# RNA Sequencing

Elisha Roberson, Ph.D.
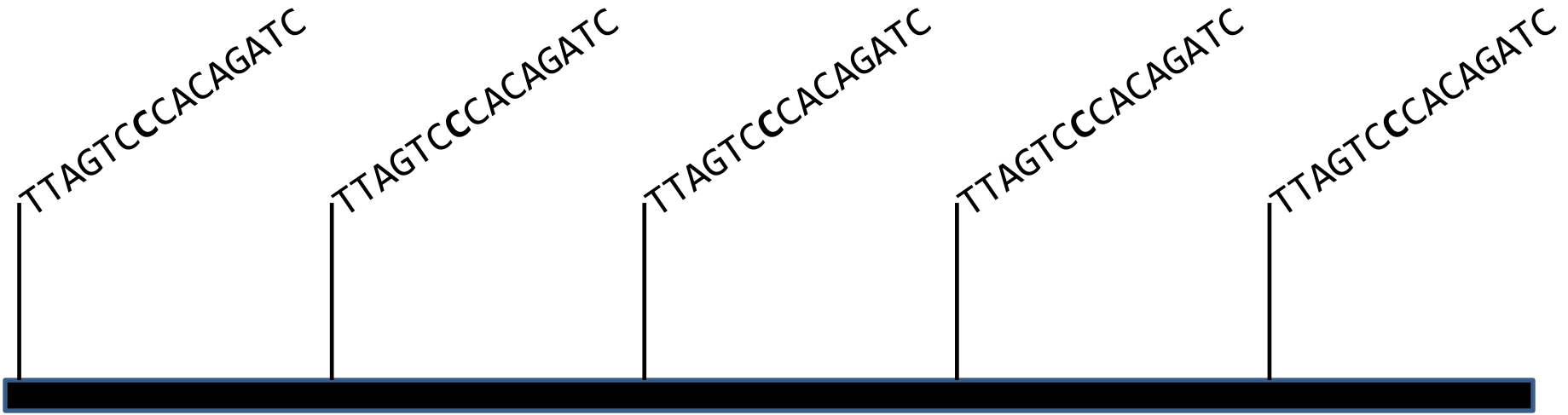
Depts. of Medicine & Genetics

eroberson@wustl.edu

2016-December-06
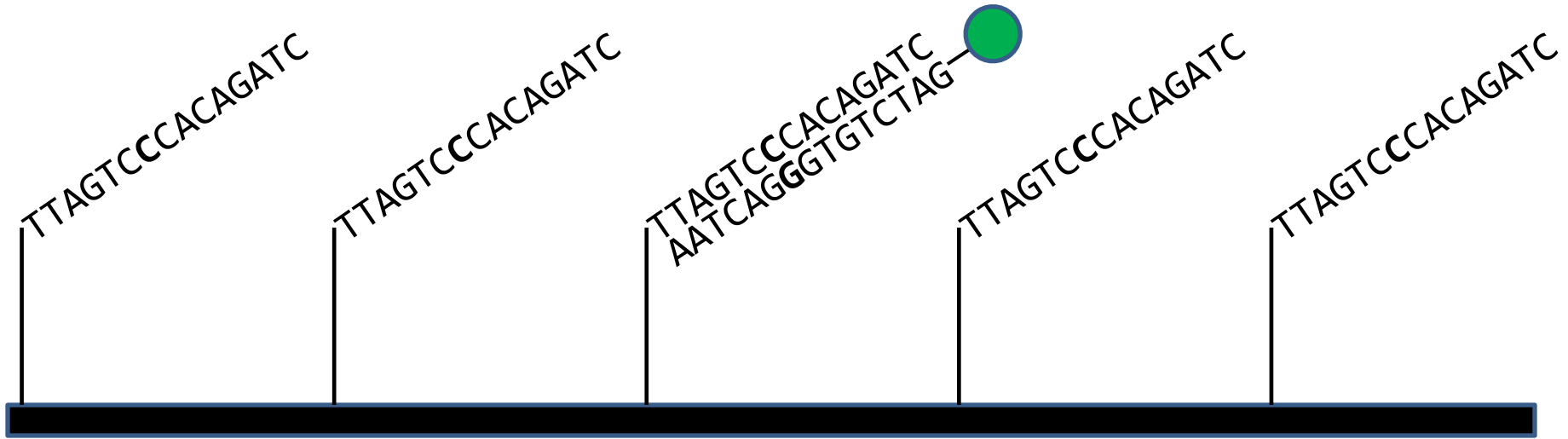
Washington University in St. Louis

The Roberson Lab

# Comparing RNA-Seq and microarrays

# RNA-Seq vs. microarrays

TTAGTC**C**CACAGATC

TTAGTC**C**CACAGATC

TTAGTC**C**CACAGATC

TTAGTC**C**CACAGATC

TTAGTC**C**CACAGATC

- Hybridization requires known targets

# RNA-Seq vs. microarrays

TTAGTC**C**CACAGATC

TTAGTC**C**CACAGATC

TTAGTC**C**CACAGATC
AATCAG**G**GTGTCTAG—

TTAGTC**C**CACAGATC

TTAGTC**C**CACAGATC

- Expression detected by fluorescence

# RNA-Seq vs. microarrays

TTAGTC**CC**ACAGATC

TTAGTCC**C**ACAGATC

TTAGTCC**C**CACAGATC
AATCAG**G**GTGTCTAG—
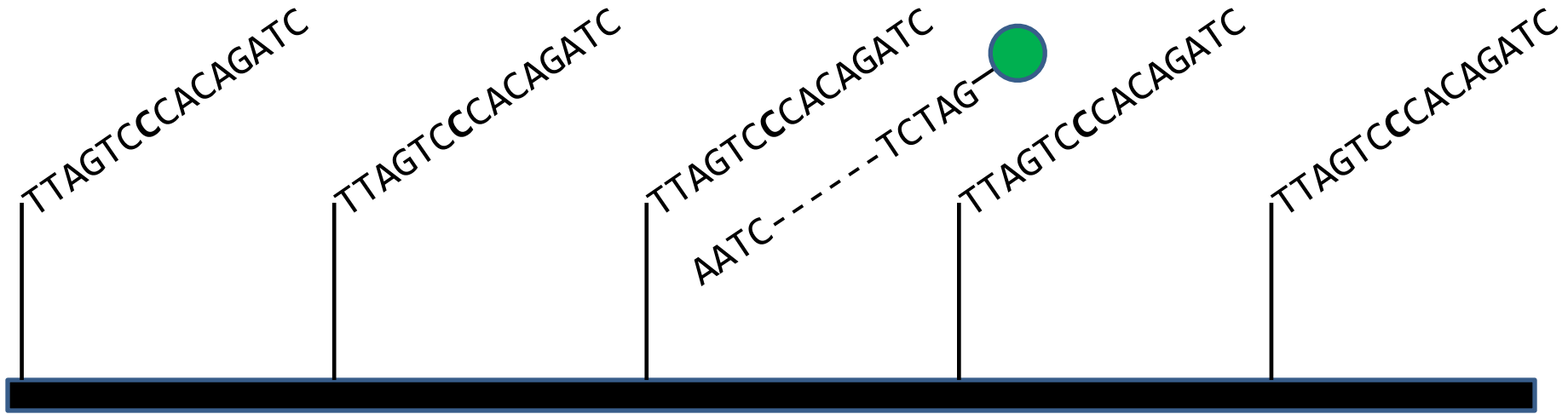
TTAGTCC**C**ACAGATC

TTAGTCC**C**ACAGATC

- Low-level expression can be difficult to detect compared to background

# RNA-Seq vs. microarrays

- High-level expression can saturate probes

# RNA-Seq vs. microarrays

TTAGTC**CC**ACAGATC

TTAGTCC**C**ACAGATC

TTAGTCC**C**ACAGATC

AATC------TCTAG—

TTAGTC**C**CACAGATC

TTAGTCC**C**ACAGATC

- Variation in the subject's RNA sequence can affect binding kinetics

# RNA-Seq vs. microarrays

**RNA-Seq differences**

• Counting – no probes to saturate (though some sequences can predominate the library) and can always sequence more to get low-expressors

• Doesn't *require* a known gene sequence

• Not affected by variation, as long as it doesn't affect transcript stability
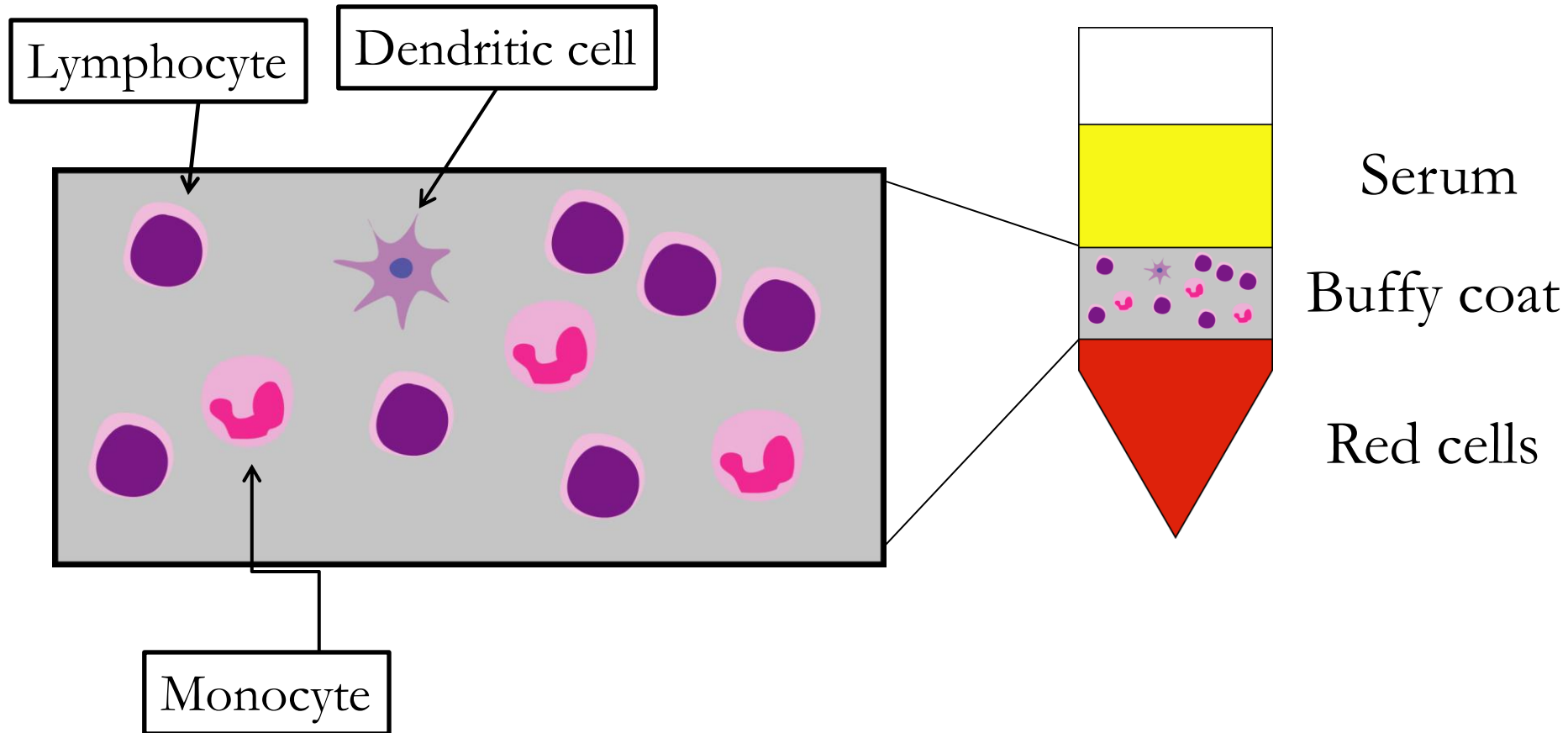
Isoform A

Isoform B

# Sample prep is *critical*

# Collection & storage



- Blood is low-risk, easy to access, and frequently used

- **But** strongly enriched for red cells. What transcripts will predominate?

# Collection & storage

Lymphocyte

Dendritic cell

Monocyte

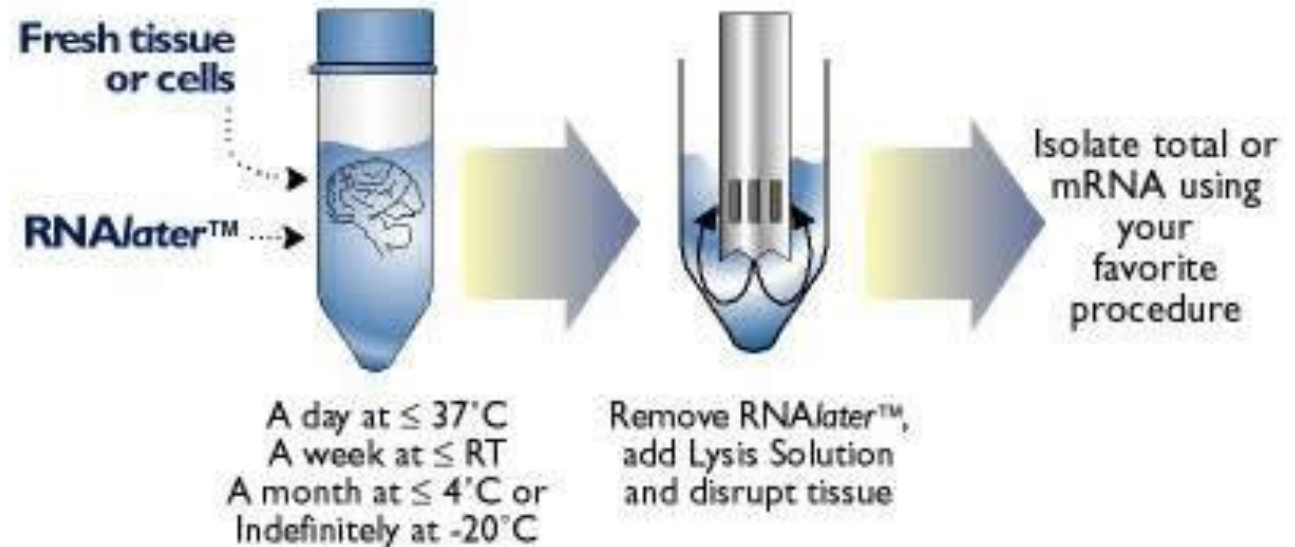Serum

Buffy coat

Red cells

- Freeze PBMCs directly
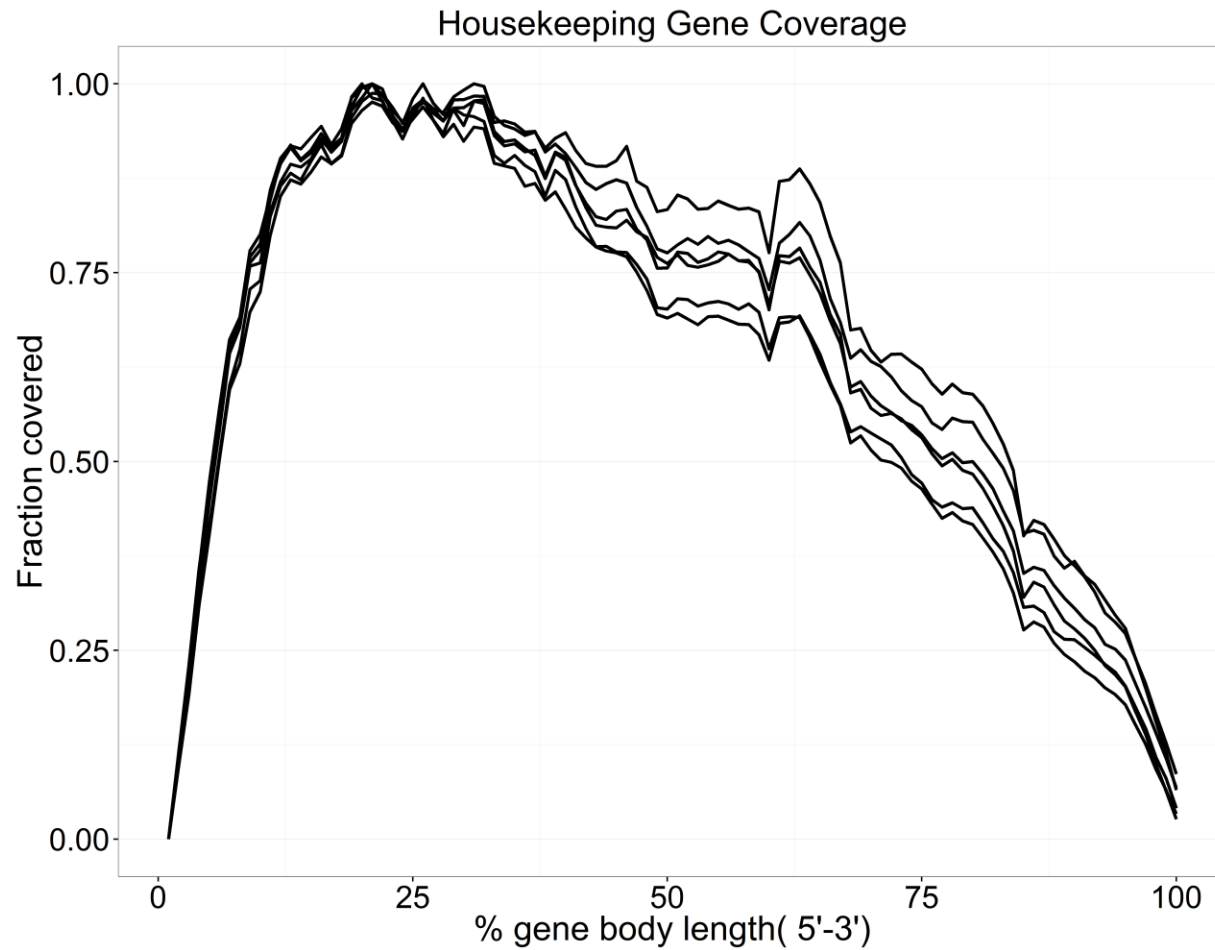- Lyse directly in Trizol

# Collection & storage



**Tissues**

- Biopsies, whole organs, etc.

- Unlike blood, the complex structure and embedded connective fibers make tissues more difficult to process and preserve.
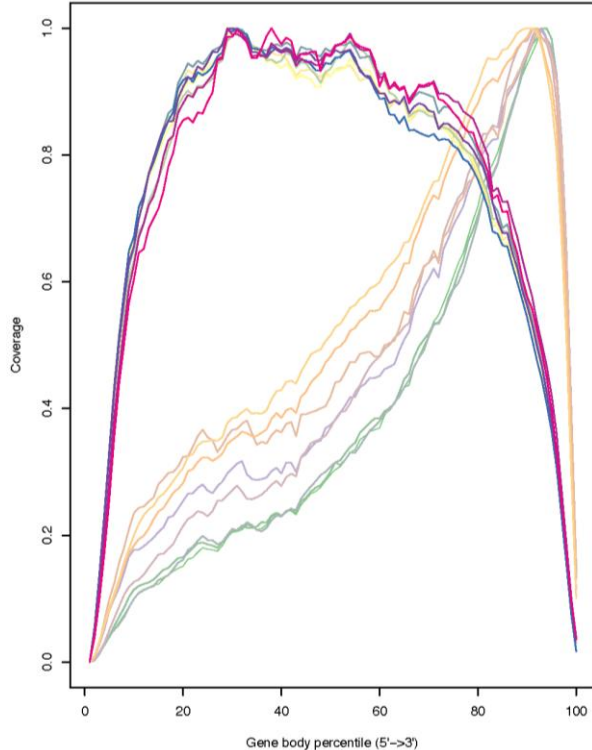
*Nischal *et al*. J Cutan Aesthet Surg. 2008 Jul-Dec; 1(2): 107–111.

# Collection & storage



- Some solutions, like RNAlater, can preserve tissue RNA for short-term storage at RT or long-term storage at -20°C

Housekeeping Gene Coverage

- PBMCs shipped on (too little) dry ice
- Degraded 3'-5'
- PolyA kit fails

# Library prep / cDNA methods

# Polyadenylation (polyA) preps



AAAAAAAAA

TTTTTTTTTTT



- polyT primed first-strand synthesis works

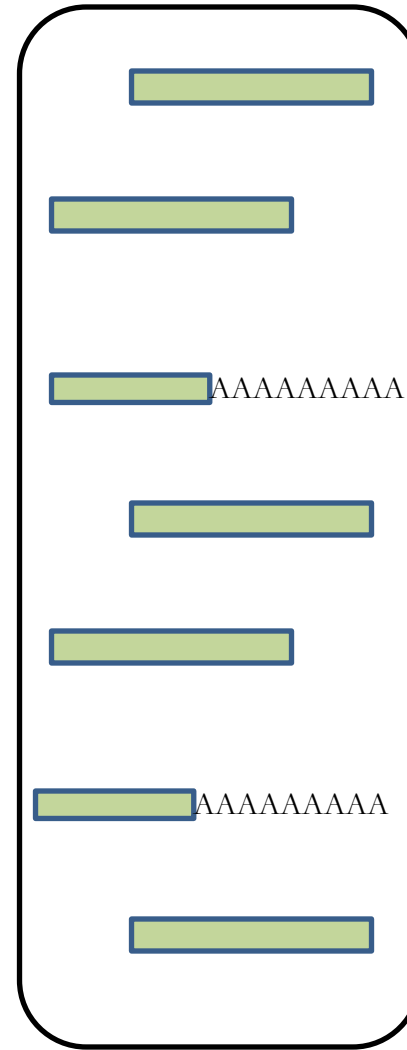- **But** can lead to 3' bias and only captures polyadenylated transcripts.
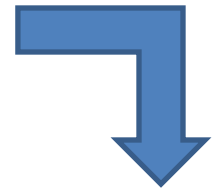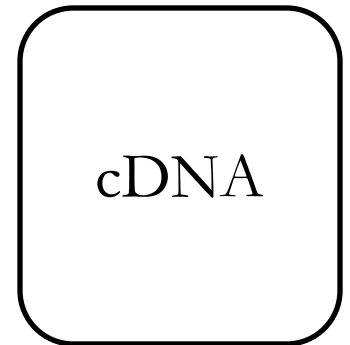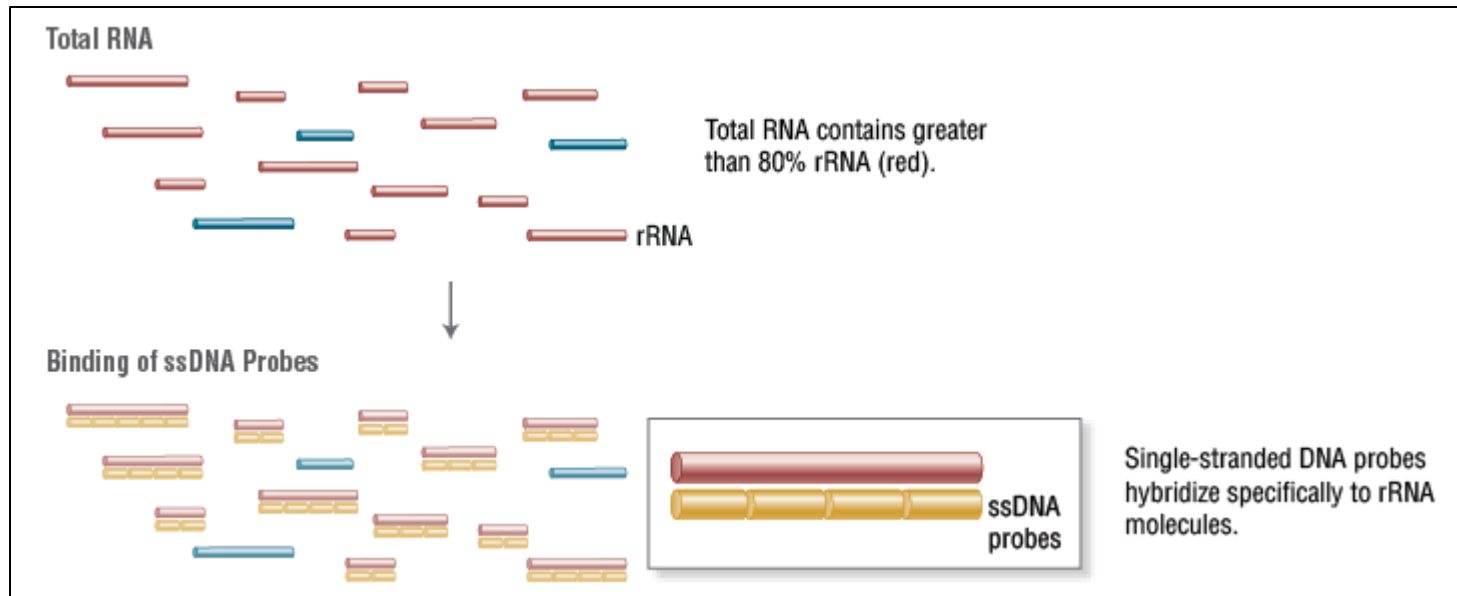
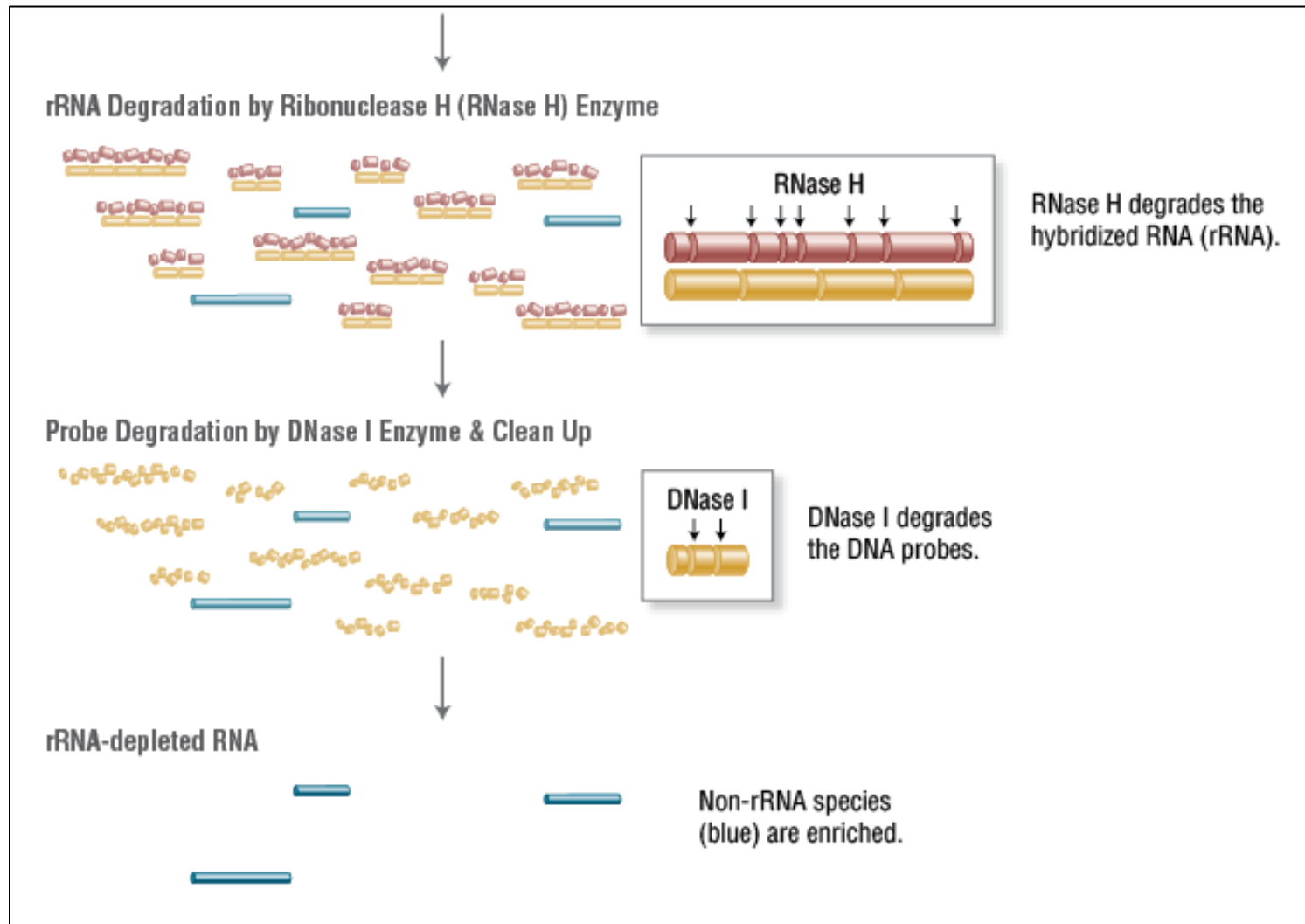*RSEQC

# Polyadenylation (polyA) preps
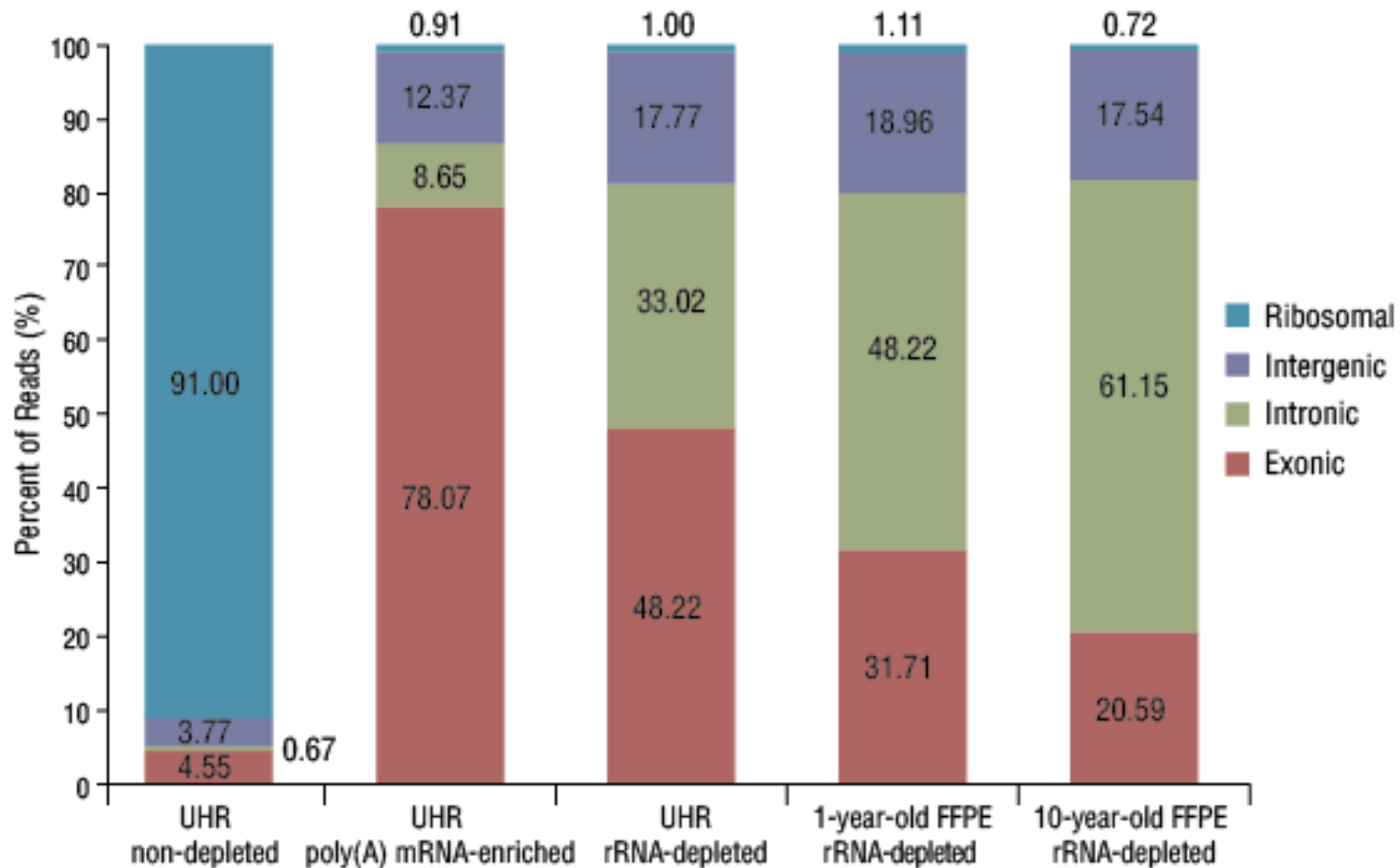
# Ribosomal depletion preps

- Majority of cellular RNA is non-coding, particularly ribosomal RNA

- RNA polymerase I (28S, 18S, 5.8S rRNA) and Pol III (5S rRNA)

- **<u>Not polyadenylated</u>**



Total RNA

Total RNA contains greater than 80% rRNA (red).

rRNA

Binding of ssDNA Probes

ssDNA probes

Single-stranded DNA probes hybridize specifically to rRNA molecules.

# Ribosomal depletion preps



rRNA Degradation by Ribonuclease H (RNase H) Enzyme

RNase H

RNase H degrades the hybridized RNA (rRNA).

Probe Degradation by DNase I Enzyme & Clean Up

DNase I

DNase I degrades the DNA probes.

rRNA-depleted RNA

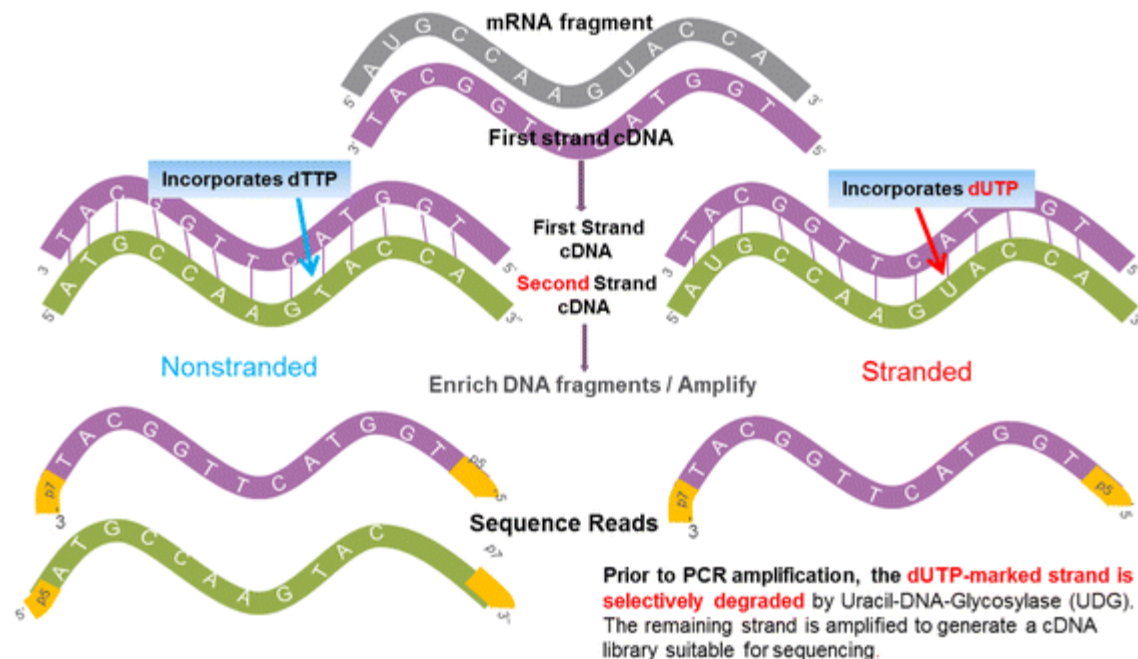Non-rRNA species (blue) are enriched.

*NEB
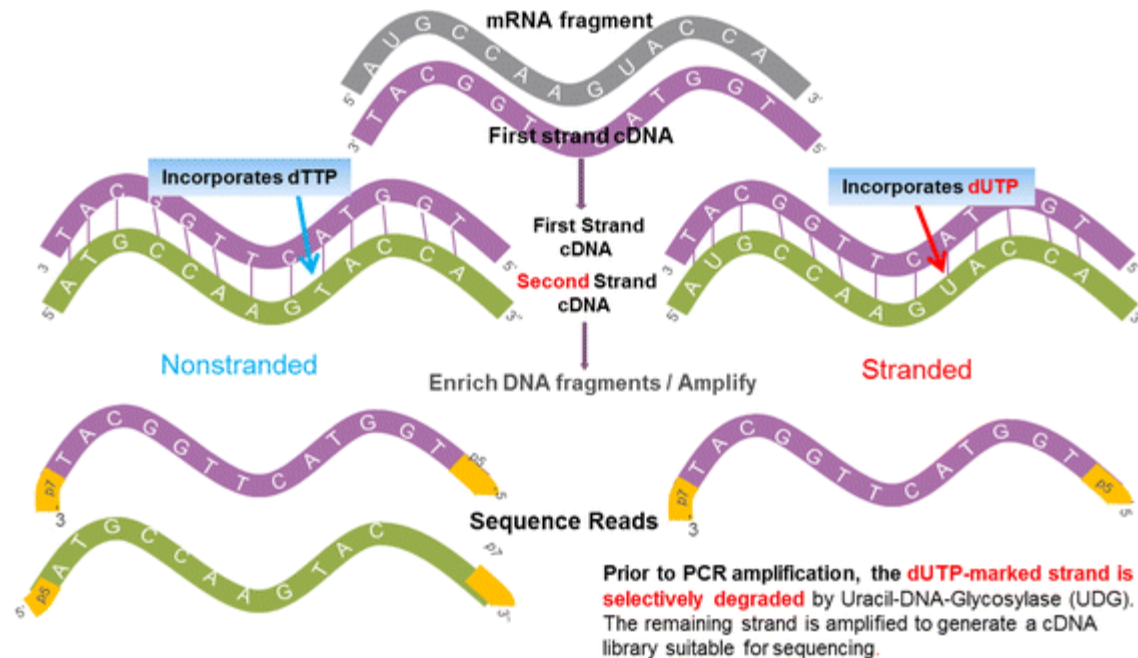
# Ribosomal depletion preps

# Stranded preps

- Standard cDNA → library prep retains no information about transcript strand.

- Some loci have antisense transcripts

# Stranded preps



- Normal first strand synthesis. 2$^{nd}$ strand incorporates uracil
- Uracil-DNA glycosylase excises U-base from DNA
- Endonuclease VIII breaks backbone at those sites

# Sequencing choices

# Short reads

- Illumina
  - Single-end vs. paired-end
    - Paired-end superior. Estimates insert size empirically

  - Read length
    - Greater cycle number preferred
    - 2x75 good compromise

  - Depth / coverage
    - Very different from DNA seq
    - Variable gene length and expression level
    - Several tools to estimate

# Long reads – Pac Bio SMRT

- Full-length isoform sequencing ($$$$$)

| Pacbio  Library Construction and Sequencing | Cost Per Sample |
|---|---|
| Sequencing SMRT Cell | $257 |
| Standard Library Prep | $560 |
| Low_input Library Prep | $603 |
| Iso-Seq Whole Transcriptome Lib_Prep | $875 |
| Iso-Seq Targeted Lib_Prep | $664 |

Min $1,100 / sample for whole transcriptome.

# Long reads – Oxford nanopore

- MinION / PromethION (also $$$$$)



FLO-MIN106

**SpotON Flow Cell Mk I (R9.4)**

$900.00

Single | 12-pack | 24-pack | 48-pack

− 1 + Buy now

SQK-LSK208

**Ligation Sequencing kit 2D (R9.4)**

$599.00

− 1 + Buy now

**Product overview**

The Ligation Sequencing Kit 2D is designed to prepare genomic, amplicon and cDNA, with or without barcoding, for sequencing on the Oxford
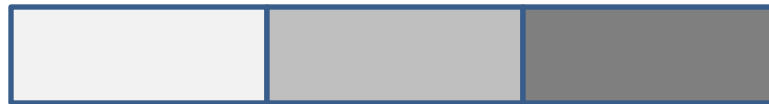
# Alignment and counting

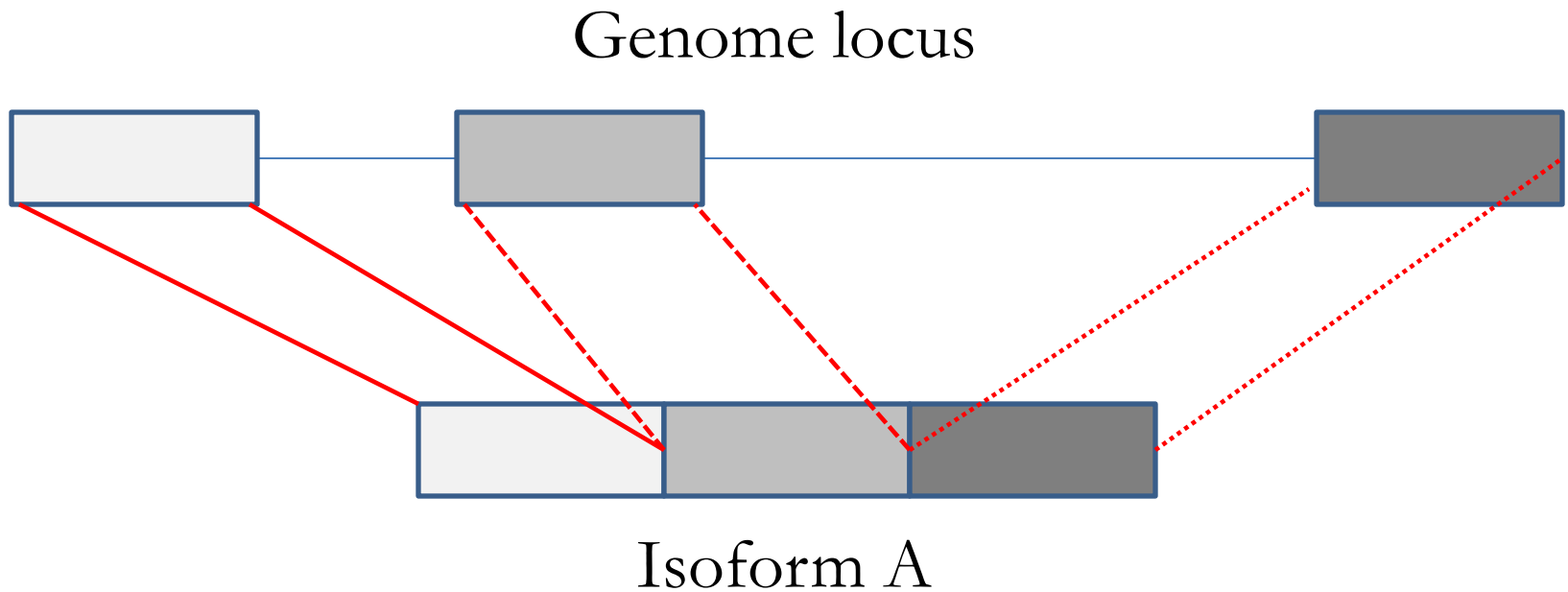# RNA-Seq does not look like the genome
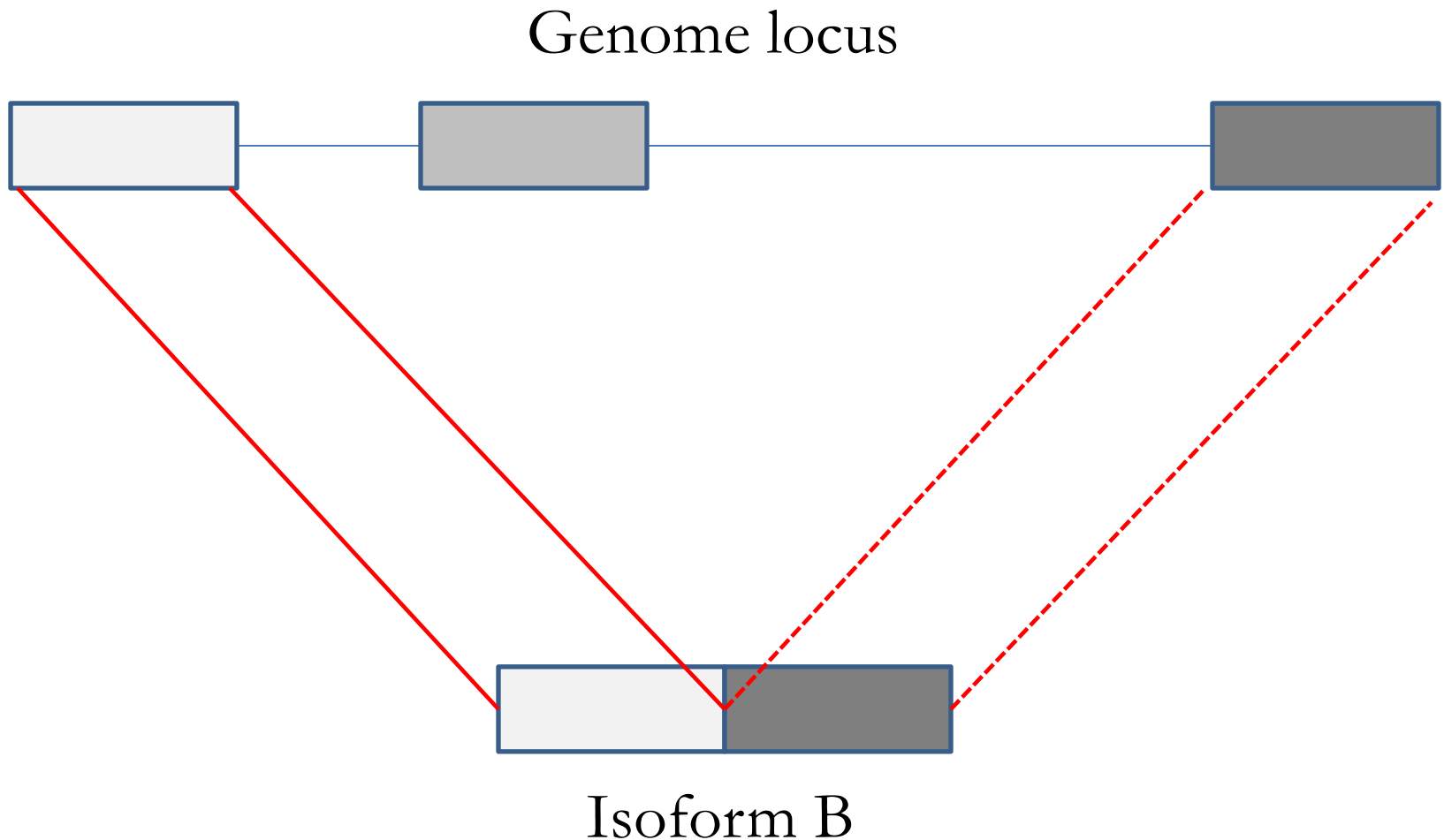
Genome locus



Isoform A



Isoform B

# RNA-Seq does not look like the genome

Genome locus



Isoform A

# RNA-Seq does not look like the genome

Genome locus



Isoform B

# RNA-Seq does not look like the genome

## Genome locus



Deletion                                    Deletion

Genomic aligners expect the library to reflect **genome** architecture. Intron splicing looks like large deletions, and can confuse aligner.

One alternative is to align to transcript FASTA rather than whole-genome.

# Splice-aware genome aligners

- Tophat2 (bowtie derived)
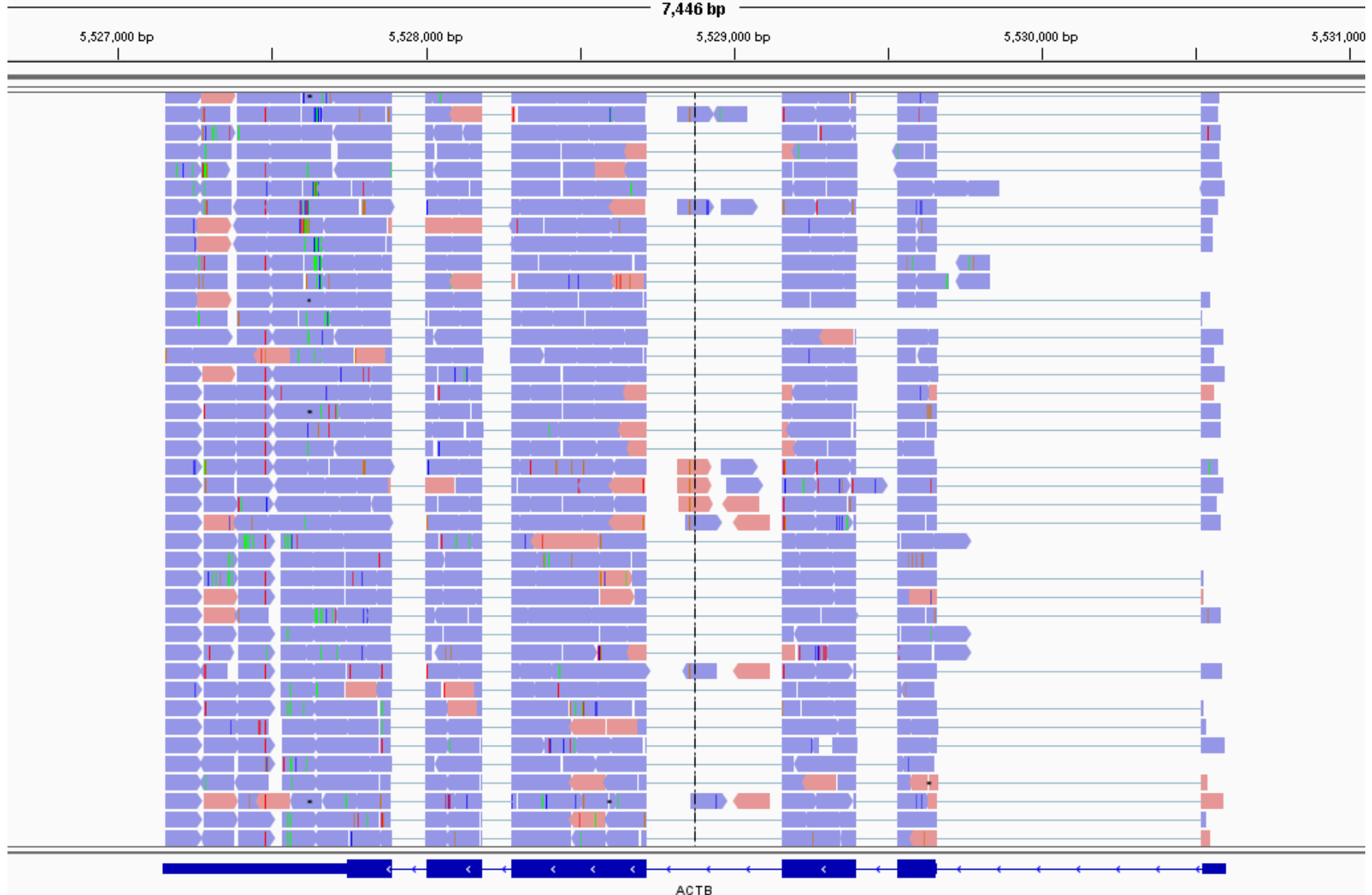
- Spliced Transcript Alignment to Reference (STAR)

# Transcript alignment

- Kallisto

- Sailfish / Salmon

# *De novo* assembly

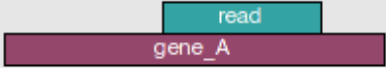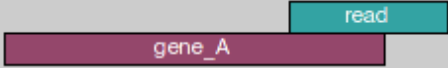- ABySS / TransABySS
- Trinity
- SOAPdenovo-Trans

# Ribosomal depletion alignment
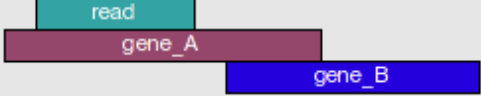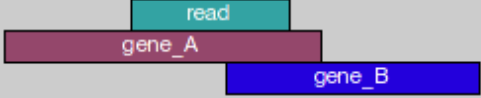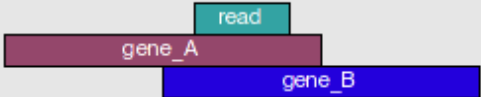
# Counting reads

**Counting tools**
HTSeq count
subread featureCount
RSEM

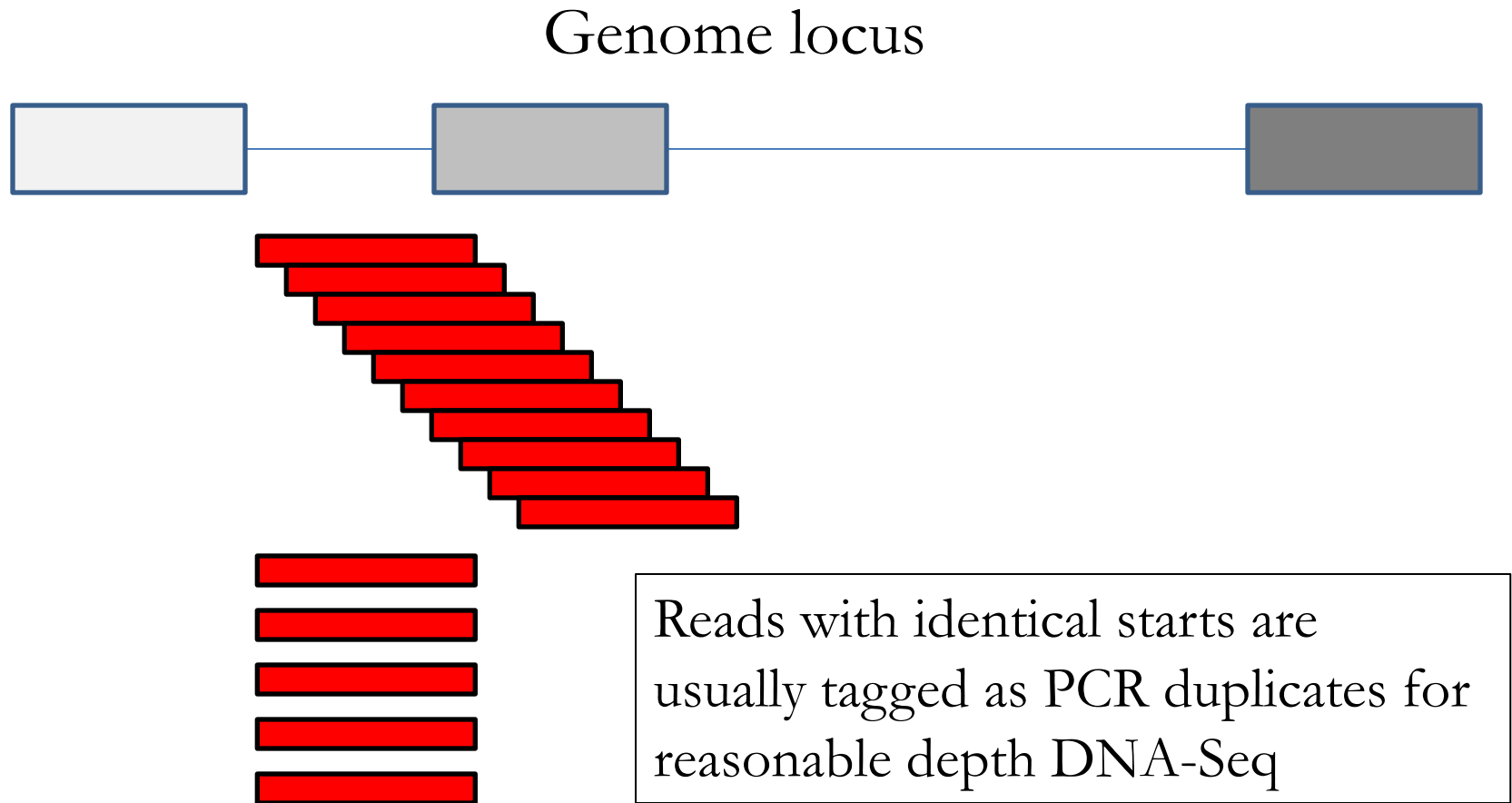| | union | intersection_strict | intersection_nonempty |
|---|---|---|---|
| | gene_A | gene_A | gene_A |
| | gene_A | no_feature | gene_A |
| | gene_A | no_feature | gene_A |
| | gene_A | gene_A | gene_A |
| | gene_A | gene_A | gene_A |
| | ambiguous | gene_A | gene_A |
| | ambiguous | ambiguous | ambiguous |

*HTSeq count

# RNA-Seq complications - duplicates

Genome locus



Highly covered genomic locus for DNA-Seq

# RNA-Seq complications - duplicates

Genome locus

Reads with identical starts are usually tagged as PCR duplicates for reasonable depth DNA-Seq

# RNA-Seq complications - duplicates

Genome locus

Best read

Reads with identical starts are usually tagged as PCR duplicates for reasonable depth DNA-Seq

# RNA-Seq complications - duplicates
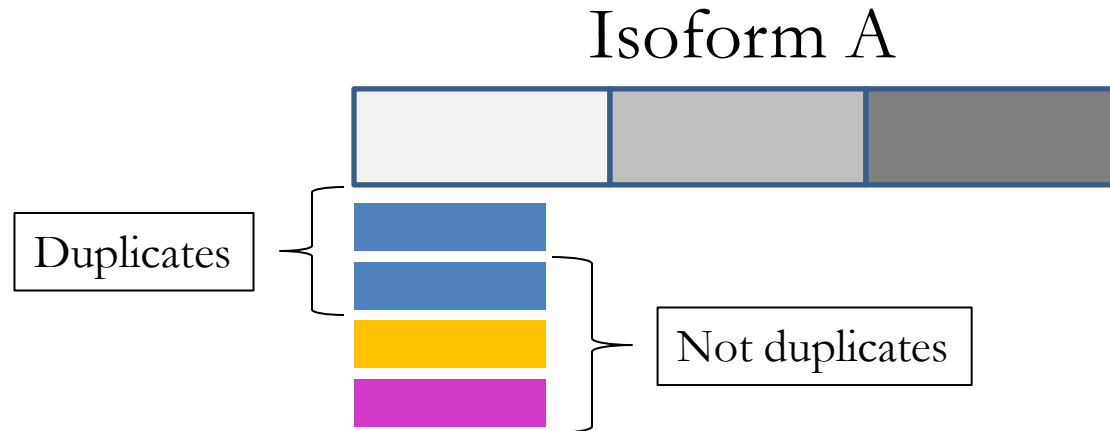


3 transcripts to library

# RNA-Seq complications - duplicates

Isoform A

• For DNA-Seq the target is equimolar, but RNA-Seq is more complicated

• Important considerations:
  • Sequencing depth
  • Gene relative expression level
  • Gene size

# RNA-Seq complications - duplicates

Isoform A



- **But** this is a solvable problem

- Adding a short barcode to each fragment during PCR, called a unique molecular identifier (UMI) we know whether reads are truly unique, even with identical 5' mapping

# Estimating gene abundance & differential expression
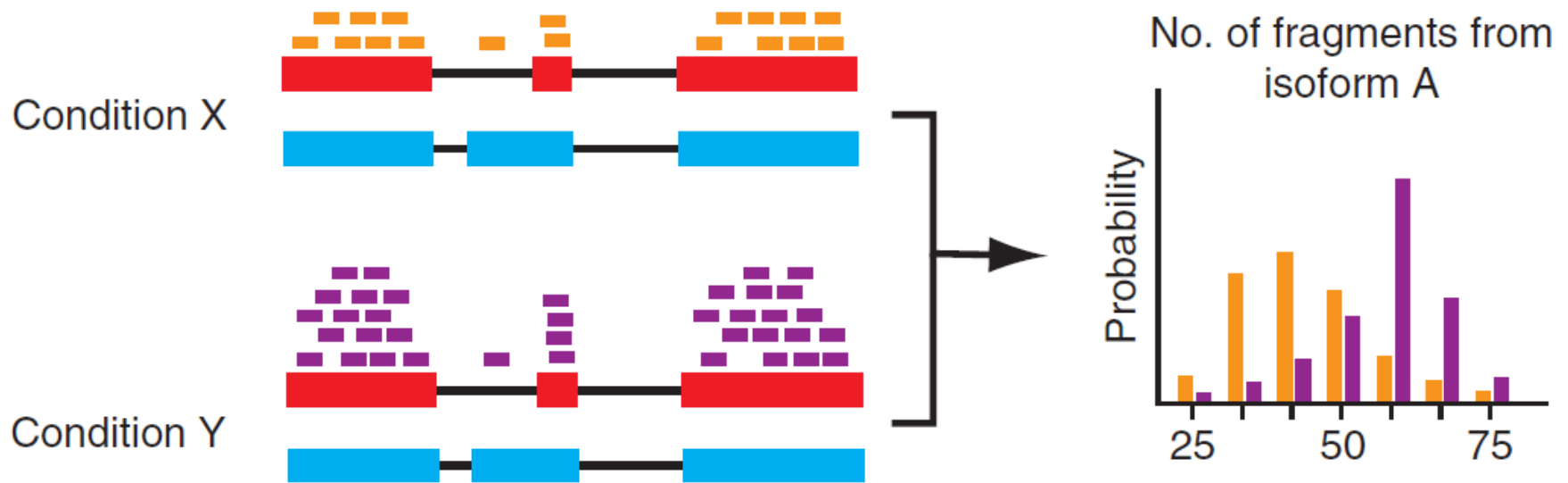
# Abundance metrics

- Fragments per Kb of exon per million mapped reads (FPKM)

- Transcripts per million (TPM)

- TPM preferable
  - Different total reads between experiments skew FPKM
  - TPM consistent, i.e. 1 TPM sample A and sample B really means similar abundance

# Modeling counts – edgeR, DESeq2

- **<u>Counts are not normally distributed</u>**

- What models counts?
  - Poisson distribution
    - But assumes variance & mean equal

  - Negative binomial
    - Mean $\neq$ variance
    - edgeR, DESeq2 R packages
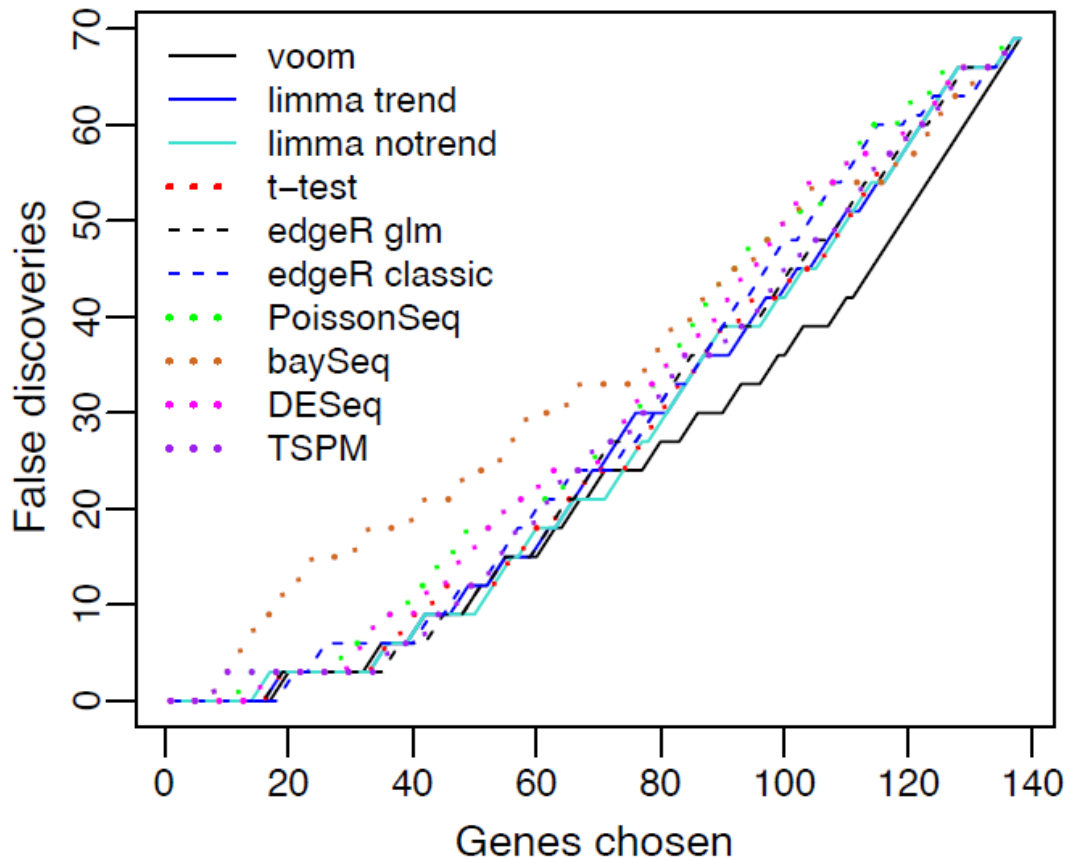
# Modeling counts – cuffdiff2



4) Combine uncertainty and overdispersion into a single model of fragment count variability (beta negative binomial)

Condition X

Condition Y

No. of fragments from isoform A

Probability

25    50    75

5) Test for signficance of changes between conditions in transcript-level counts
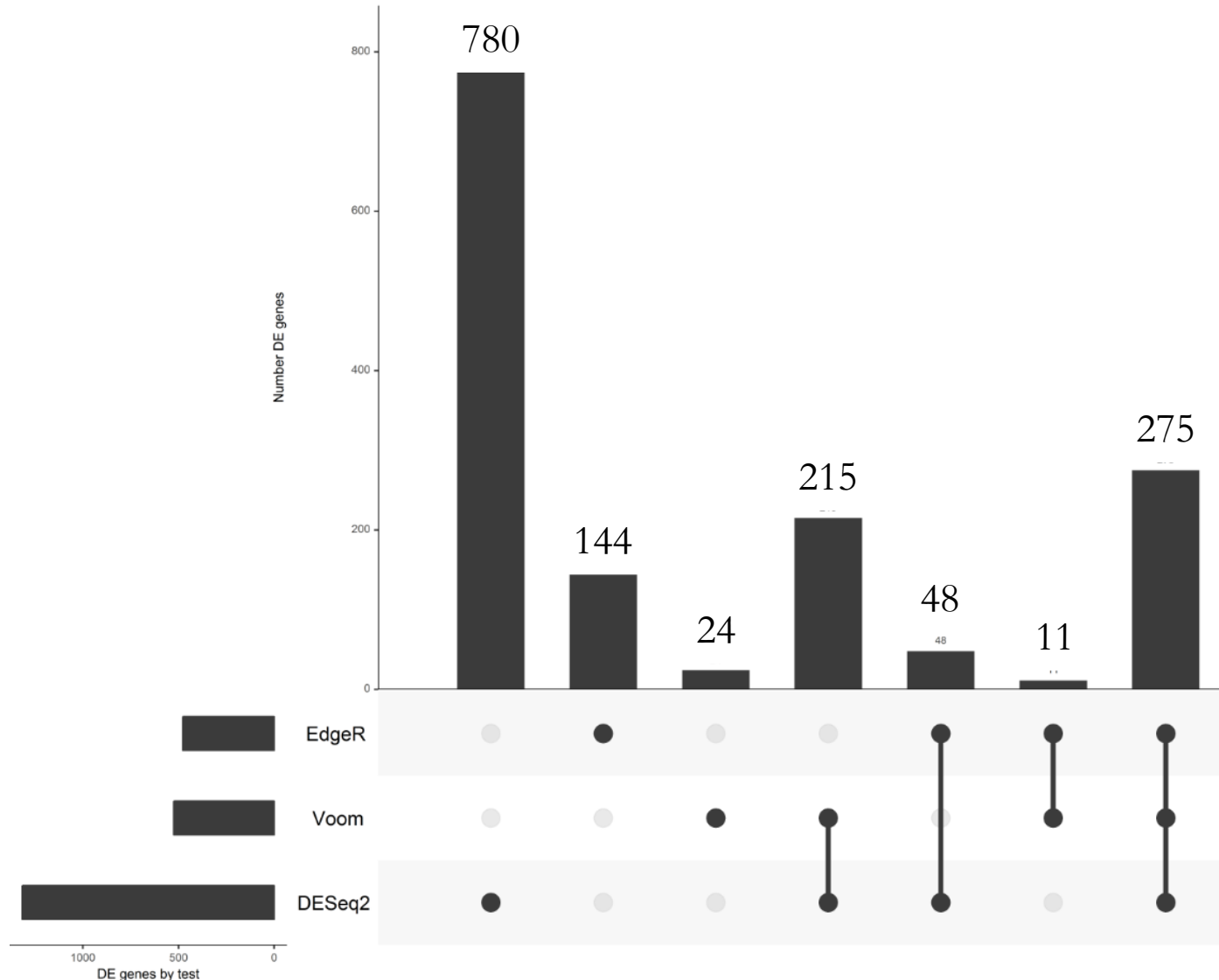
- edgeR & DESeq2 model **gene-level** differential expression
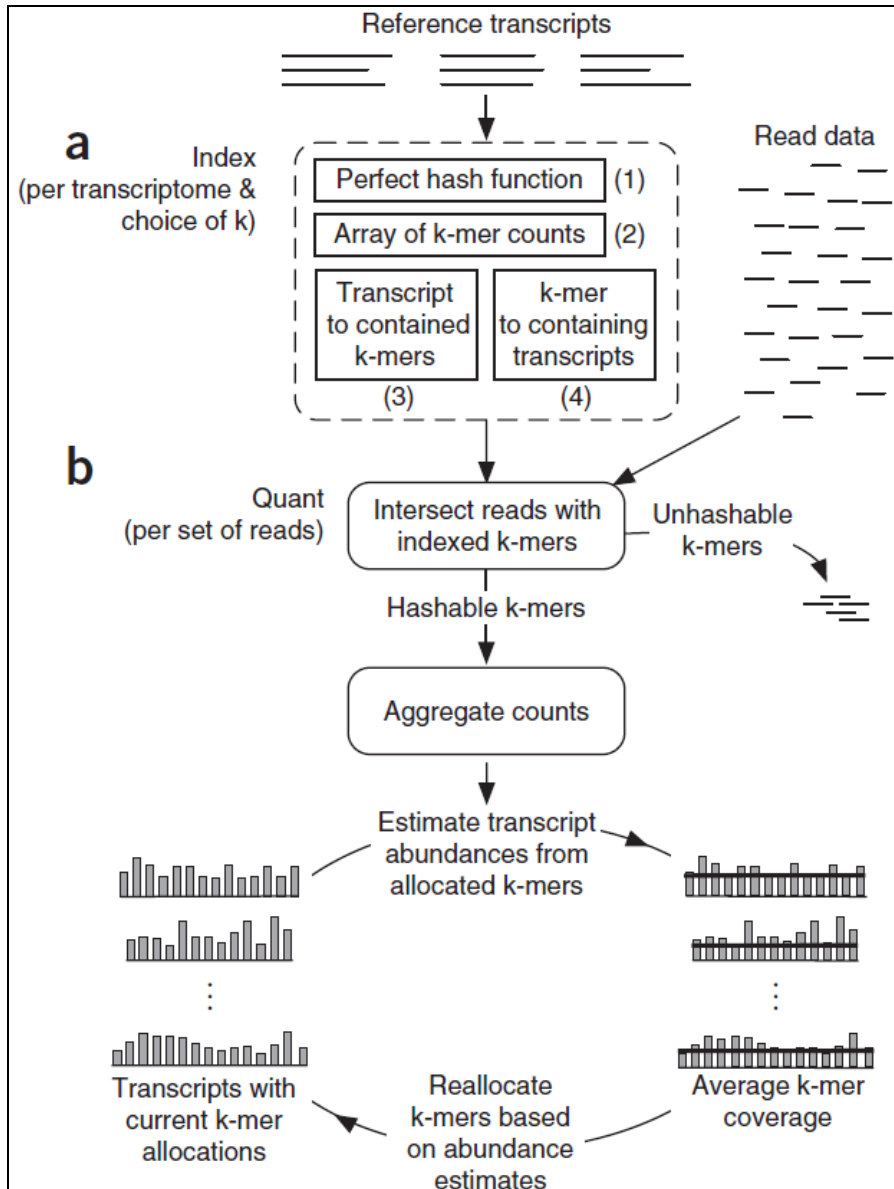- cuffdiff2 tests for significant isoform-level DE

# Modeling counts - VOOM



- VOOM uses log2 of counts per million normalization factor

- Differential expression using the empirical Bayes limma pipeline

*2014 Law *et al.* PMCID: PMC4053721

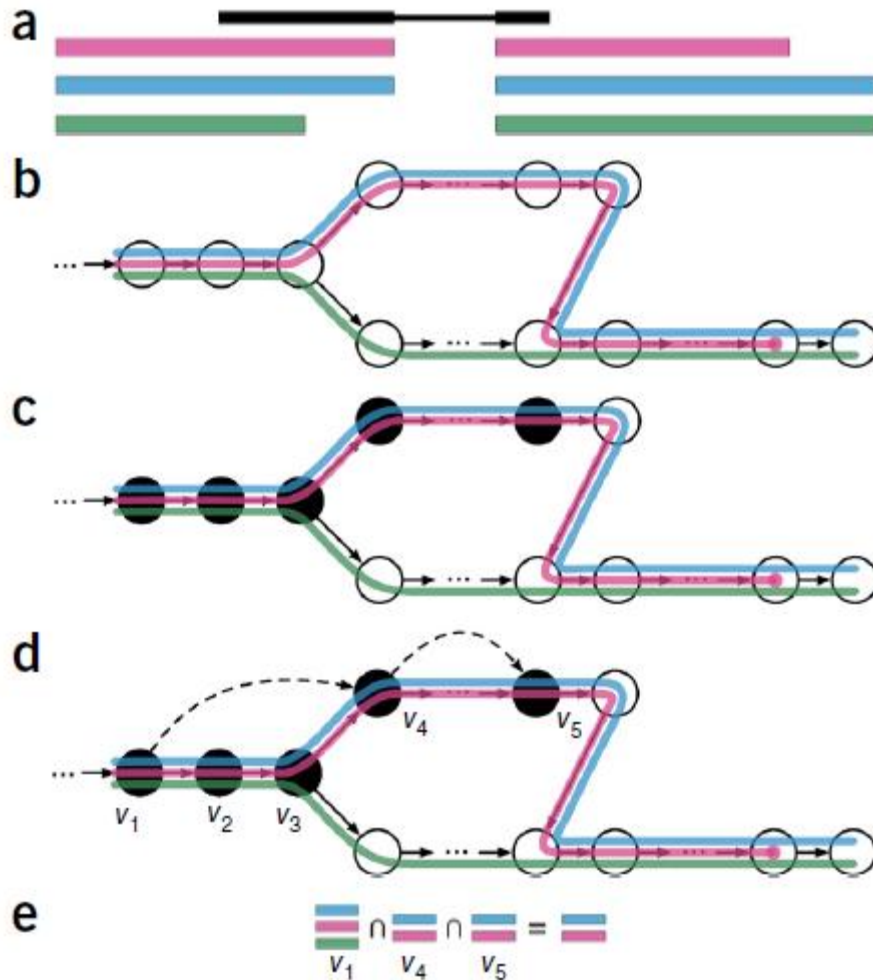# How similar are gene DE algorithms?

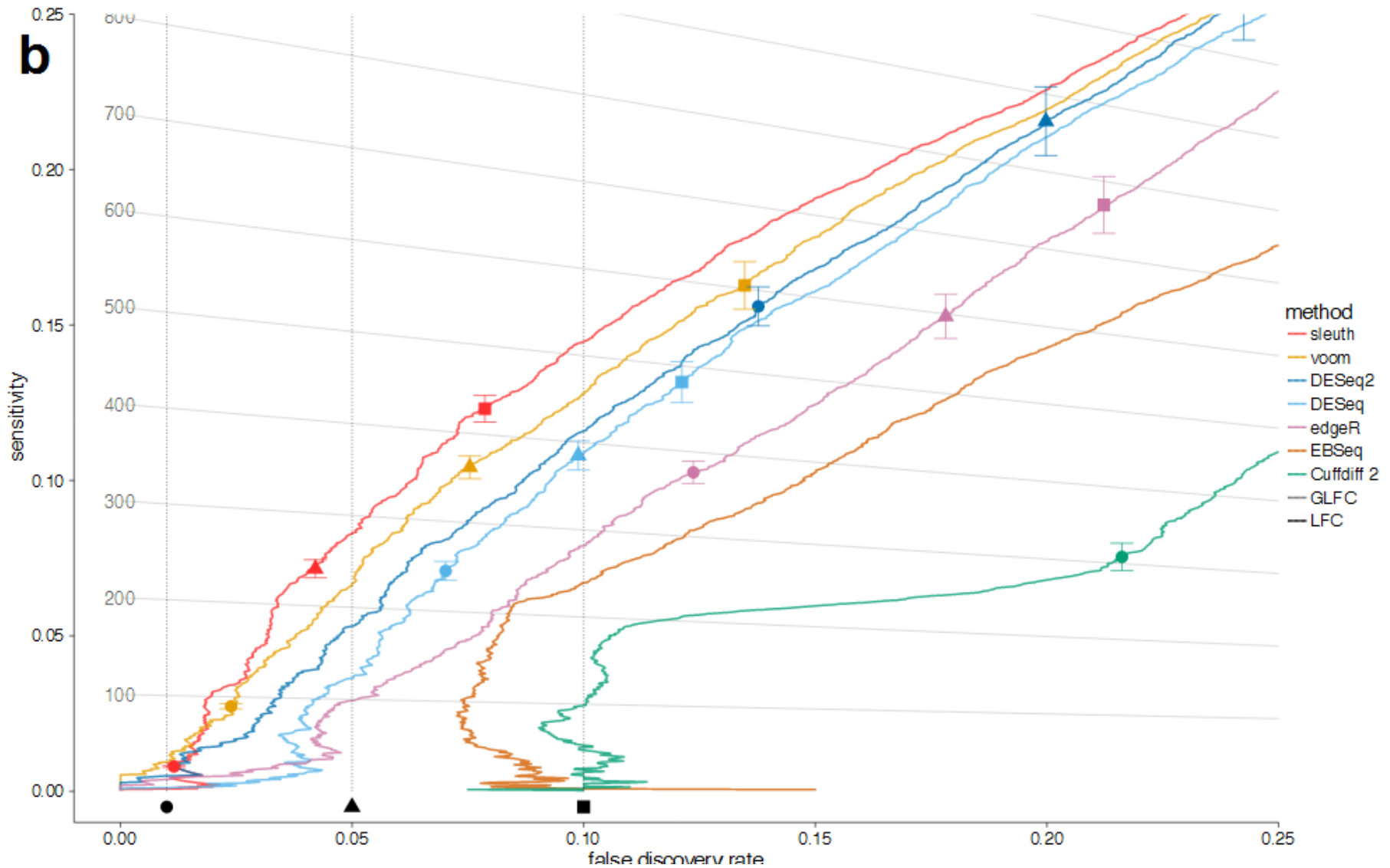# Transcript k-mer modeling - sailfish



- Hashing uses a large amount of memory

- But lookups are blazingly fast

- Calculates TPMs

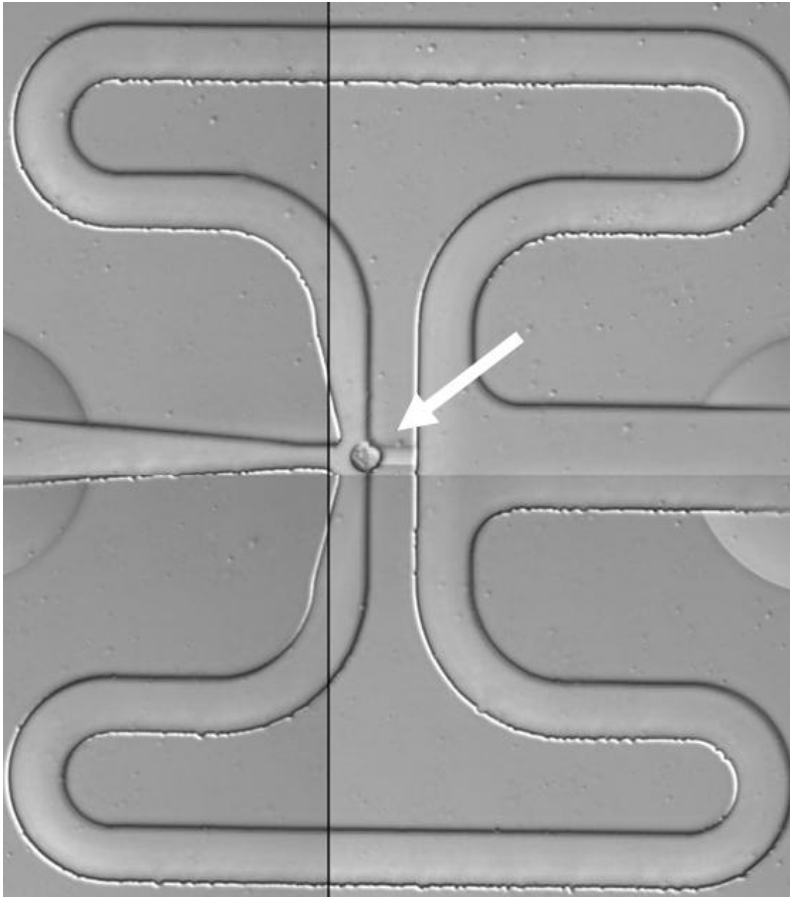# Transcript pseudoalignment - kallisto



- Builds a de Bruijn graph of transcript sequences
- Pseudoalignment – compatible transcripts, not where in transcript

- Very fast and efficient

- Allows for bootstrapping

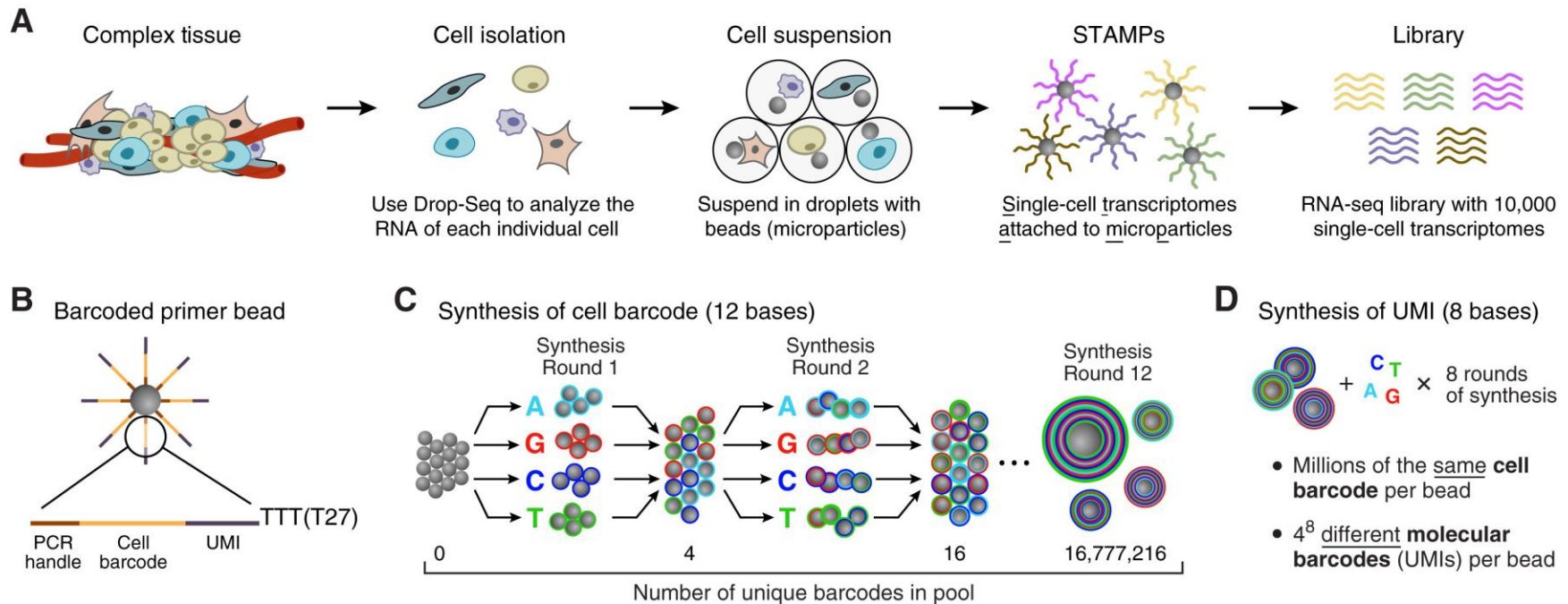# Transcript DE from kallisto - Sleuth

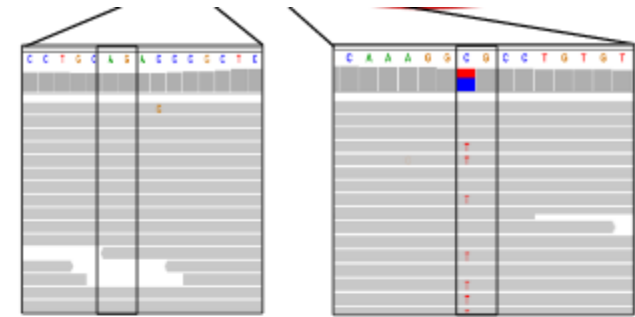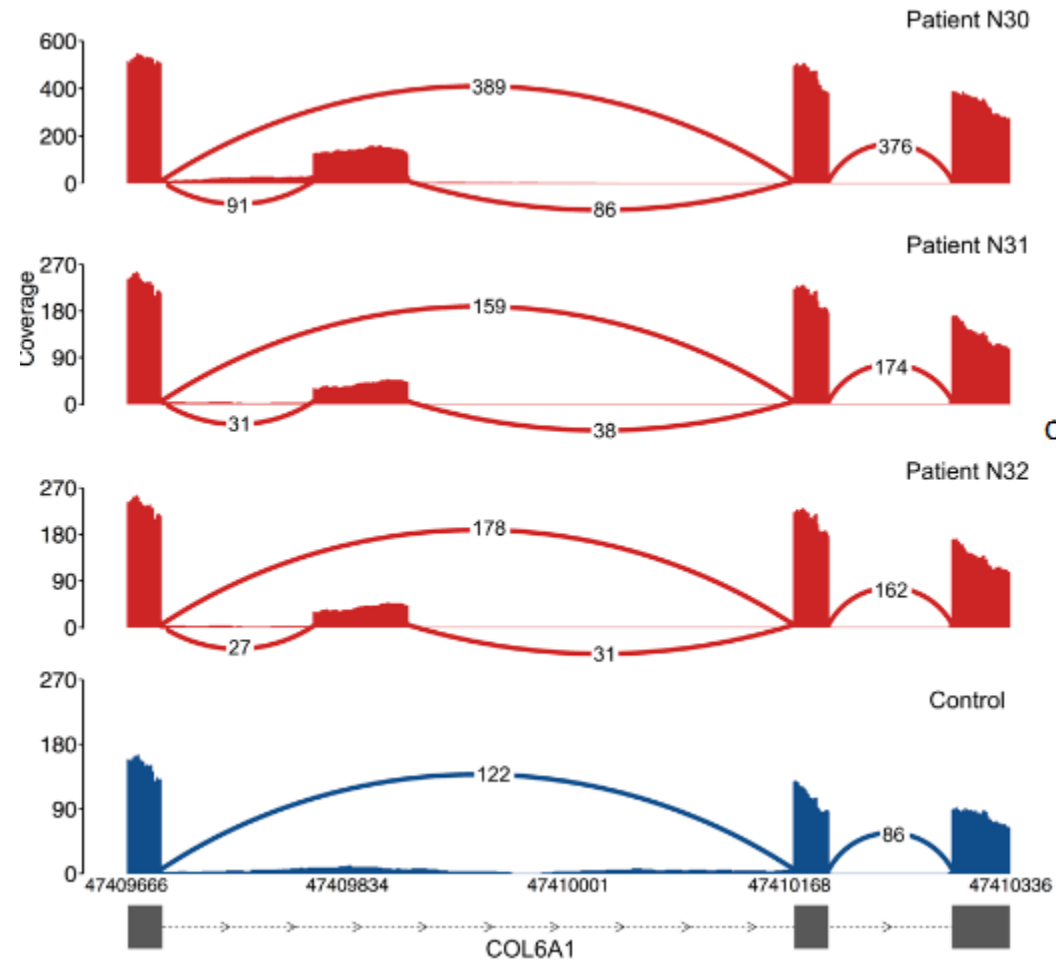# Single cell RNASeq – Fluidigm C1



- Microfluidics capture single cells

- Lyse cells and generate cDNA

- Requires live cells

- 96-800 cells / chip
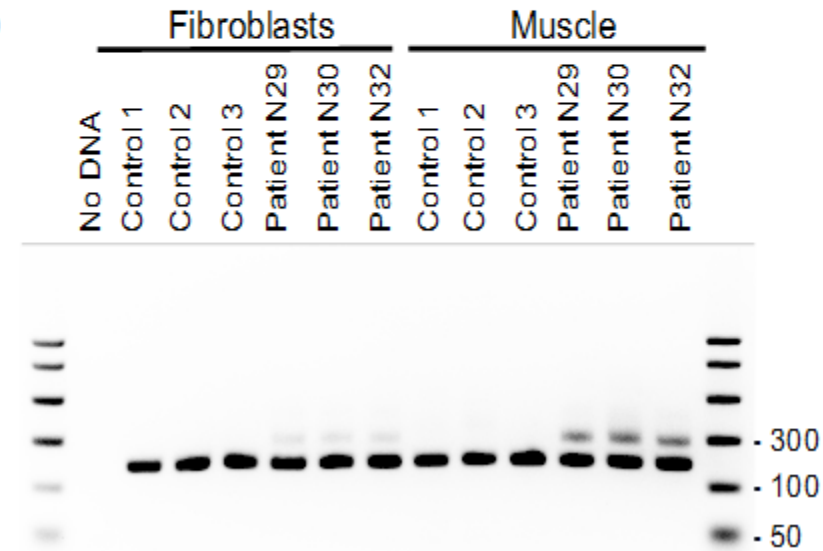
# Single cell RNASeq – DropSeq



- Flows beads in a droplet.
- Cells, usu singles, merge into a droplet for library prep.
- 10s of thousands of cells.

* Macosko *et al.* 2015. PMCID: PMC4481139

# RNA-Seq to diagnose Mendelian disease

# Summary

- RNA-Seq vs. microarrays
  - Microarray requires knowing target sequence
  - Poor dynamic range (hard to detect low-level expressors, saturate at high-levels)

- Sample collection & storage
  - Blood should be spun down to PBMCs
    - Can be directly lysed
  - Tissue should have RNAlater applied soon as possible, followed by disruption.

# Summary

- cDNA / library prep
  - Polyadenylation library methods selectively capture polyA transcripts

  - Ribosomal depletion methods degrade ribosomal RNA, but leave non-polyadenylated

  - Strandedness
    - There are an appreciable number of genes with antisense transcripts.
    - Also useful for identifying genes in species without a reference genome

# Summary

- Sequencing technologies
  - Long reads (expensive), but sequence full isoform
  - Short read. Reasonable price.

- Aligner
  - Must use an aligner that is aware of introns
  - May align to either genes (STAR, etc) or transcriptome (Tophat 2 and kallisto)

- Counting
  - Subread featureCounts, HTSeq count, RSEM

# Summary

- Differential expression
  - Gene, negative binomial: edgeR, DESeq2
  - Gene, log2 counts per million: VOOM
  - Transcript, TPM, kallisto.
  - Transcript, TPM, Sailfish

- Single-cell
  - Isolate a single cell and make a libraries (low-output)
  - Spike-ins help
  - UMIs help also