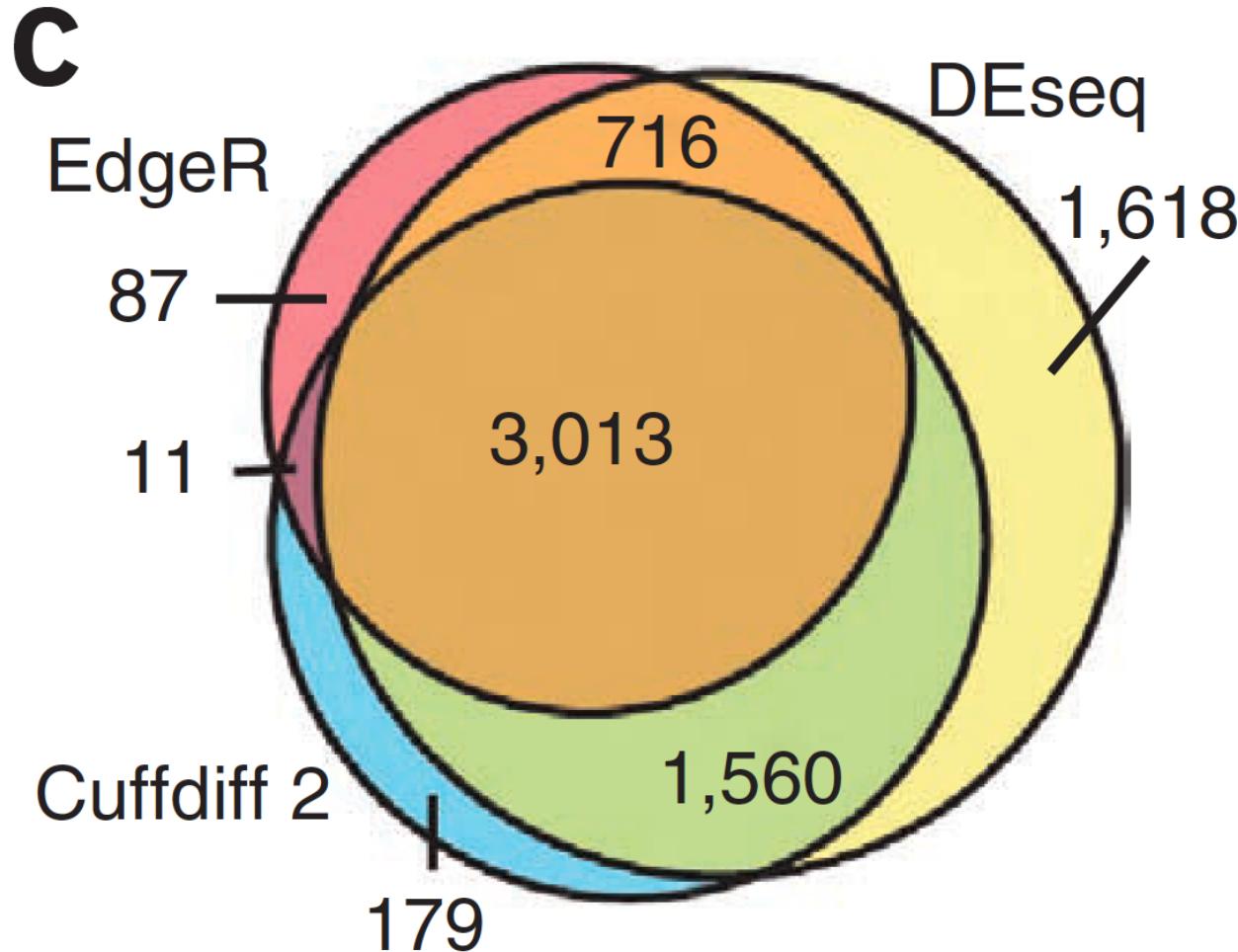
'FPKM' expression estimates vs. 'raw' counts

- Which should I use?
- FPKM
 - When you want to leverage benefits of tuxedo suite
 - 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

Alternative differential expression methods

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

Multiple approaches advisable



Multiple testing correction

- As more attributes are compared, it becomes more likely that the treatment and control groups will appear to differ on at least one attribute by random chance alone.
- 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
 - Almost infinite number of potential features
 - Genes, transcripts, exons, junctions, retained introns, microRNAs, lncRNAs, etc, 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 cufflinks/cuffdiff (or alternative) can be fed into many analysis pipelines
- See supplemental R tutorial for how to format cufflinks data and start manipulating in R
- Clustering/Heatmaps
 - Provided by cummeRbund
 - 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
 - David
 - IPA
 - Cytoscape
 - Many R/BioConductor packages: <http://www.bioconductor.org/help/search/index.html?q=pathway>

We are on a Coffee Break &
Networking Session