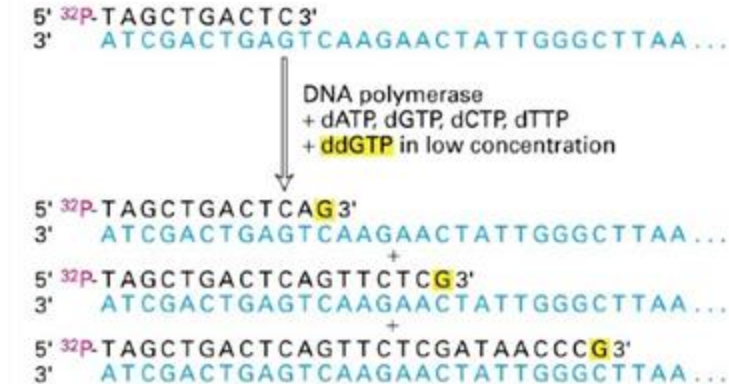# Long Read Sequencing

Chris Miller, Ph.D.
Washington University in St Louis

# How to sequence a human genome: Sanger method
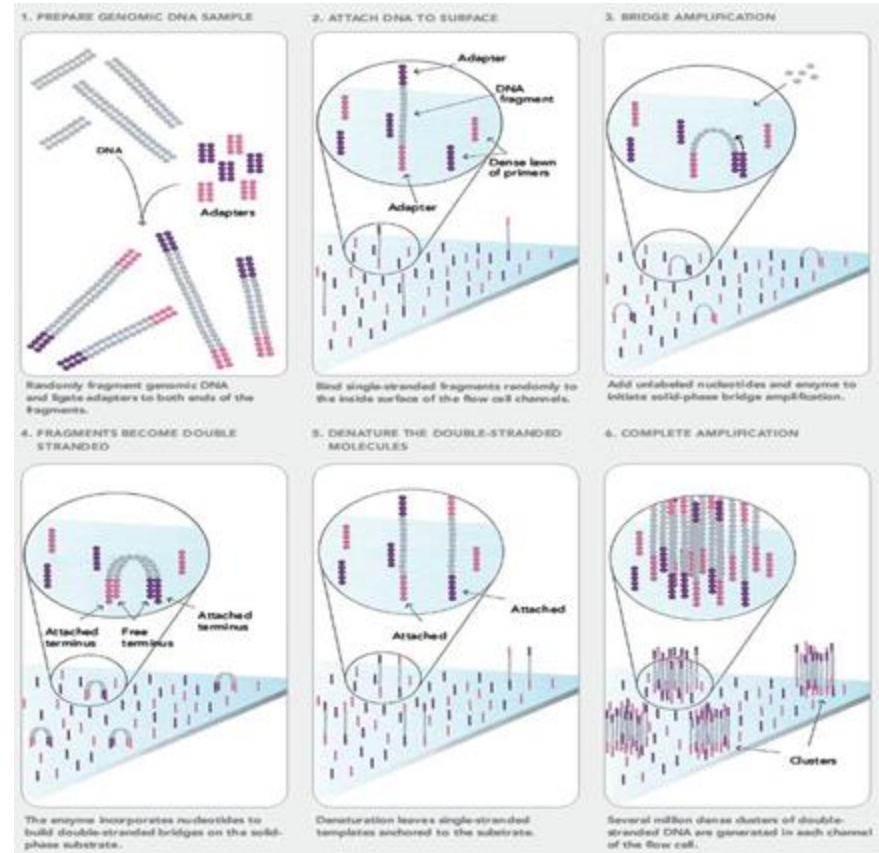
**Key points:**

1) sequencing by synthesis (not degradation)

2) primers hybridize to DNA

3) polymerase + dNTPS + ddNTP terminators at low concentration
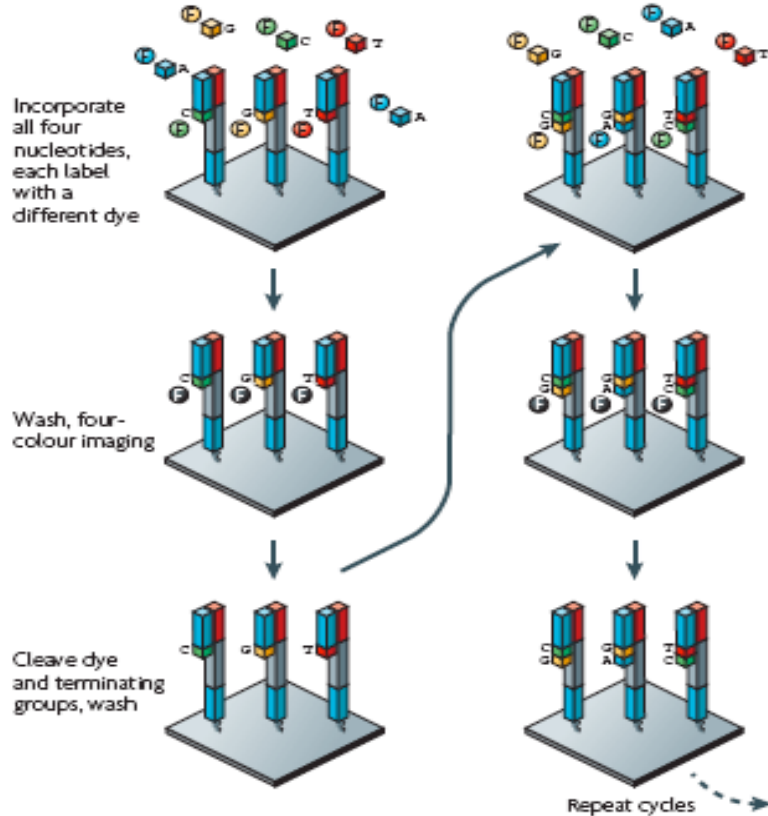
4) 1 lane per base, visually interpret ladder
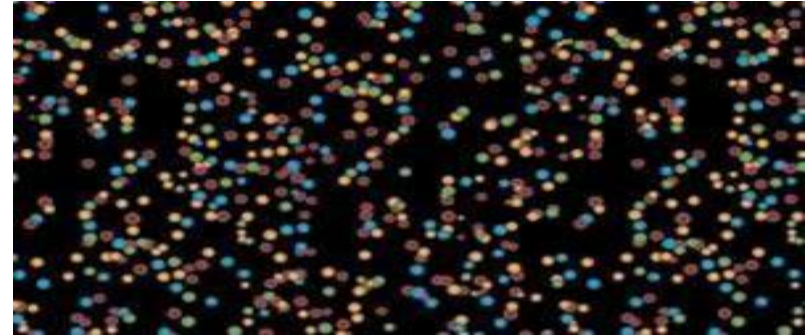
# Solexa (Illumina) sequencing (2006)

- **PCR amplify sample (opt.)**

- **Immobilize and amplify single molecules on a solid surface**

- **Reversible terminator sequencing with 4 color dye-labelled nucleotides**



*Slide from Ira Hall*

# Illumina sequencing (2005)



Incorporate all four nucleotides, each label with a different dye

Wash, four-colour imaging

Cleave dye and terminating groups, wash

Repeat cycles

**4 different images merged**

**6 cycles w/ base-calling**
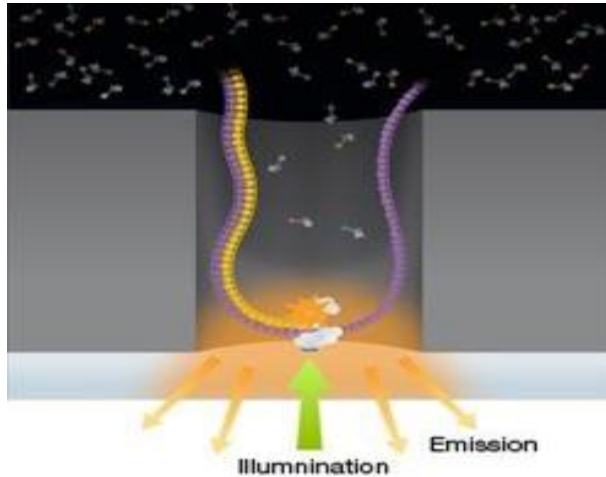
C 🟢  A 🔵
T 🔴  G 🟡
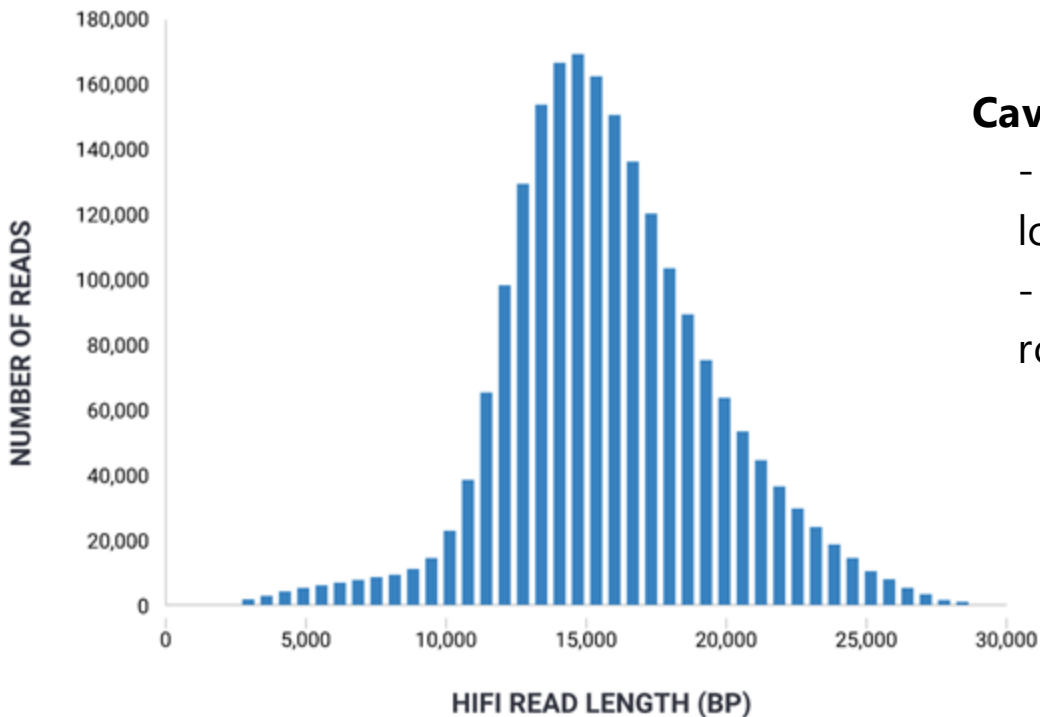
Top: CATCGT
Bottom: CCCCCC

*Slide from Ira Hall*

# Pacific Biosciences



**Key Points:**
- 1 DNA molecule and 1 polymerase in each well (zero-mode waveguide)
- 4 colors flash in real time as polymerase acts
- Methylated cytosine has distinct pattern
- No *theoretical* limit to DNA fragment length



*Slide from Ira Hall*

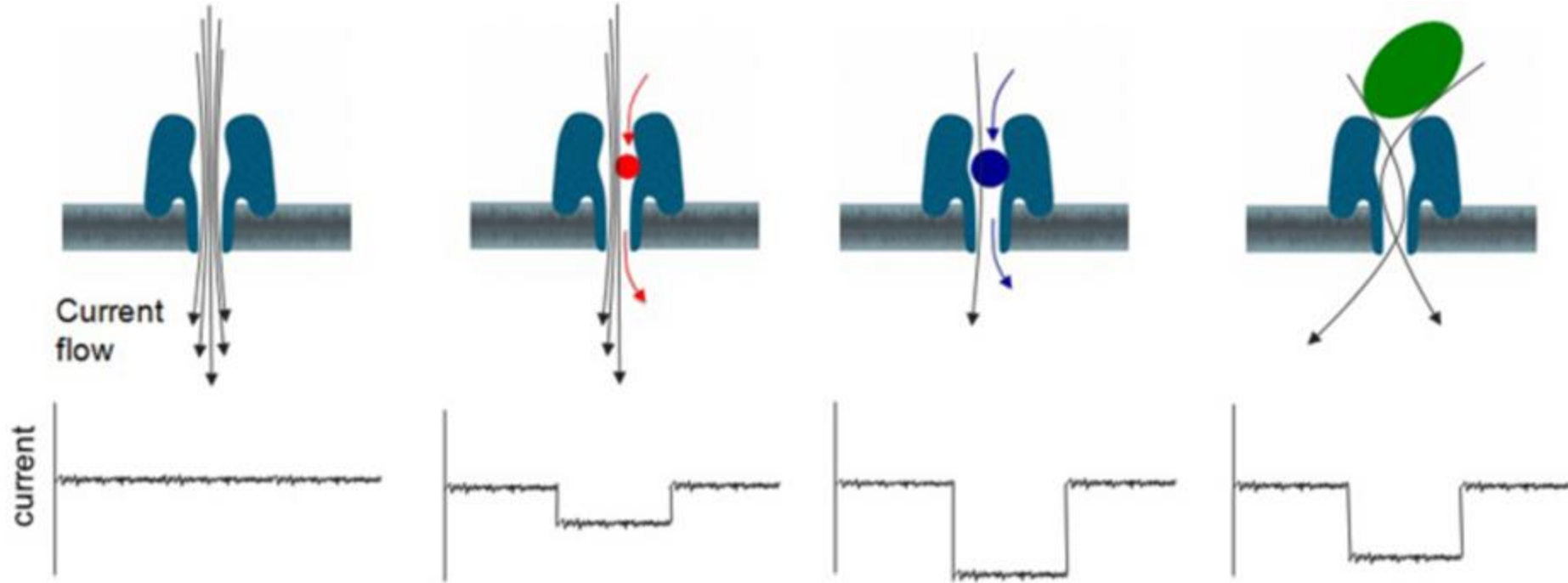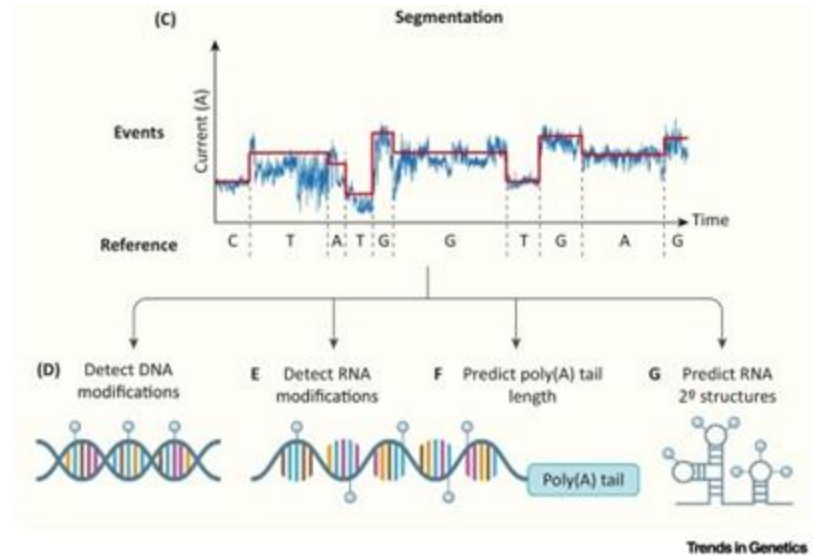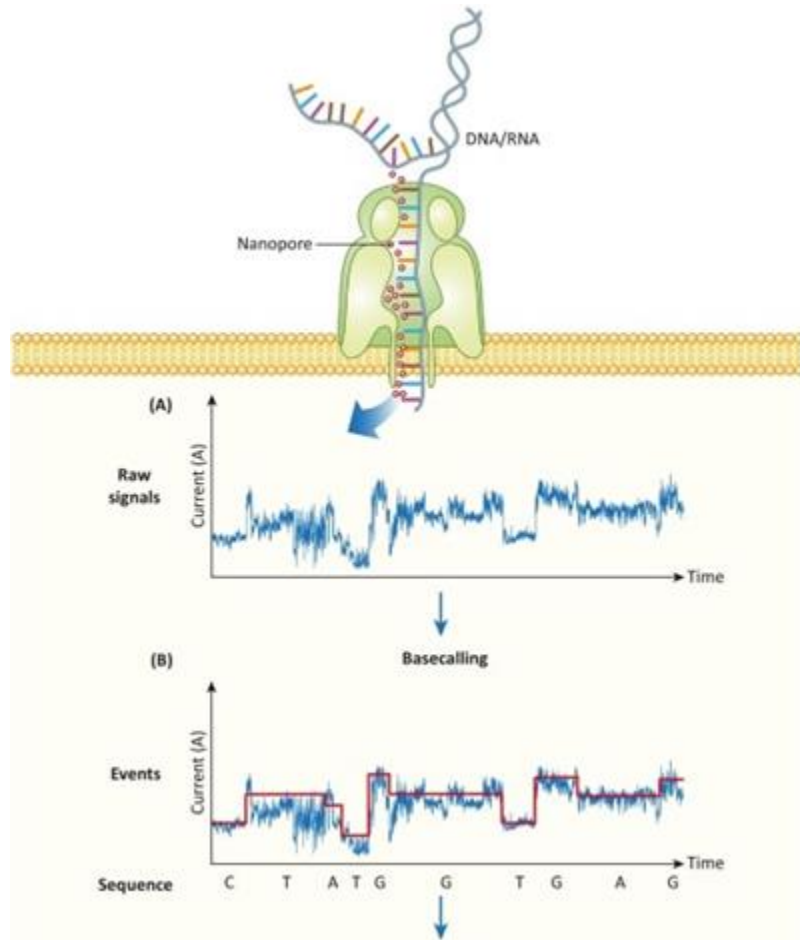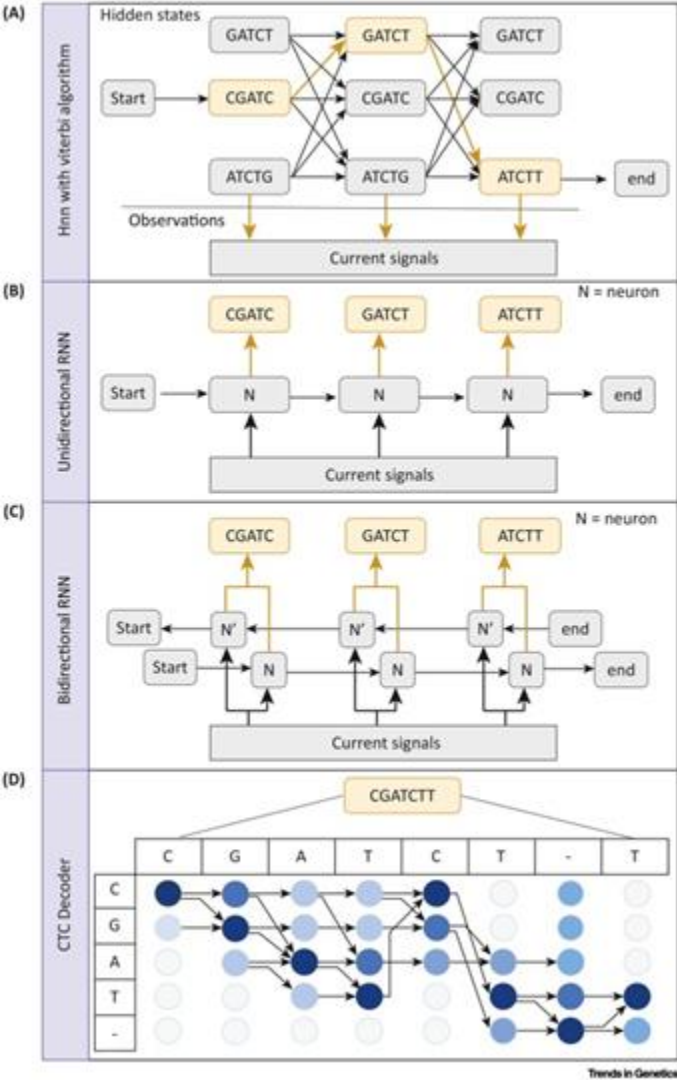# Pacific Biosciences: long reads. Great for genome assembly



**Caveats:**
- higher error rate (1-2%), lower with Duplex runs
- lower throughput : roughly 90 gigabases per run

About $4,000 for a 30x human genome on the RevIO machine
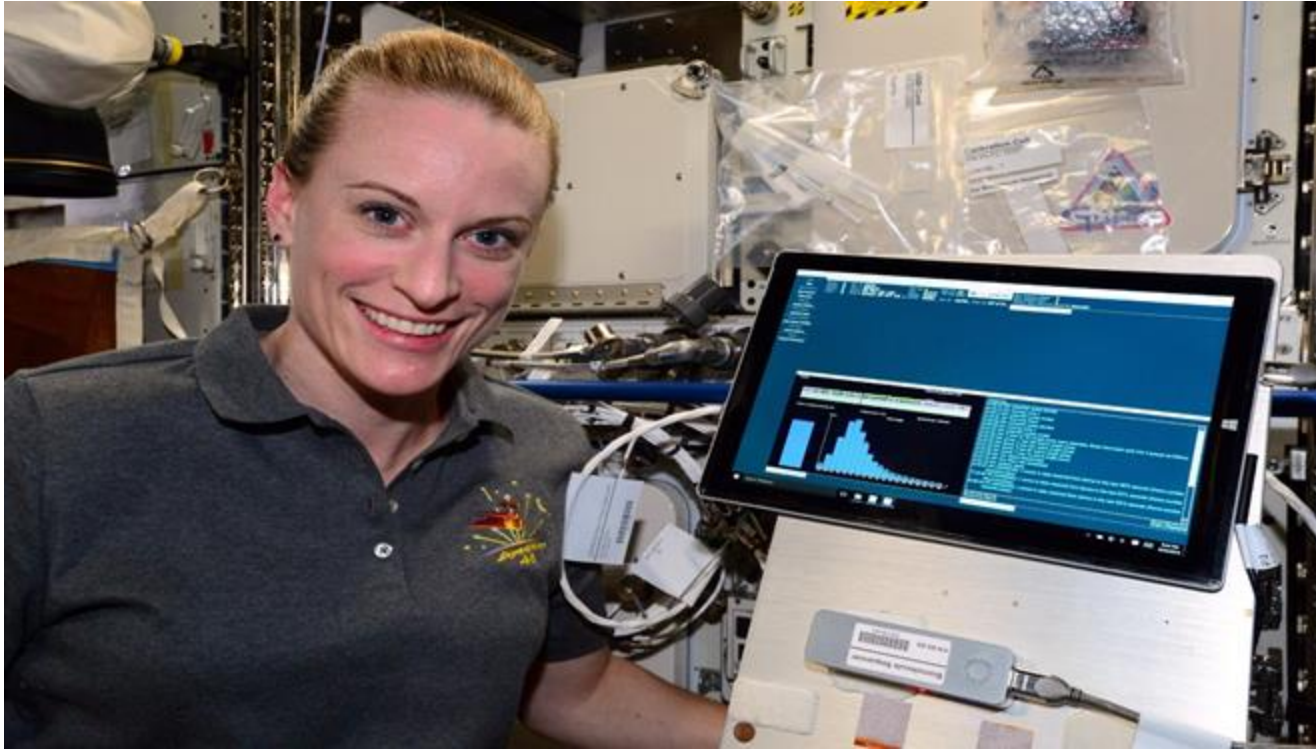
# Oxford Nanopore Technologies

Neural networks to translate signal into base calls

- Guppy (many versions)

- Dorado (v0.4, eventual guppy replacement)

- many others

Practically, that means that we can't yet throw away our raw signal intensities. (1 Tb or more per run)

# Nanopore sequencing is *extremely* portable



Kate Rubins sequencing DNA on the ISS

# ONT sequence length distribution



Recent run of a tumor sample
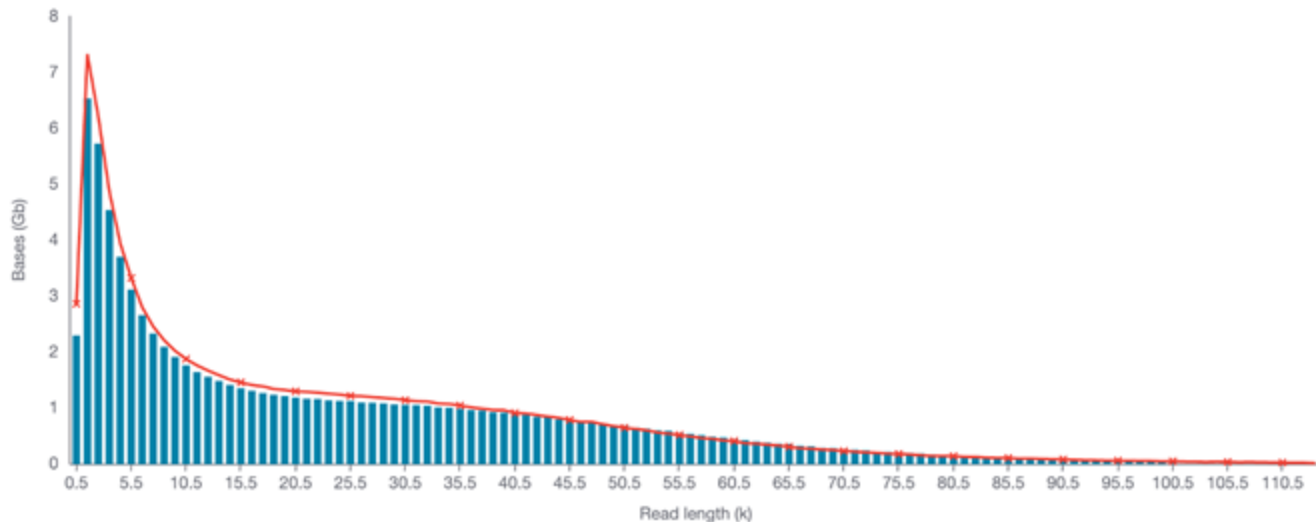
About $3,500 for a
30x human genome
on a PromethION

# What does the data look like?



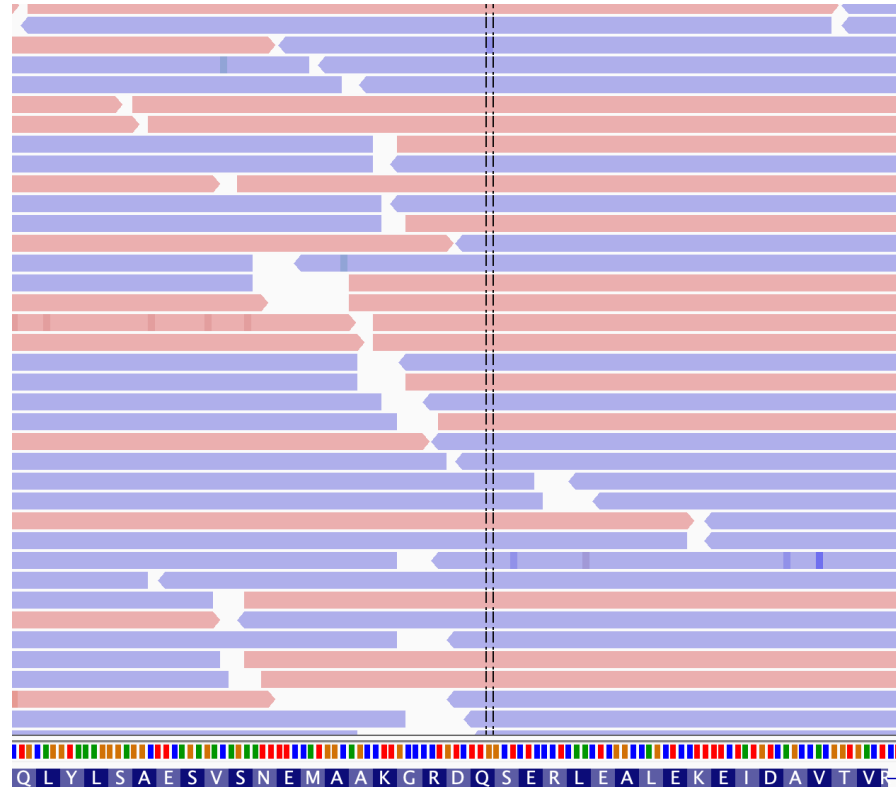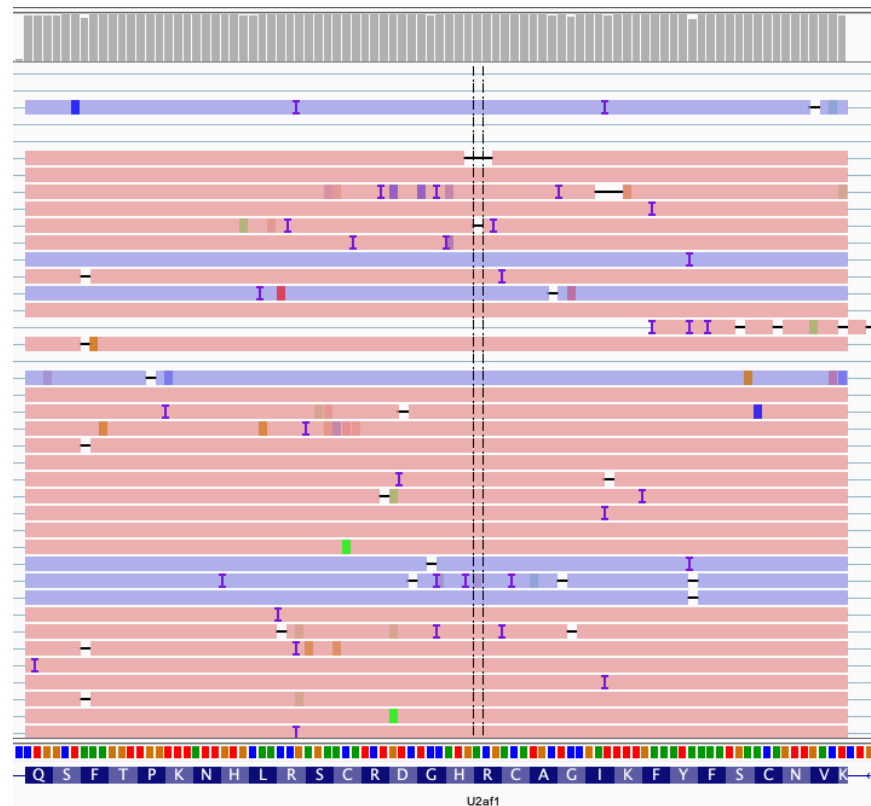Long-read ONT    ~5% base error rate

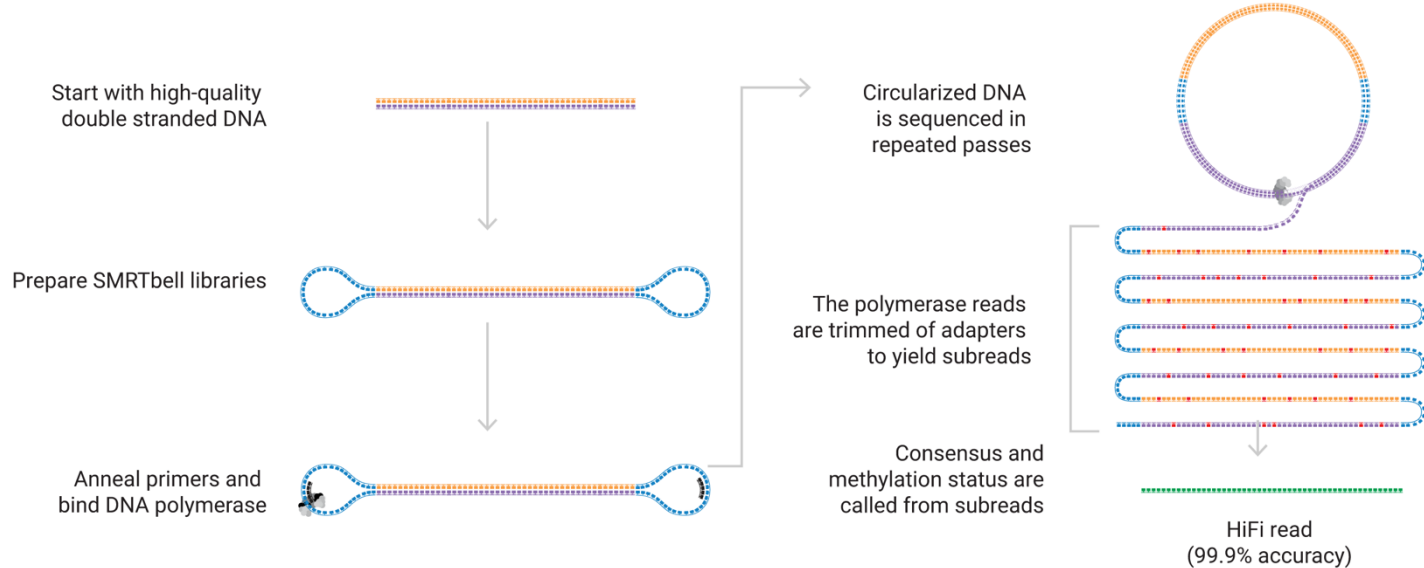Short-read Illumina    ~0.3% base error rate

# Error rates are contentious and confusing

- How do you calculate error?
  Per base?
  Per read?
  Per variant call?
  (after collapsing all of the data?)

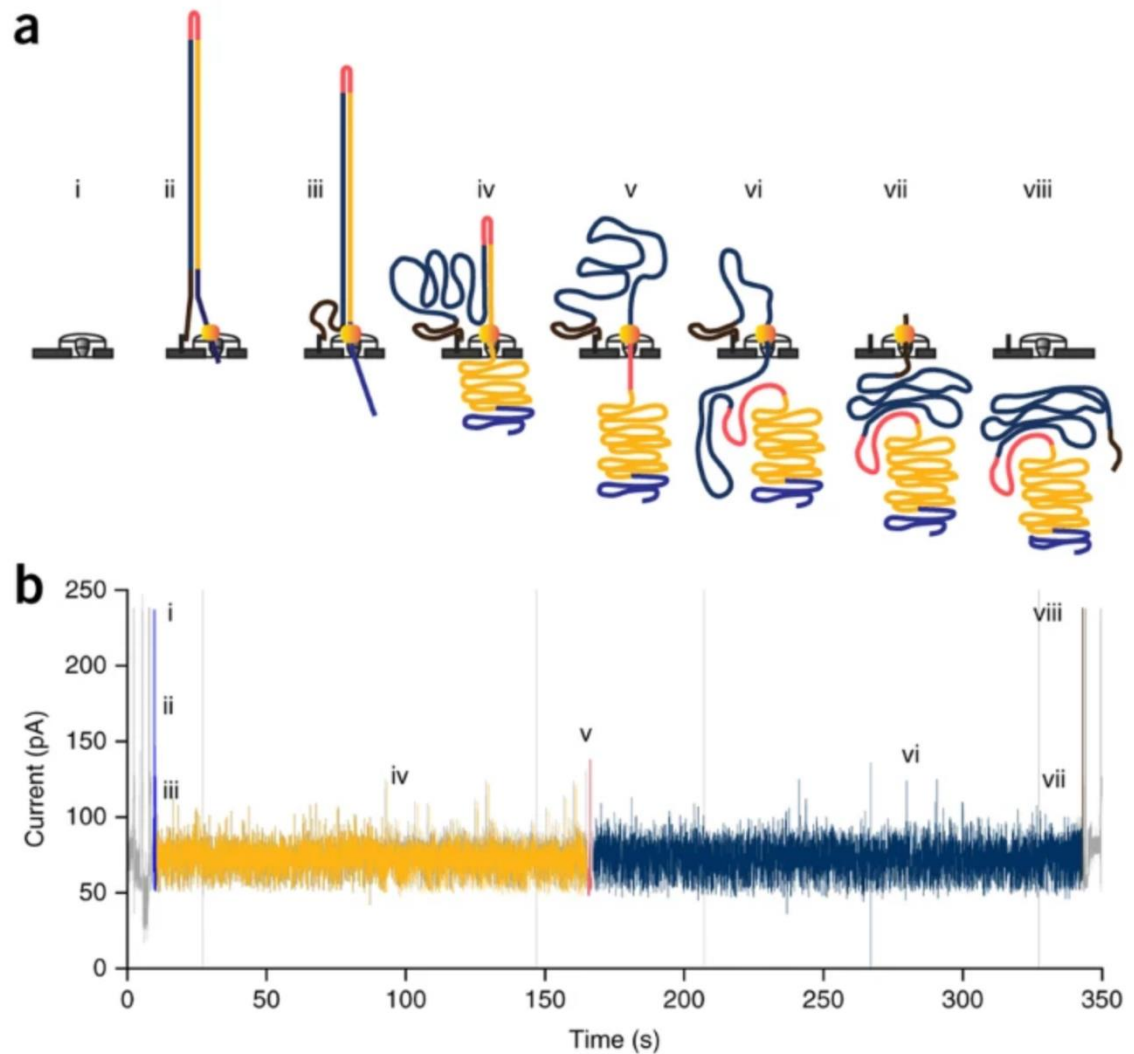# PacBio HiFi Sequencing



**How are HiFi reads generated?**

Start with high-quality double stranded DNA

Prepare SMRTbell libraries

Anneal primers and bind DNA polymerase

Circularized DNA is sequenced in repeated passes

The polymerase reads are trimmed of adapters to yield subreads

Consensus and methylation status are called from subreads

HiFi read
(99.9% accuracy)

Improved error rates

higher cost/lower throughput

https://www.pacb.com/technology/hifi-sequencing/

# ONT Duplex sequencing



Improved error rates
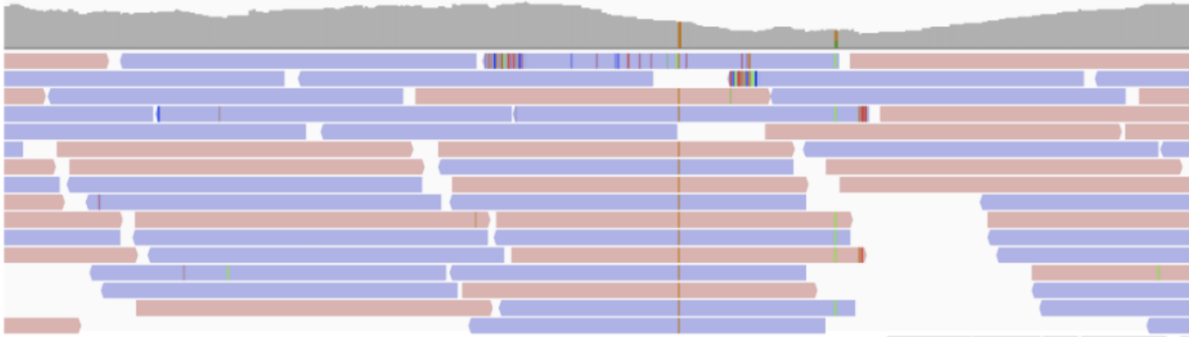
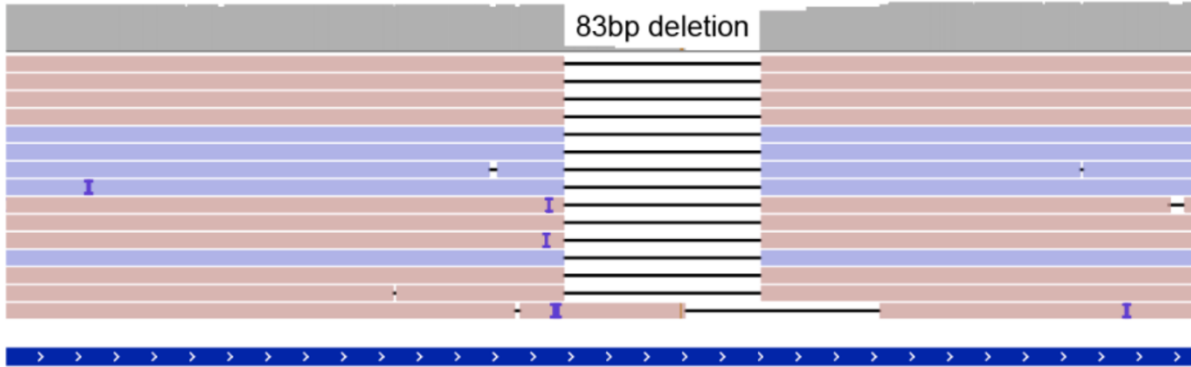higher cost/lower throughput

# Genomic DNA advantages



**Figure 1:** Blue-labeled genomic regions are accessible to long reads but not short, and have functional annotations (e.g. genes or enhancers)

David Spencer

# Large Indel detection

No indel detectable - Short-read sequencing - Illumina

83bp deletion - Long-read sequencing - Oxford Nanopore

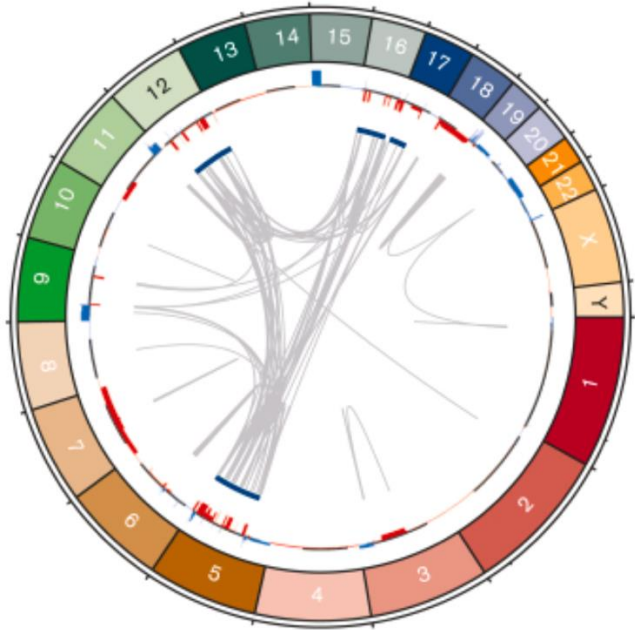83bp deletion

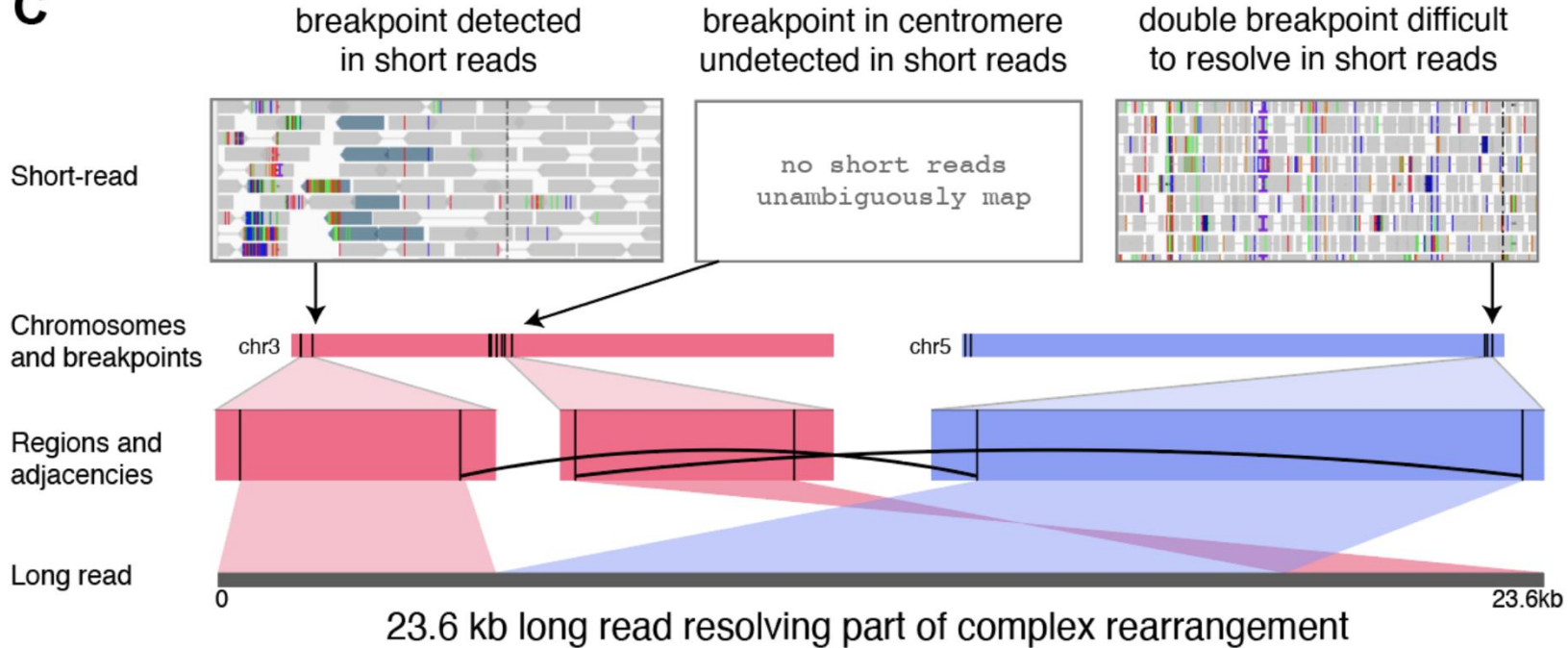chr3: 31990200-31990700          ZNF860          (protein-coding sequence)

Haley Abel

# Structural variant resolution

TP53-mutated AML

# Structural variant resolution



breakpoint detected in short reads

breakpoint in centromere undetected in short reads

double breakpoint difficult to resolve in short reads

Short-read

no short reads unambiguously map

Chromosomes and breakpoints

chr3

chr5

Regions and adjacencies

Long read

0

23.6kb

23.6 kb long read resolving part of complex rearrangement

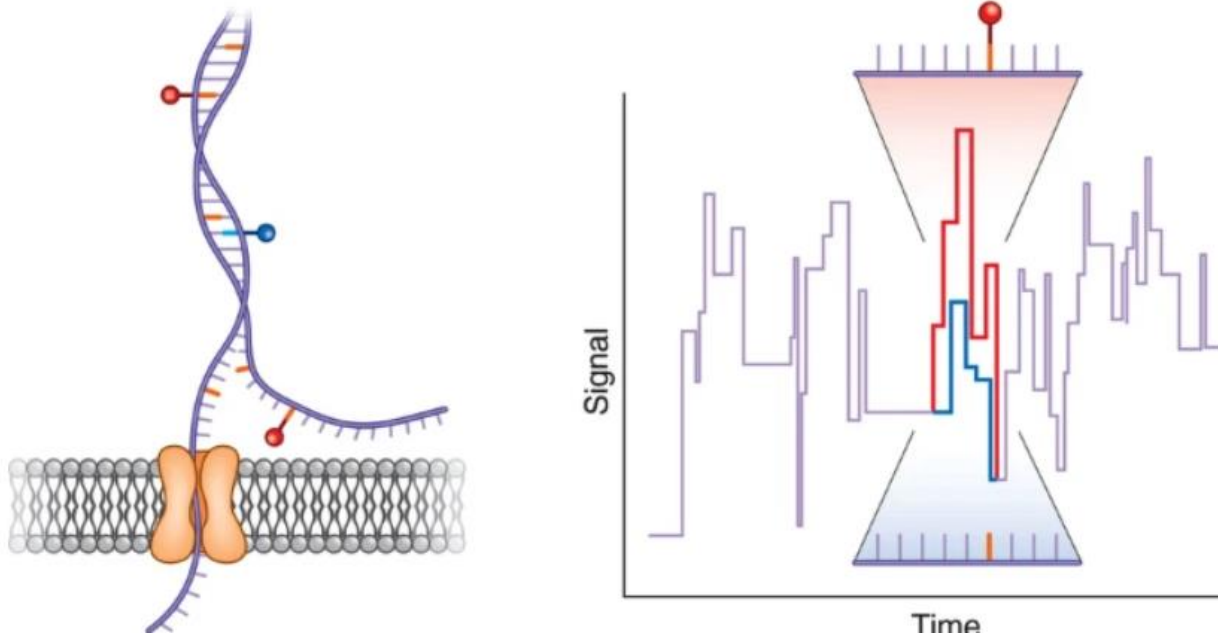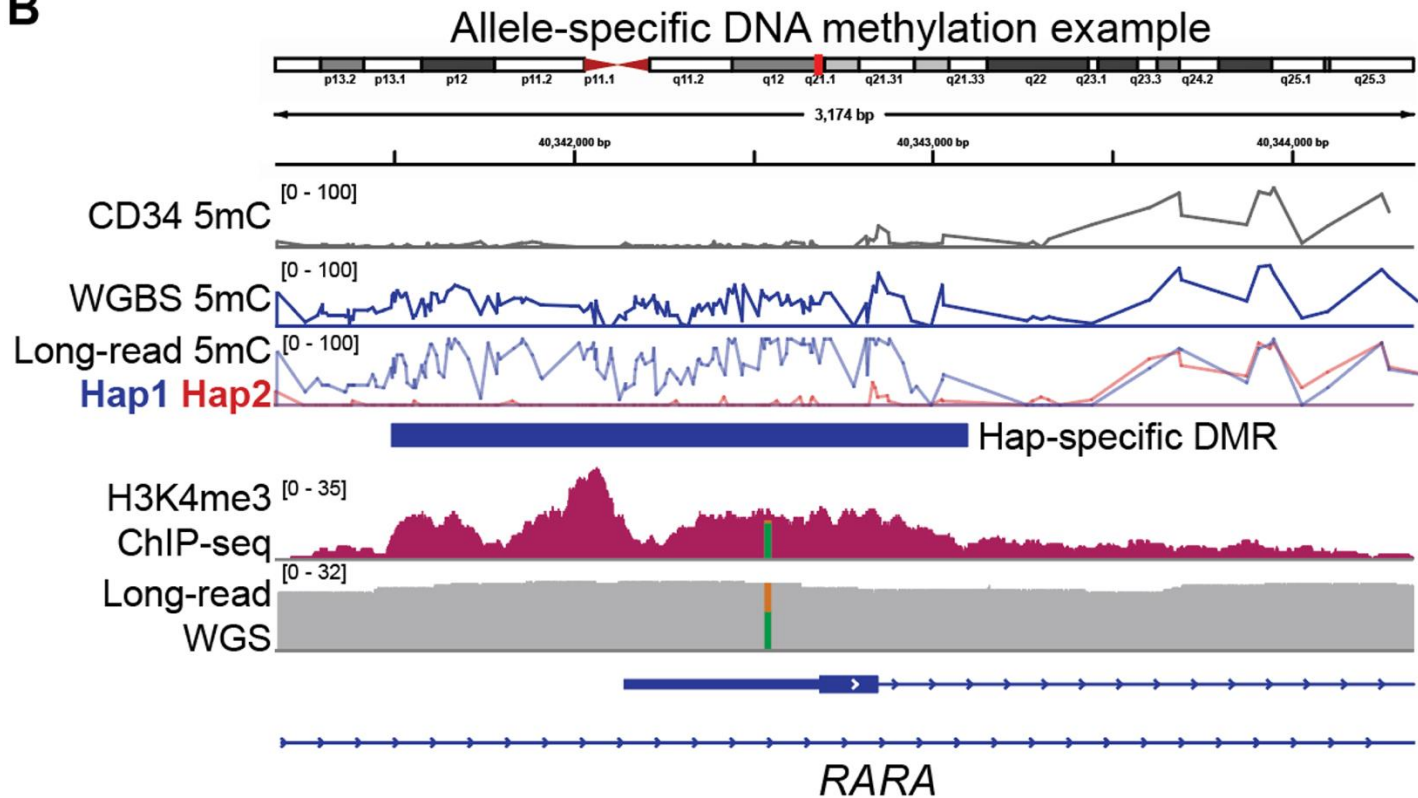Haley Abel

# Base modification detection



Can be used for 5mC as well as m6A in direct RNAseq

# Phasing of reads/modifications



**B**

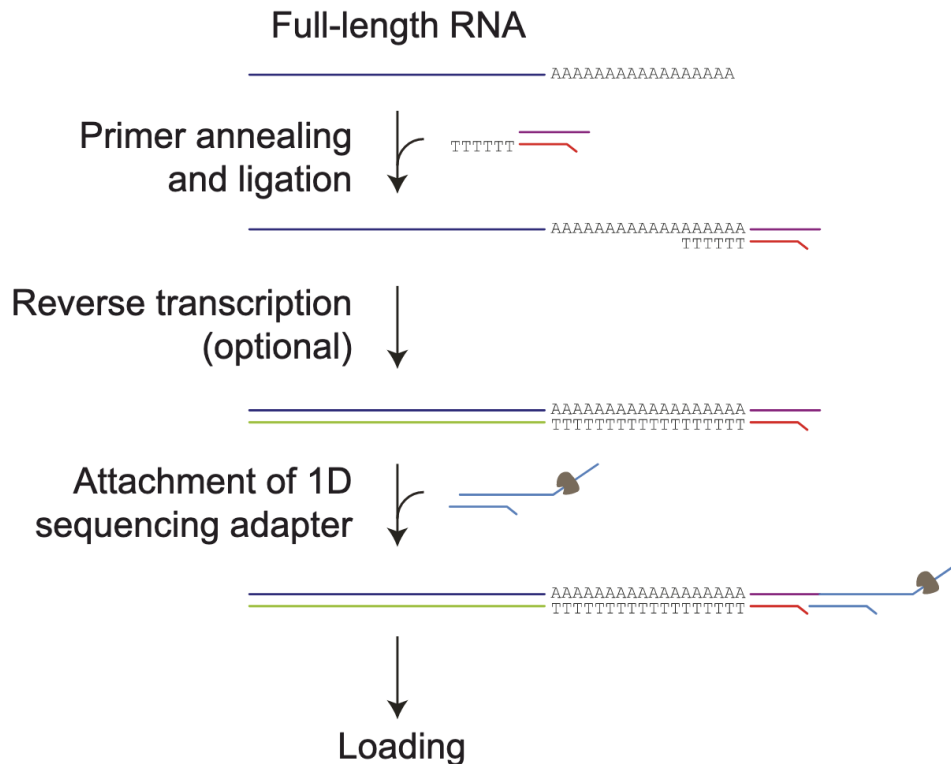Allele-specific DNA methylation example

Dave Spencer

# Genome assembly

- Assembly of personal genomes



CBFB-MYH11 fusion

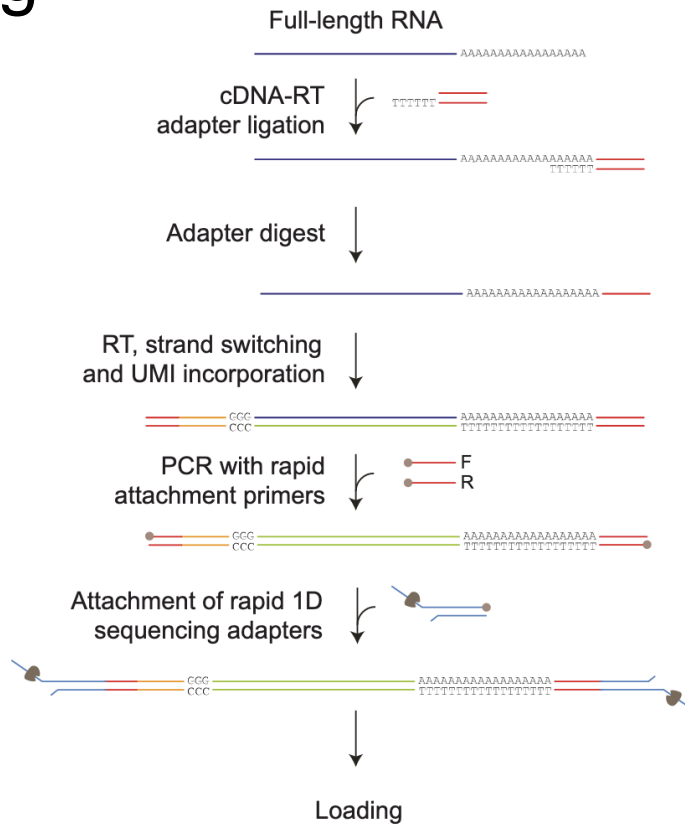| sample | hap | metric | value |
|---|---|---|---|
| Presentation | hap1 | Number_of_contigs | 402 |
| Presentation | hap1 | N50 | 88.3 Mbp |

Dave Spencer, Haley Abel

# Long-read RNA sequencing

- Direct RNA

- No amplification, less bias

- Preserves base modifications (m6a, etc)

Full-length RNA

AAAAAAAAAAAAAAAAA

Primer annealing and ligation

TTTTTT

AAAAAAAAAAAAAAAAAA
TTTTTT

Reverse transcription (optional)

AAAAAAAAAAAAAAAAAA
TTTTTTTTTTTTTTTTTT

Attachment of 1D sequencing adapter
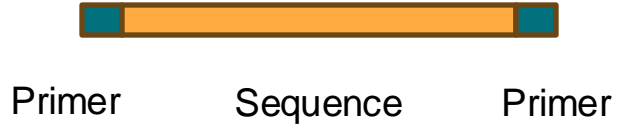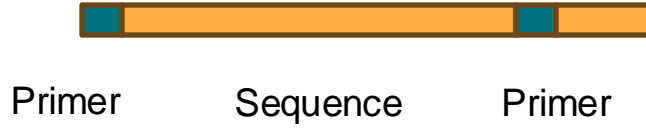
AAAAAAAAAAAAAAAAAA
TTTTTTTTTTTTTTTTTT

Loading

# Long-read RNA sequencing

- cDNA sequencing

- much higher yields

# Pychopper



Primer          Sequence          Primer

# Pychopper



Primer       Sequence       Primer

# Pychopper



Primer      Sequence      Primer
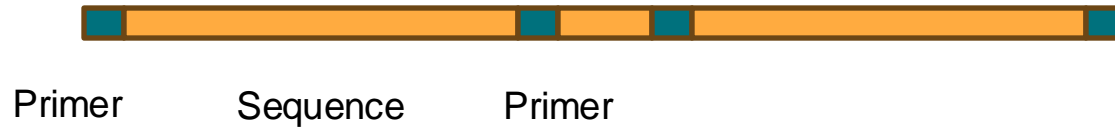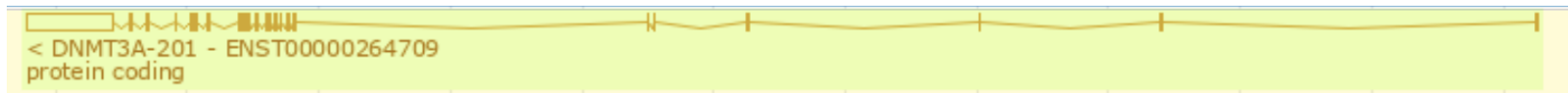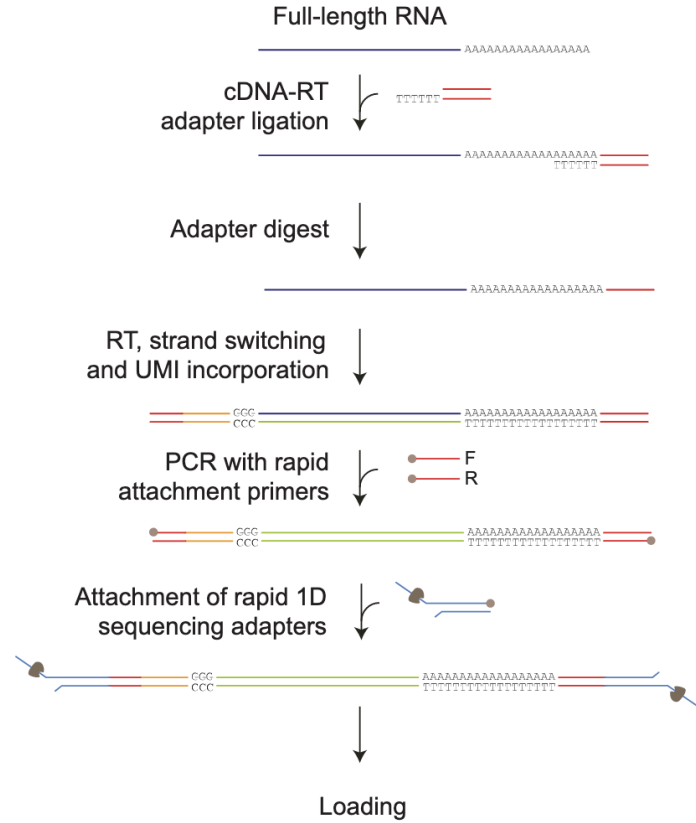
# How to estimate duplication rates

- In short read data, reads at the same position are assumed to be duplicates

- How do we know if we're saturating our libraries?



< DNMT3A-201 - ENST00000264709
protein coding

# UMIs



Full-length RNA

cDNA-RT adapter ligation

Adapter digest

RT, strand switching and UMI incorporation

PCR with rapid attachment primers

F
R

Attachment of rapid 1D sequencing adapters

Loading

# UMIs

- UMI at the 3' end of the read

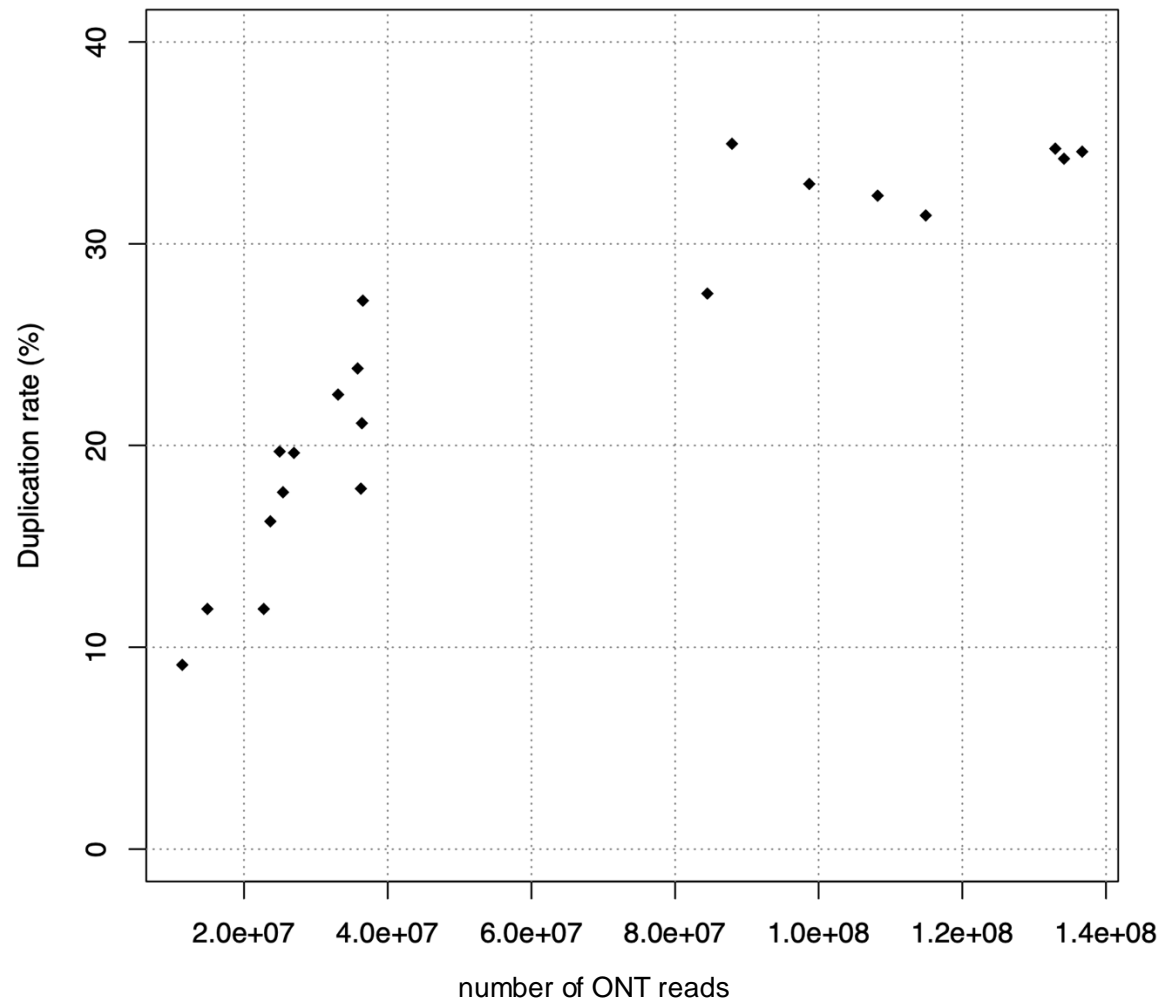TTT GGGG TT GGAA TT GGCG TT GGCA TTT

# UMIs

- UMI at the 3' end of the read

```
30 TTTCACCCTCCACTTCCCGTCTCAGAATT
29 TTTGAAACAGCTTCACCTTGAACTTT
29 TTTCCAATAAAAAAAATTACAATTT
29 TTTCAGCAAAATAAAATTCCGGTTT
27 TTTGGAGTTGGGGTTGCGCTTGGGGTTT
27 TTTGAGGTTGGAGTTGGGGTTGGCGTTT
24 TTTGGAGTTGGCGTTGCGGTTGGGGTTT
23 TTTGGGGTTGGAATTGGCGTTGGCATTT
23 TTTGGGATTAAGATTGGCATTGCGGTTT
23 TTTAGGGTTCGCGTTGGGGTTGCAGTTT
23 TTTAGGGTTAGCGTTGGAGTTGGGGTTT
22 TTTGGCGTTGGGGTTGGCGTTGGCGTTT
22 TTTGGCGTTGGAGTTCAGCTTACGGTTT
22 TTTGCGGTTGGAGTTGGGCTTGGCGTTT
22 TTTACACTTGTGCTCTCCTTAGCCTTT
21 TTTGGGGTTGGAGTTGGCGTTGGCATTT
21 TTTGGCGTTGGCATTGGCGTTGGGGTTT
21 TTTGGCGTTCGGGTTGGAATTCGCGTTT
```

# UMIs

- UMI at the 3' end of the read

- Different lengths indicative of high
  error rate

- only 47% of reads have fully intact UMI
- 7% have no UMI at all

- Even using some error correction with
  Levenshtein distance, it's ugly

$0.98\^28 = 0.56$

```
30  TTTCACCCTCCACTTCCCGTCTCAGAATT
29  TTTGAAACAGCTTCACCTTGAACTTT
29  TTTCCAATAAAAAAAATTACAATTT
29  TTTCAGCAAAATAAAATTCCGGTTT
27  TTTGGAGTTGGGGTTGCGCTTGGGGTTT
27  TTTGAGGTTGGAGTTGGGGTTGGCGTTT
24  TTTGGAGTTGGCGTTGCGGTTGGGGTTT
23  TTTGGGGTTGGAATTGGCGTTGGCATTT
23  TTTGGGATTAAGATTGGCATTGCGGTTT
23  TTTAGGGTTCGCGTTGGGGTTGCAGTTT
23  TTTAGGGTTAGCGTTGGAGTTGGGGTTT
22  TTTGGCGTTGGGGTTGGCGTTGGCGTTT
22  TTTGGCGTTGGAGTTCAGCTTACGGTTT
22  TTTGCGGTTGGAGTTGGGCTTGGCGTTT
22  TTTACACTTGTGCTCTCCTTAGCCTTT
21  TTTGGGGTTGGAGTTGGCGTTGGCATTT
21  TTTGGCGTTGGCATTGGCGTTGGGGTTT
21  TTTGGCGTTCGGGTTGGAATTCGCGTTT
```

# What does the data look like?



Genomic DNA – standard prep

# What does the data look like?



RNA/cDNA – standard prep

# What does the data look like?



Observed Read lengths (Gb)

Transcript Lengths (knownGene hg38)
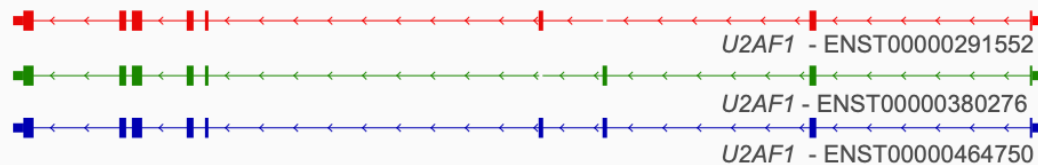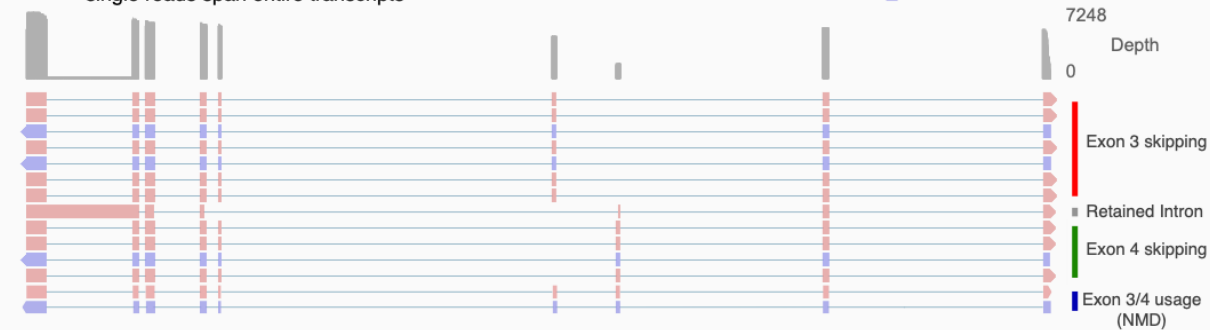
cDNA – standard prep

[0 - 1138]

U2af1

U2AF1

Short-read coverage (65M reads)
single reads span 1-3 exons

716
Depth
0

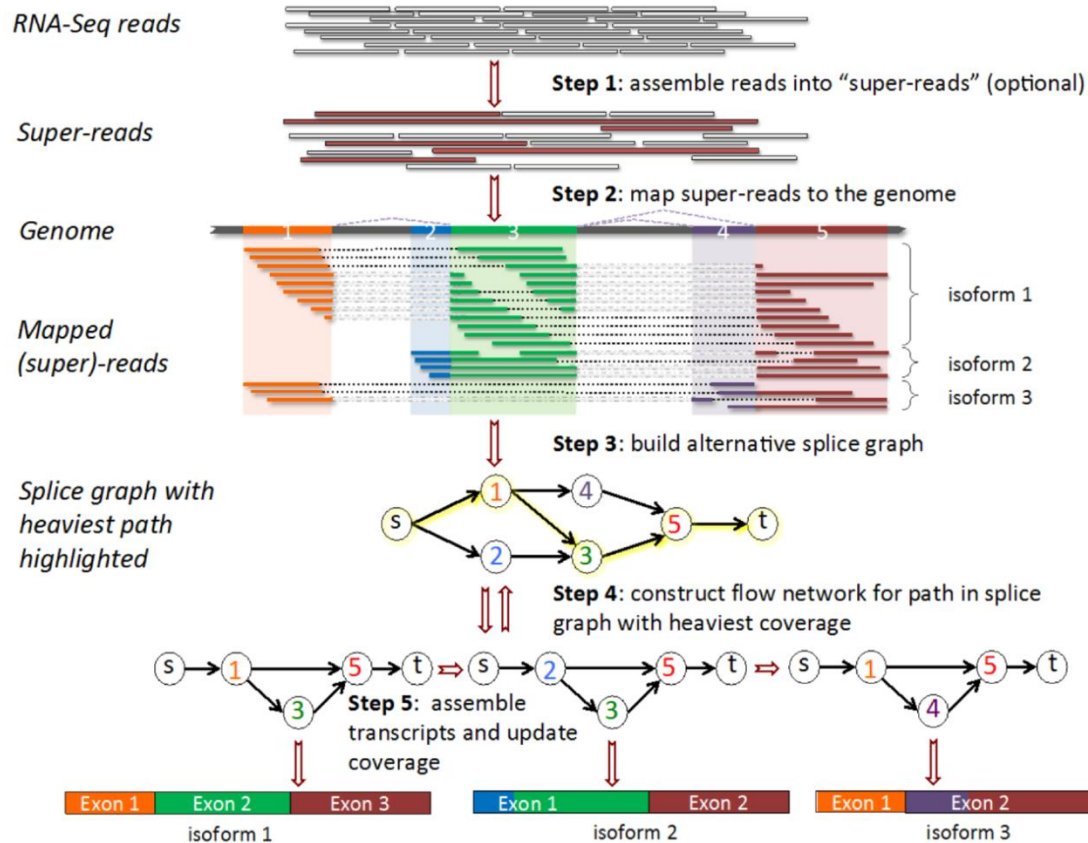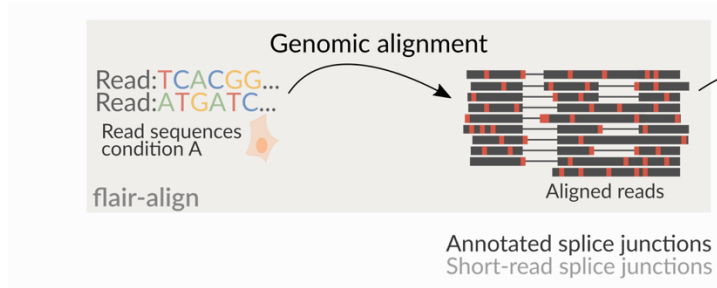Long-read coverage (48M reads)
single reads span entire transcripts

7248
Depth
0

Exon 3 skipping

Retained Intron

Exon 4 skipping

Exon 3/4 usage
(NMD)

*U2AF1* - ENST00000291552

*U2AF1* - ENST00000380276

*U2AF1* - ENST00000464750

chr21    43,095,000    43,100,100    43,105,000

# Estimating transcript abundance – short-read, Stringtie



Pertea et al. Nature Biotechnology, 2015
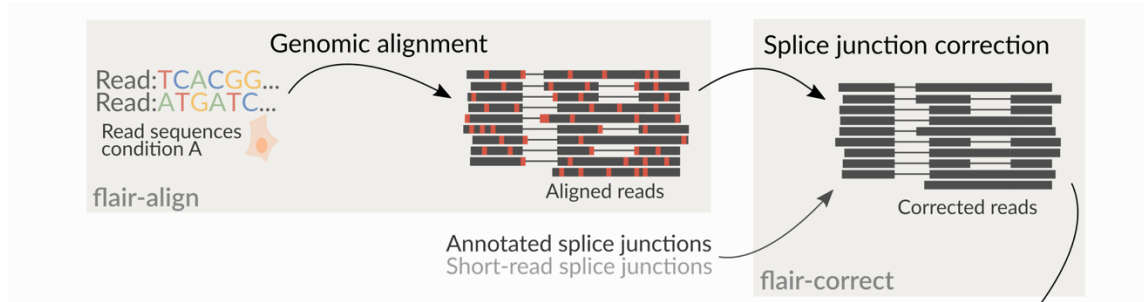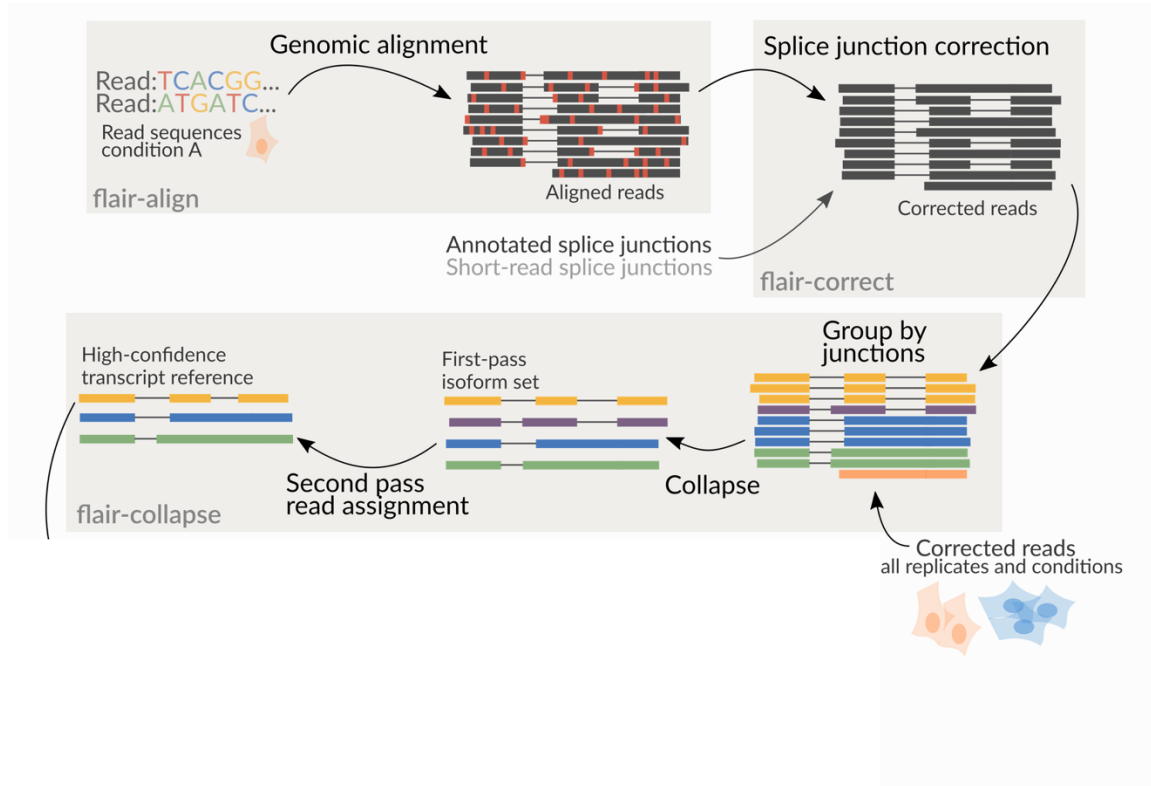
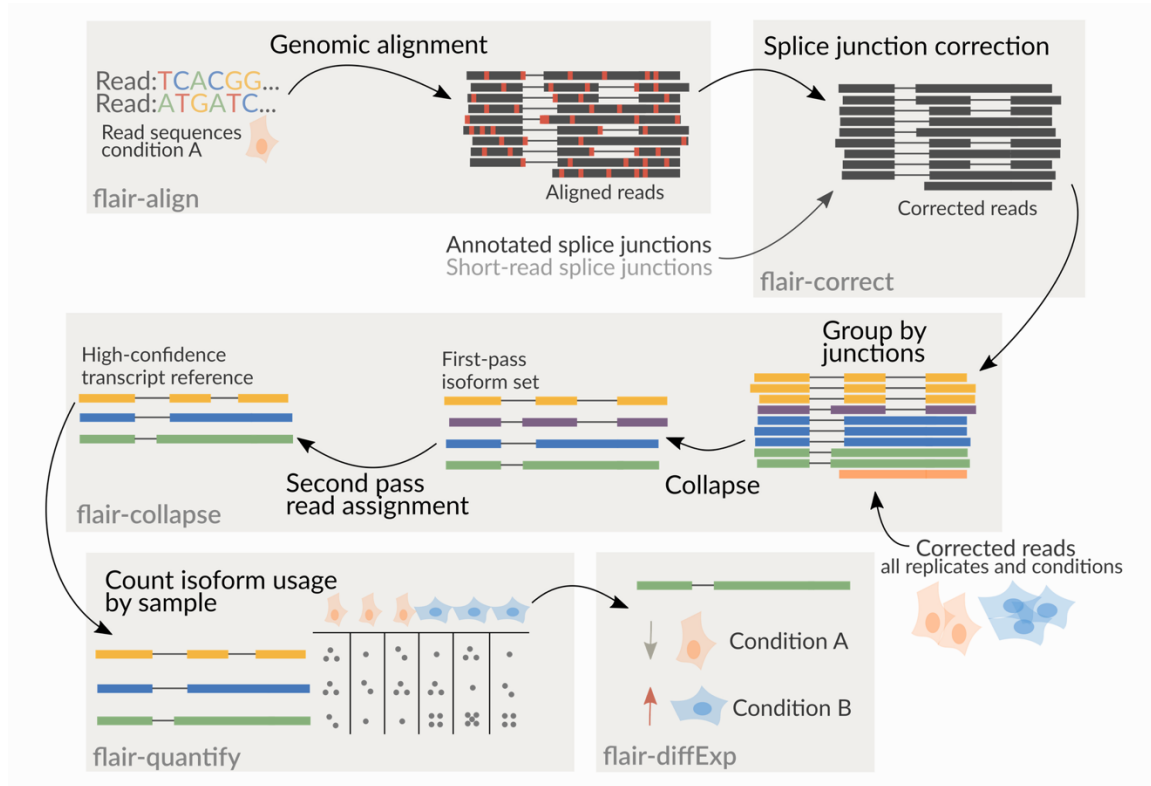# Estimating transcript abundance – long read

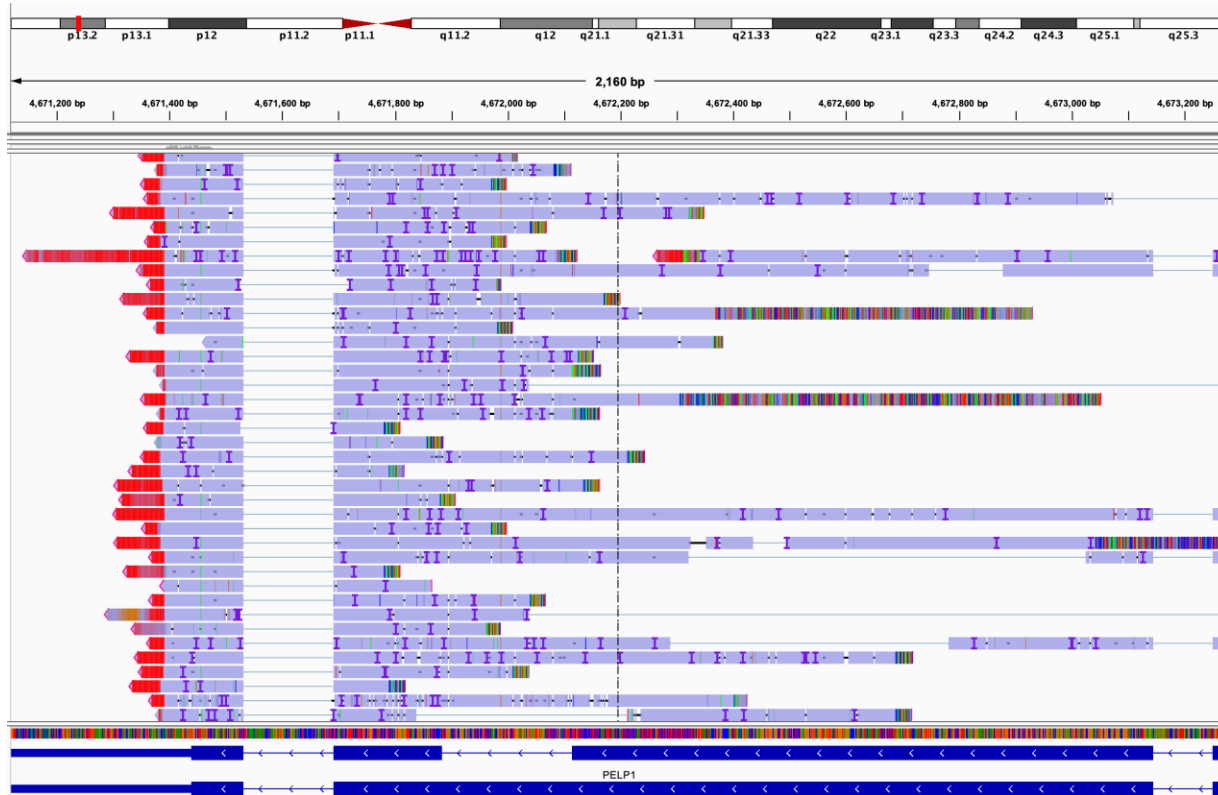# Estimating transcript abundance – long read

# Estimating transcript abundance – long read

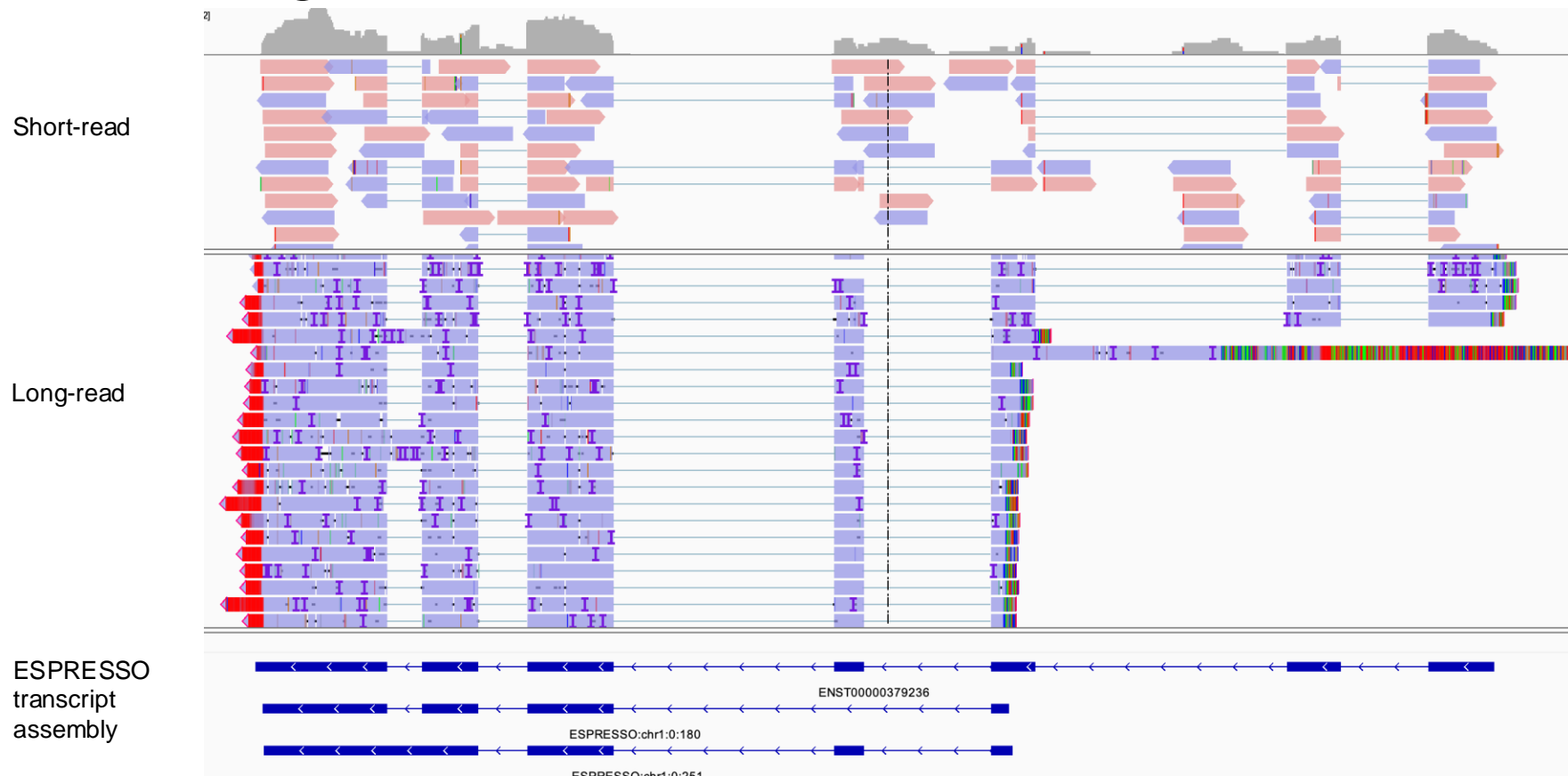# Estimating transcript abundance – long read

# Technical artifacts



Looks like fragmentation of this RNA throughout the long exon. (This is an egregious case – most are more mild)
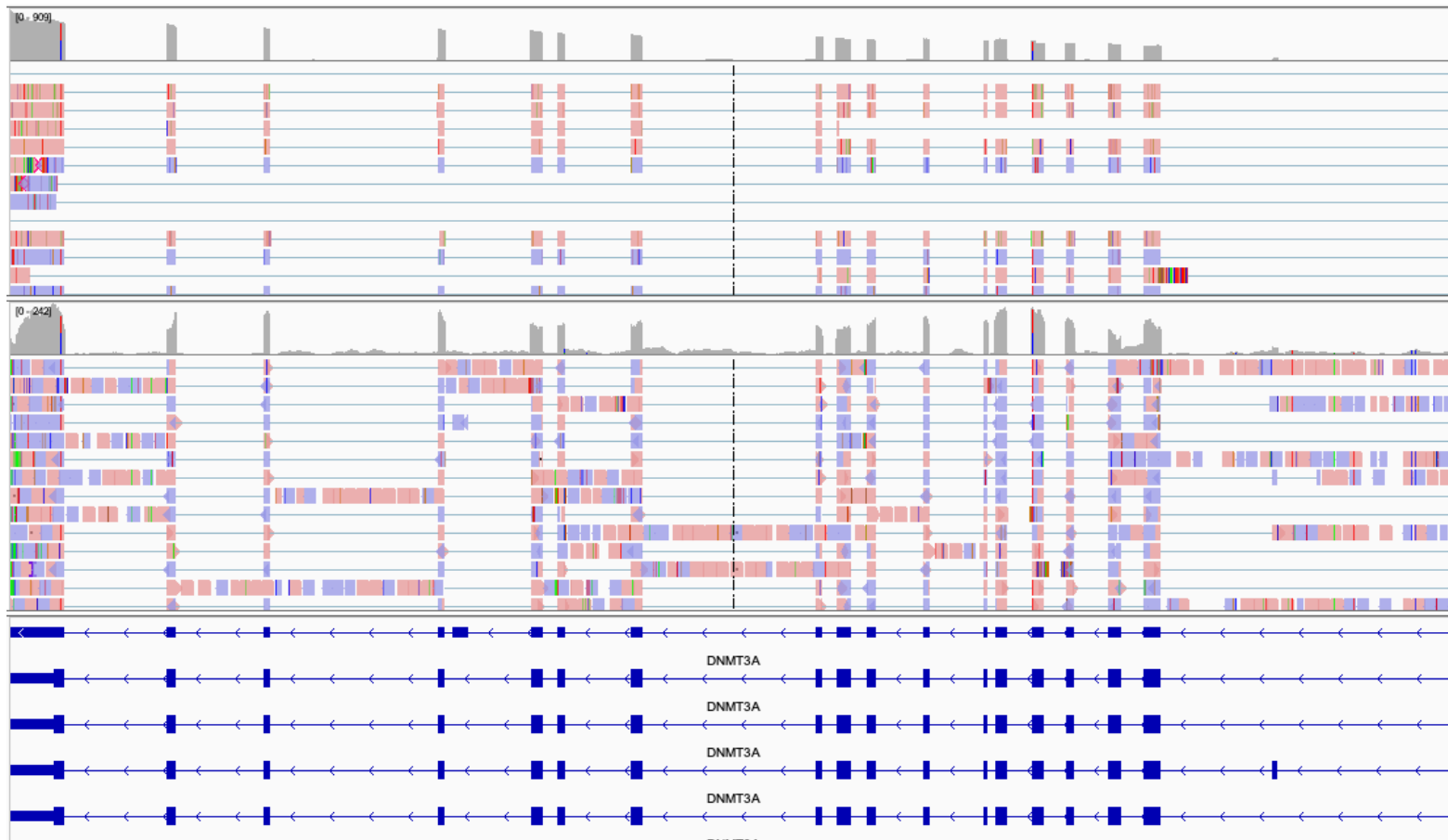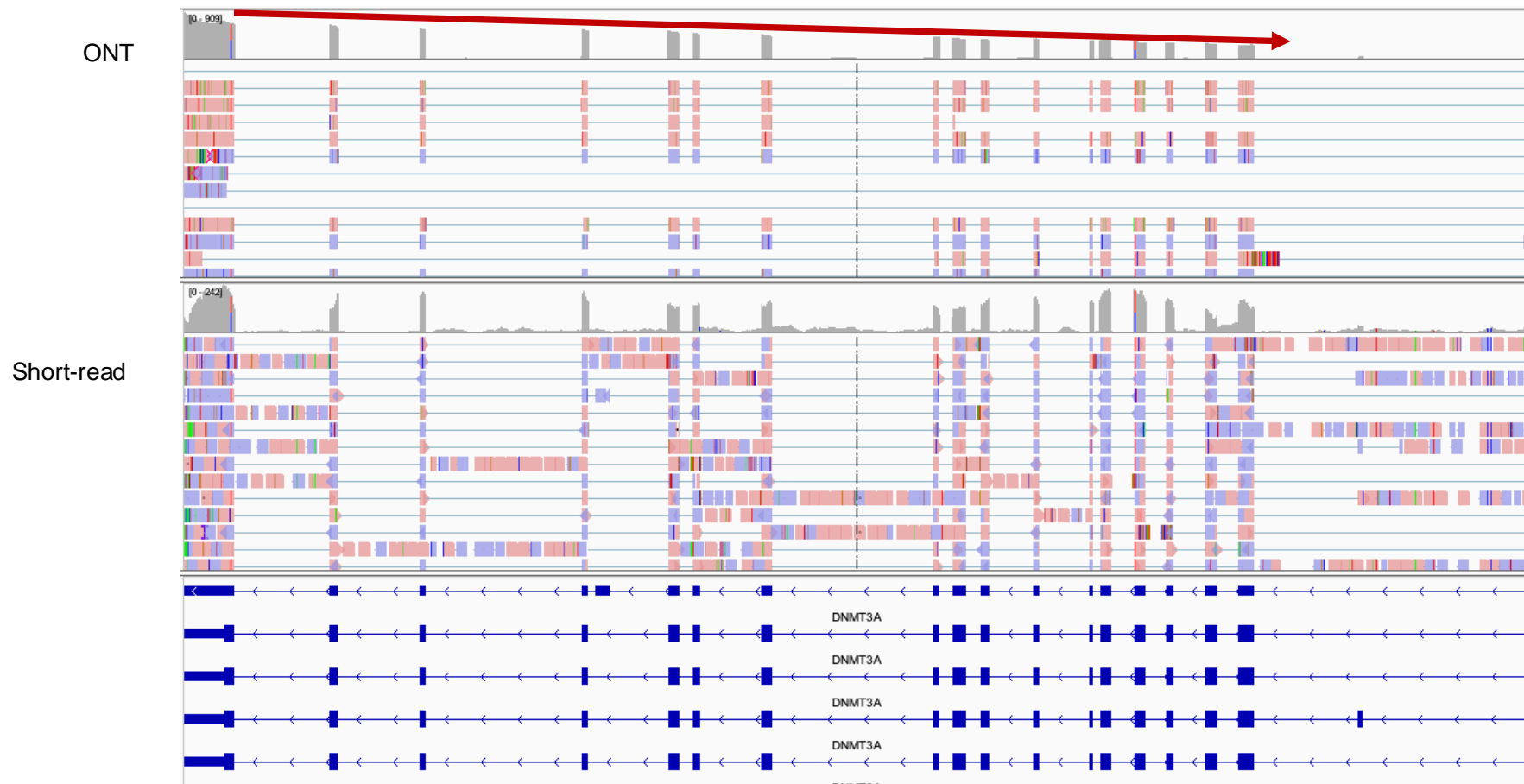
# "Full-length" reads



Short-read

Long-read

ESPRESSO
transcript
assembly

ENST00000379236
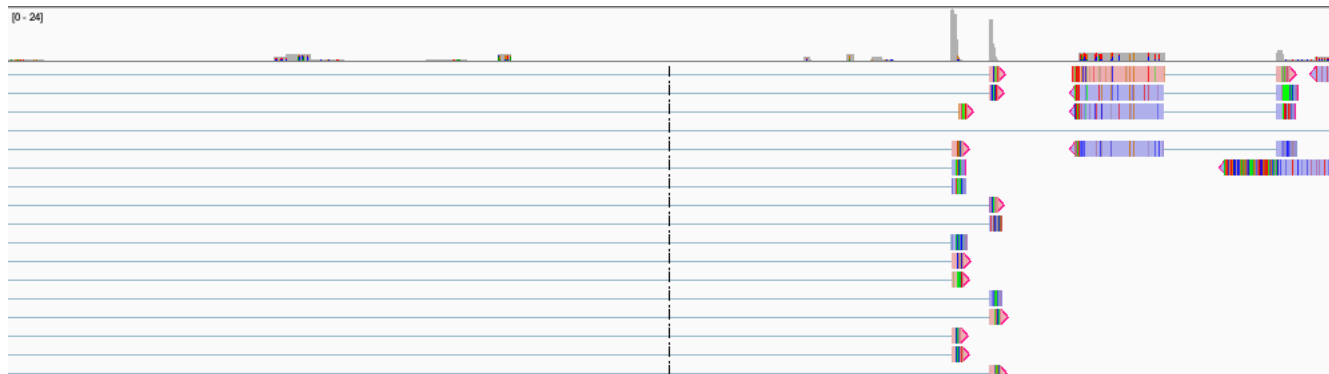
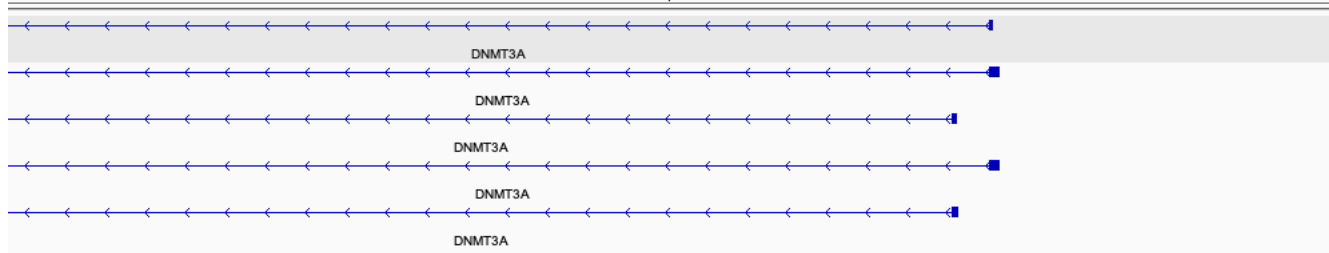ESPRESSO:chr1:0:180

ESPRESSO:chr1:0:251

# human DNMT3A

# human DNMT3A

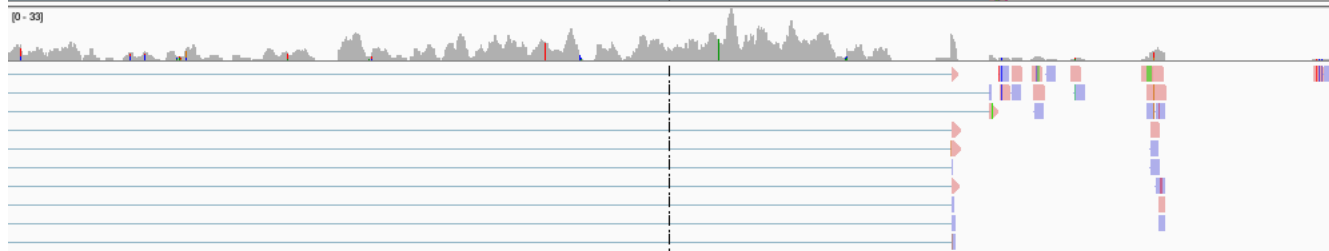# human DNMT3A

# Truncated reads

- Appear to be caused by RNA fragmentation

- assessing RIN values of your samples can help – choose clean ones when possible

- When not possible, iteratively assemble transcripts and remove non-full-length reads

# Assignment

- Start with some long-read RNAseq data from a cell line

- QC the data, trim adapters

- Align the reads

- Examine a few genes to see how the data looks