

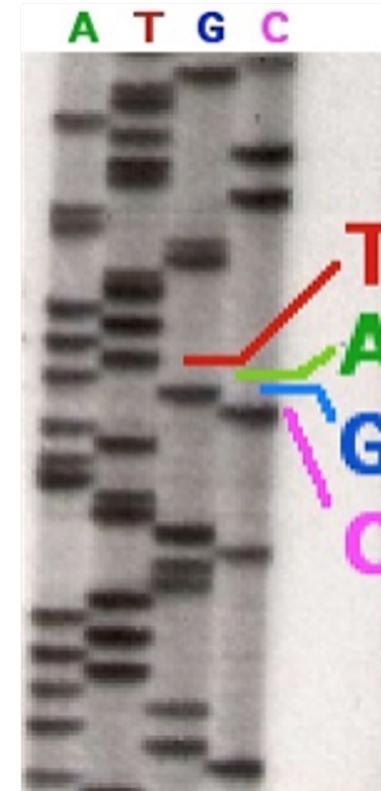
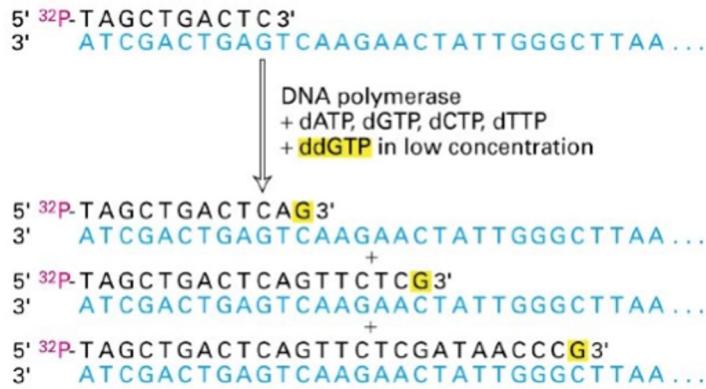
Long Read Sequencing

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

How to sequence a human genome: Sanger method

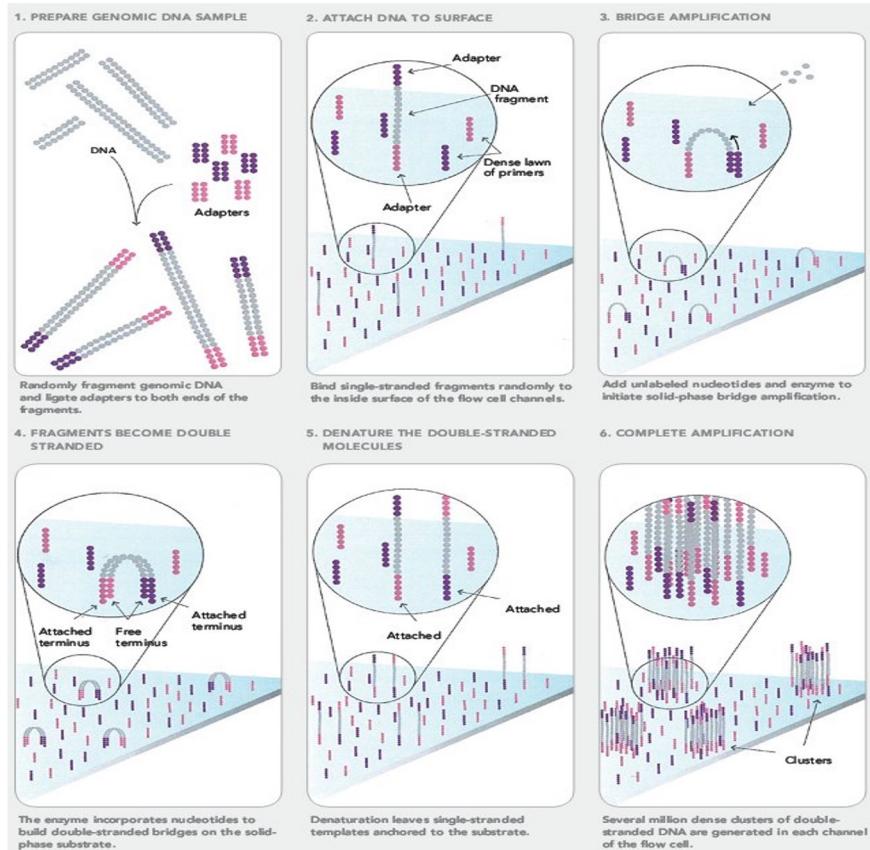
Key points:

- 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

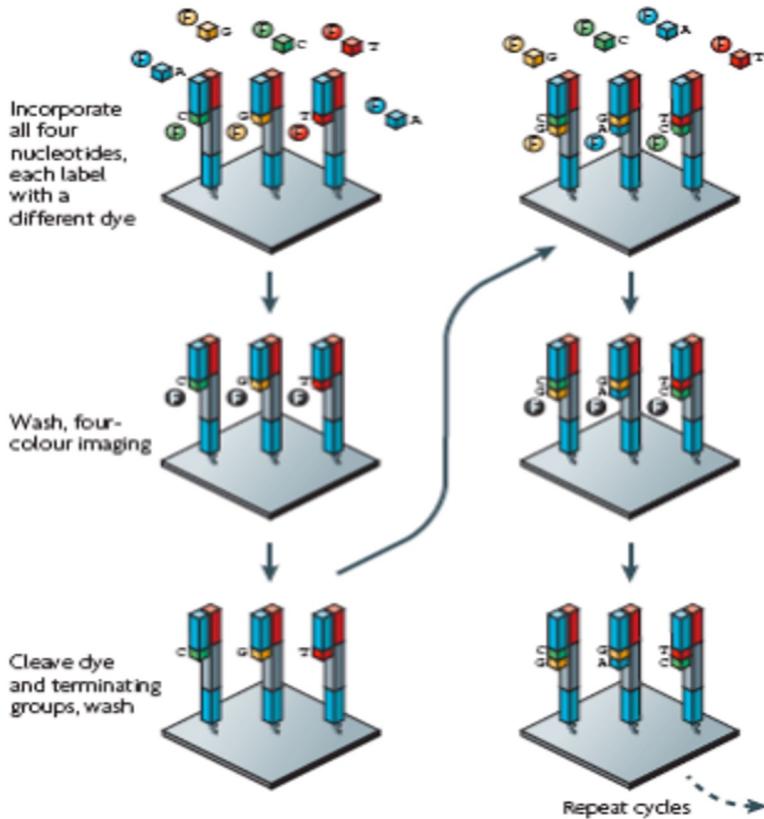


Solexa (Illumina) sequencing (2006)

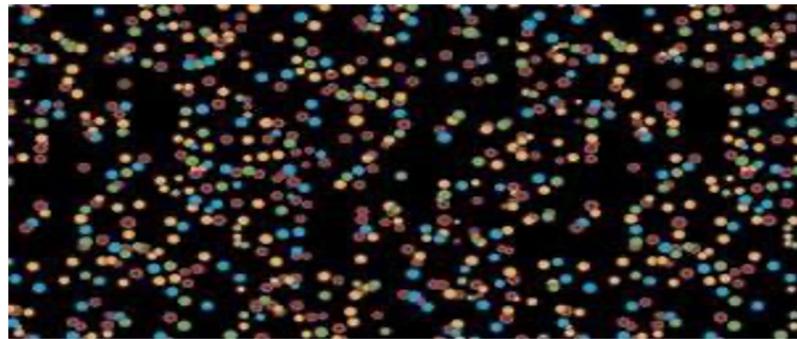
- PCR amplify sample (opt.)
- Immobilize and amplify single molecules on a solid surface
- Reversible terminator sequencing with 4 color dye-labelled nucleotides



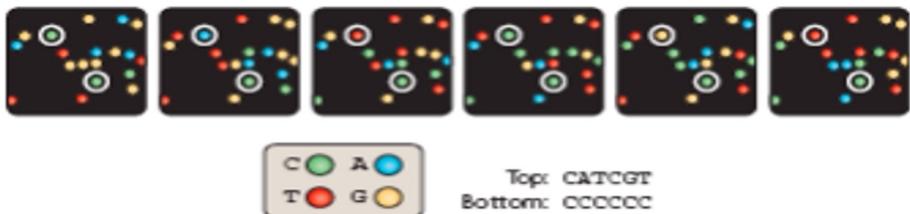
Illumina sequencing (2005)



4 different images merged



6 cycles w/ base-calling



Paired-end sequencing: A molecular hack to sequence longer fragments

genomic DNA



Shear to desired length (~400bp)

DNA fragments

ligate adapters, size select

sequencing library



Illumina GA2



clusters on a flow-cell



millions to billions of paired-end reads (readpairs)

~150bp

~200bp

~150bp

5' GGTGTACGAATAGTTCCCTTACACTCCTGACCATCCTAGC

GGACTGAAACTTCATCTGTCTTATAGATATGCGTGCAGCAGC 5'

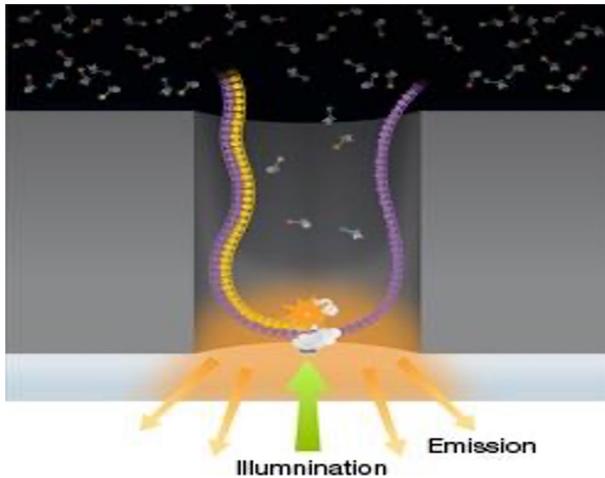
// //

// //

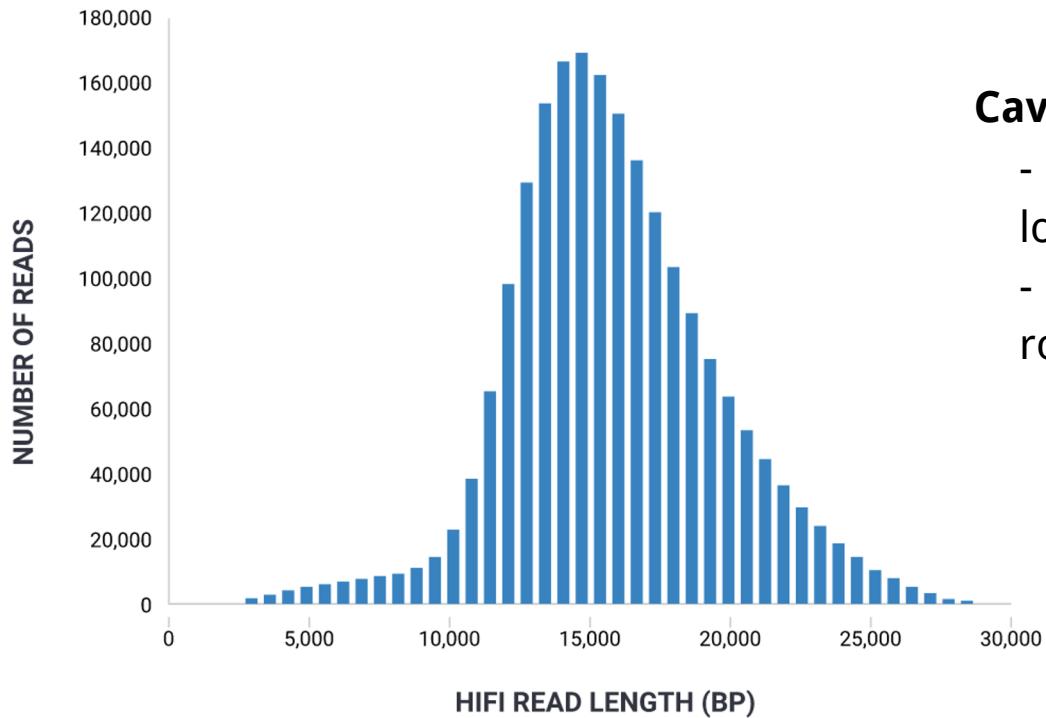
Pacific Biosciences

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



Pacific Biosciences: long reads. Great for genome assembly

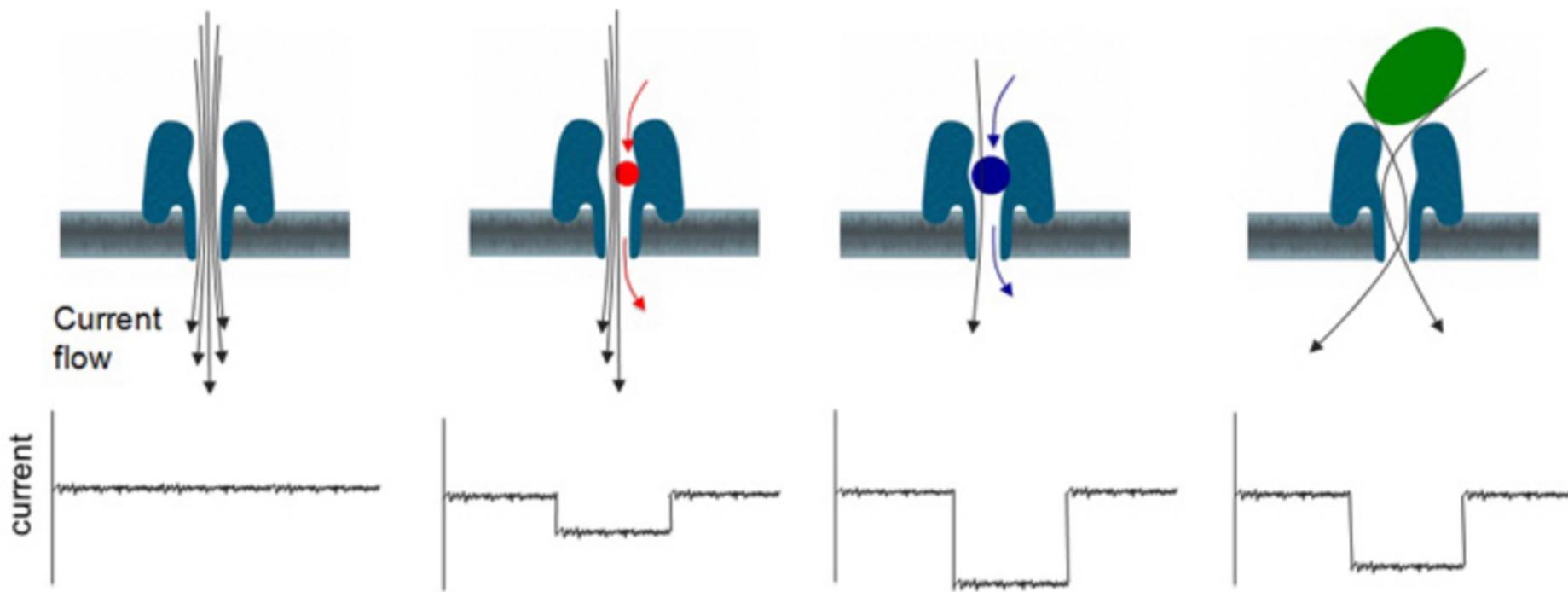


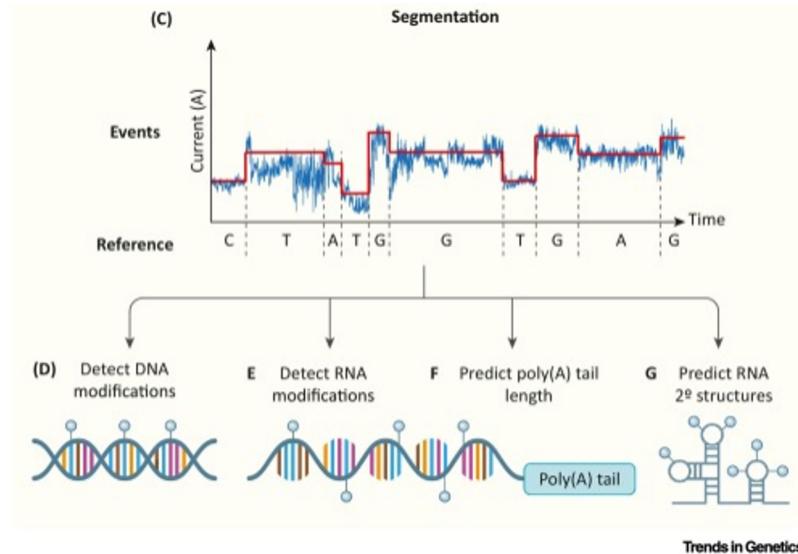
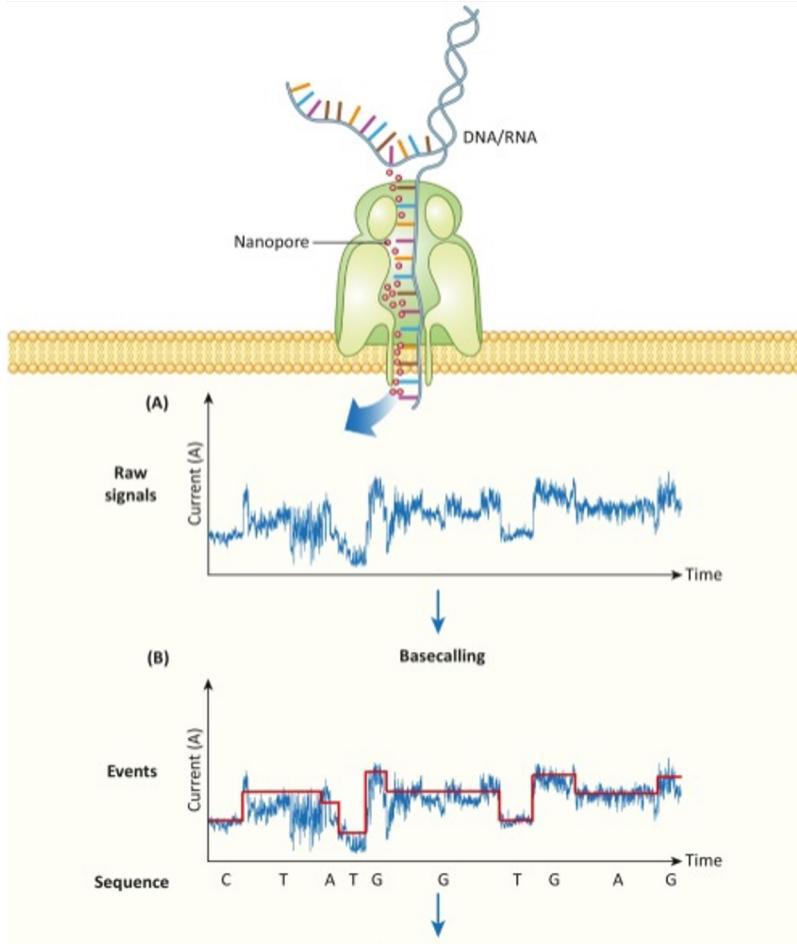
Caveats:

- 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

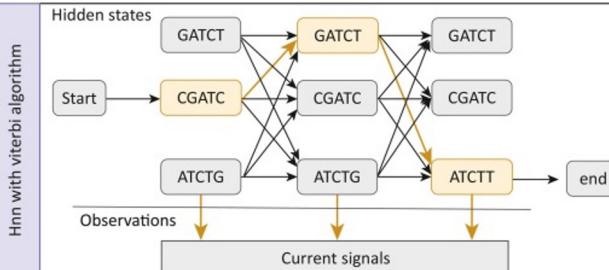
Oxford Nanopore Technologies



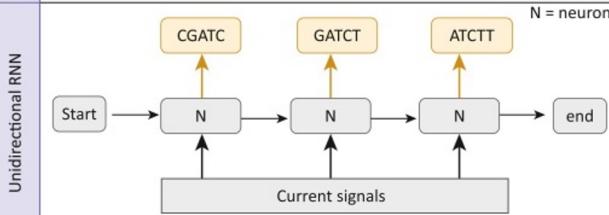


doi.org/10.1016/j.tig.2021.09.001

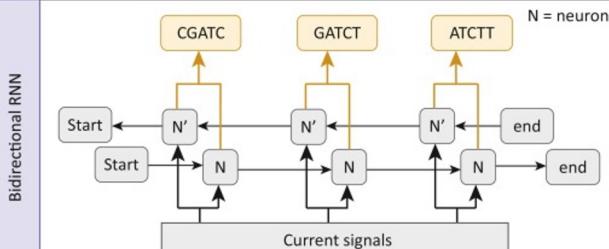
(A)



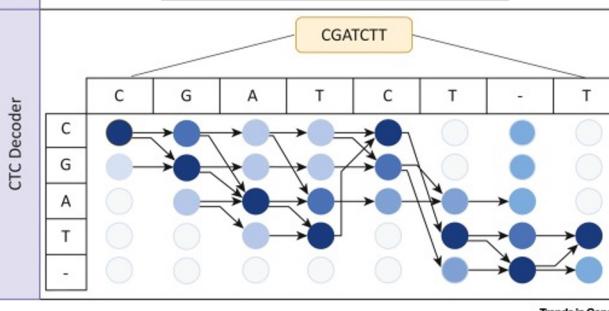
(B)



(C)



(D)



Neural networks to translate signal into base calls

- Guppy (many versions)

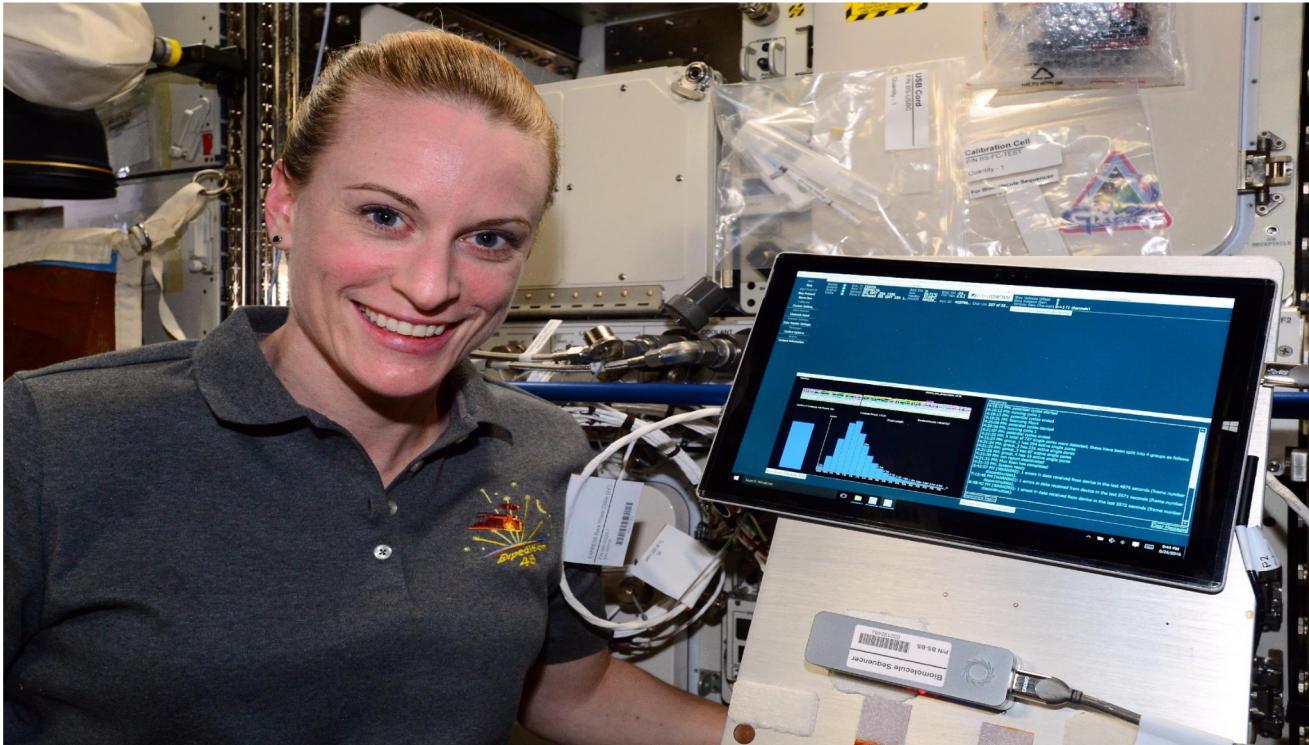
- Dorado (v0.4, eventual guppy replacement)

- many others

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

doi.org/10.1016/j.tig.2021.09.001

Nanopore sequencing is *extremely* portable

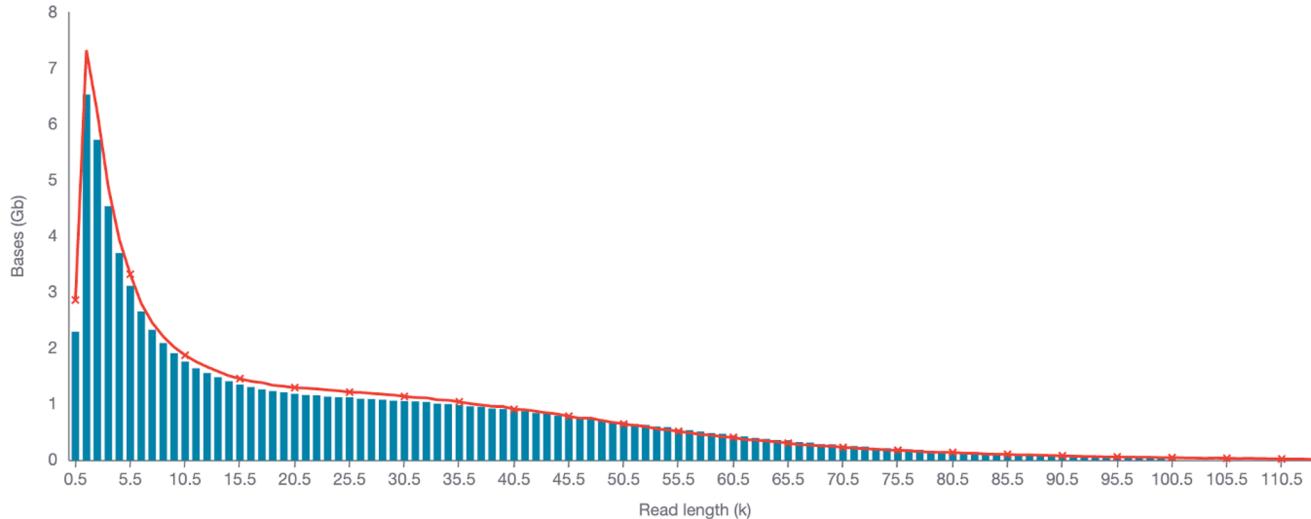


Kate Rubins sequencing DNA on the ISS

ONT sequence length distribution

Legend

Basecalled Estimated



Estimated N50

17.75 kb

% Basecalled

100%

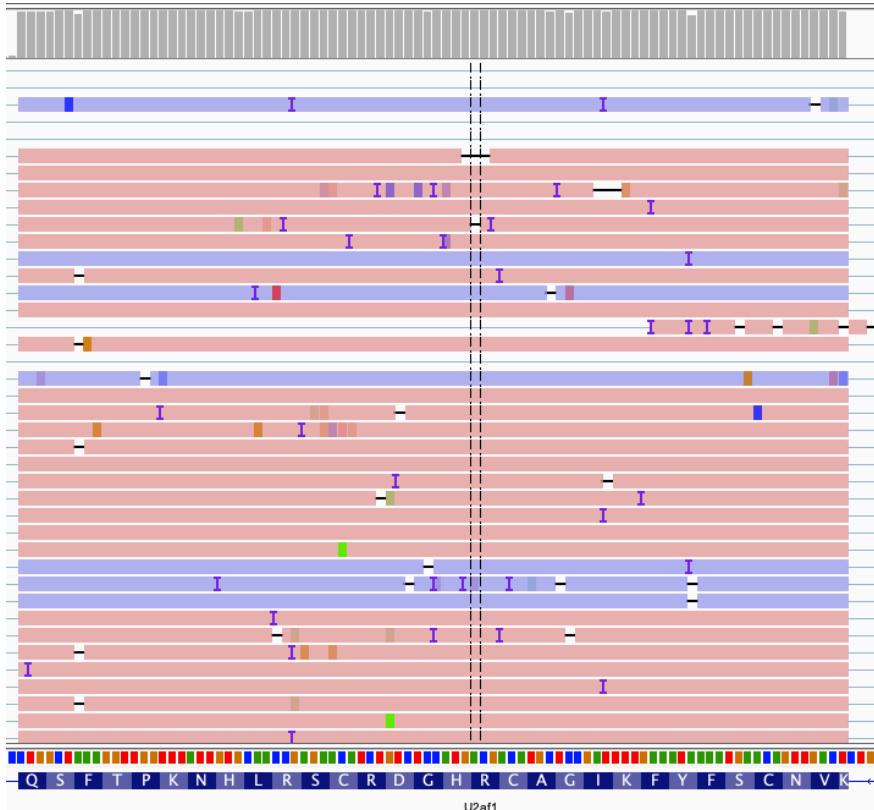
their relative amounts.

Read length (kb)	Aggregated reads (Mb)
100 - 164	886.98
164 - 228	36.06
228 - 292	4.02
292 - 344	0.35

Recent run of a tumor sample

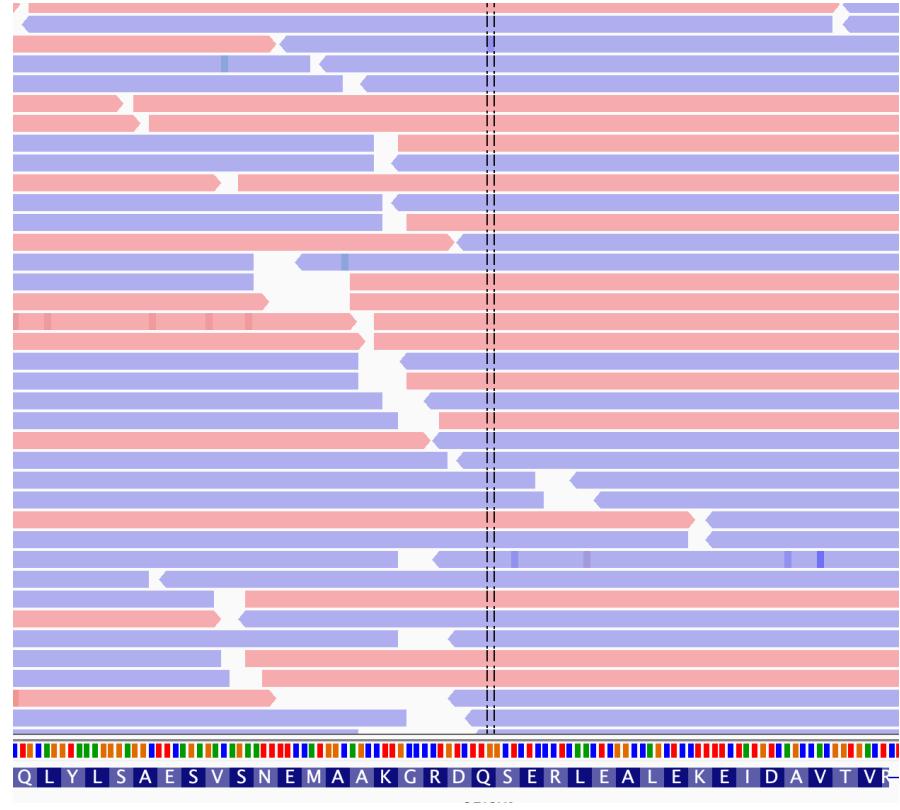
About \$3,500 for a
30x human genome
on a PromethION

What does the data look like?



Long-read ONT

~5% base error rate

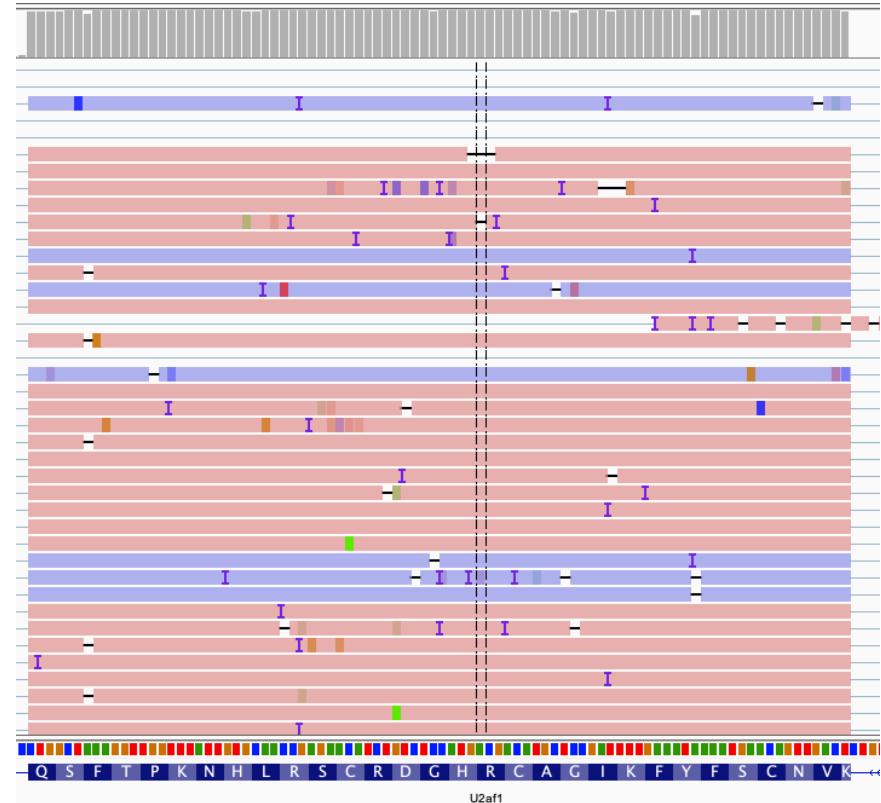


Short-read Illumina

~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?

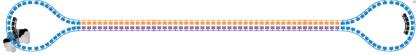
Start with high-quality double stranded DNA



Prepare SMRTbell libraries



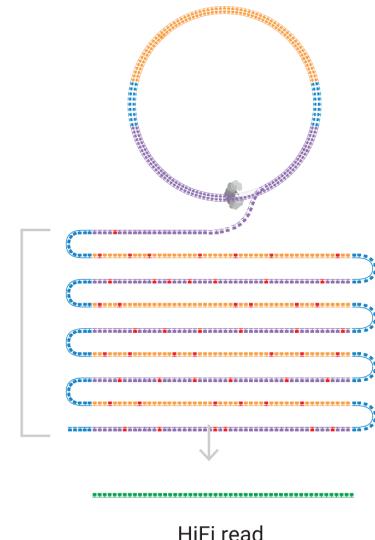
Anneal primers and bind DNA polymerase



Circularized DNA is sequenced in repeated passes

The polymerase reads are trimmed of adapters to yield subreads

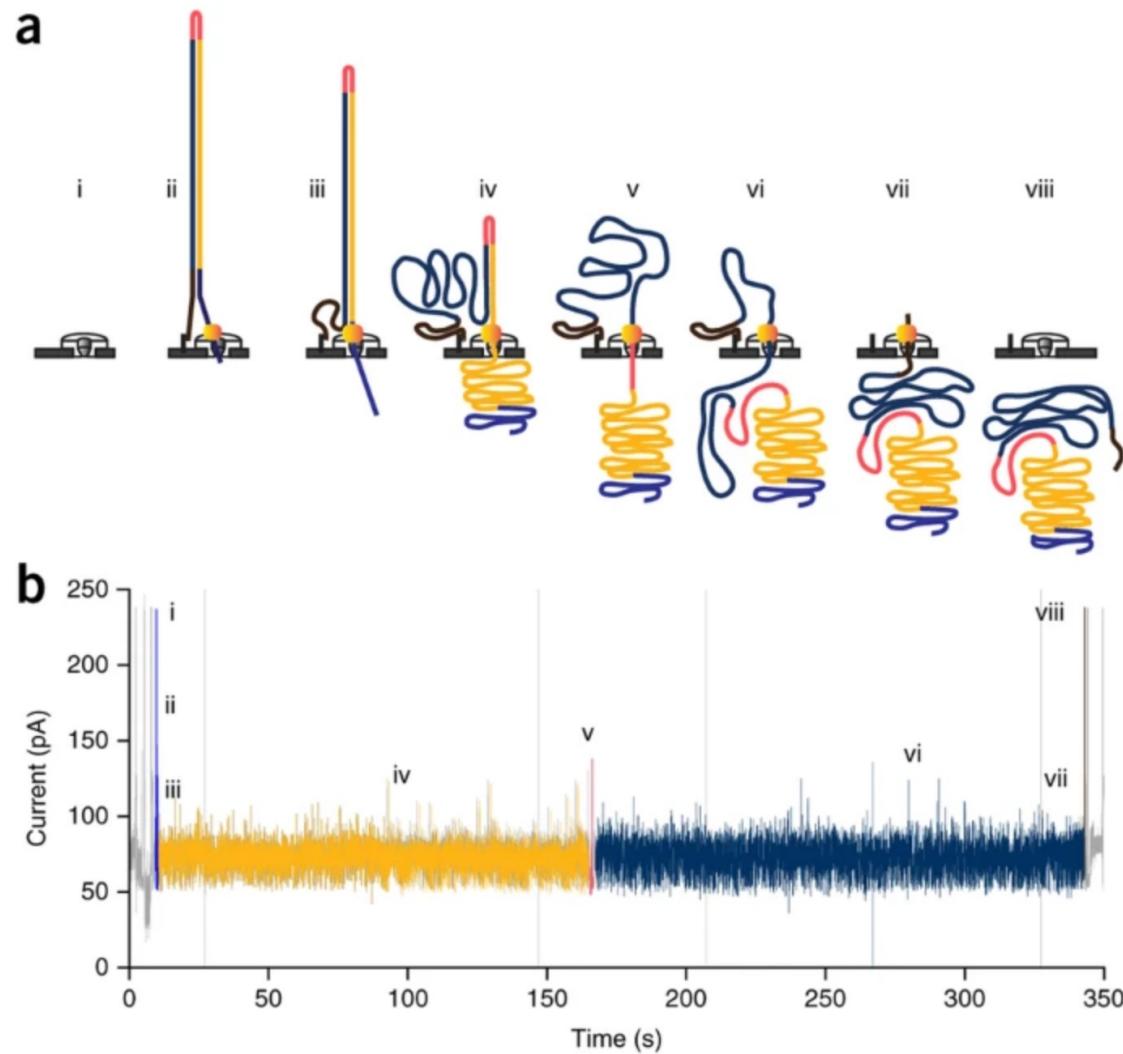
Consensus and methylation status are called from subreads



Improved error rates

higher cost/lower throughput

ONT Duplex sequencing



Improved error rates

higher cost/lower throughput

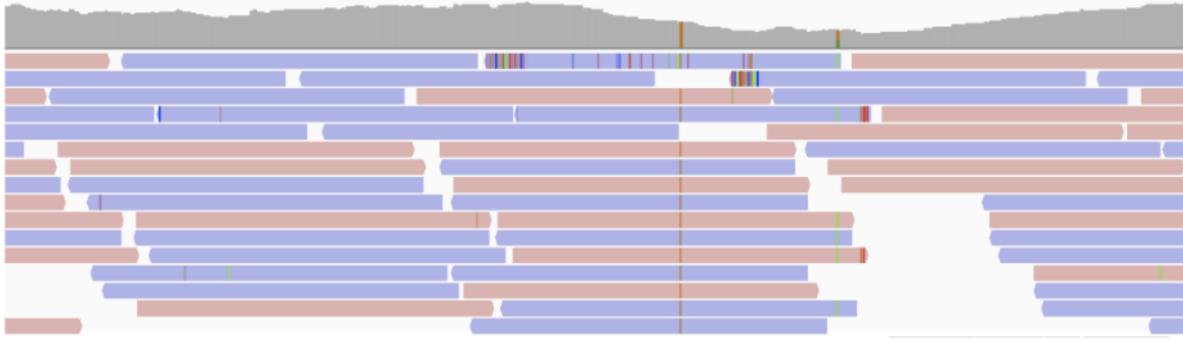
Genomic DNA advantages



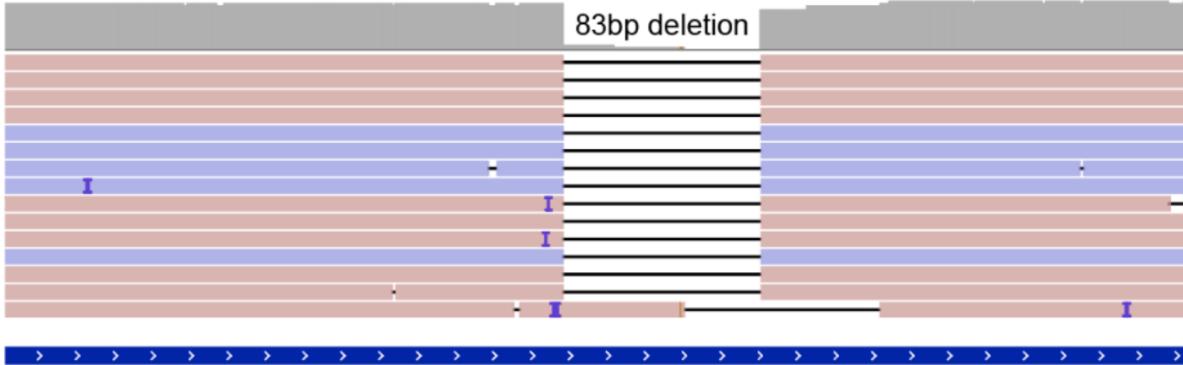
Figure 1: Blue-labeled genomic regions are accessible to long reads but not short, and have functional annotations (e.g. genes or enhancers)

Large Indel detection

No indel detectable - Short-read sequencing - Illumina



83bp deletion - Long-read sequencing - Oxford Nanopore



chr3: 31990200-31990700

ZNF860

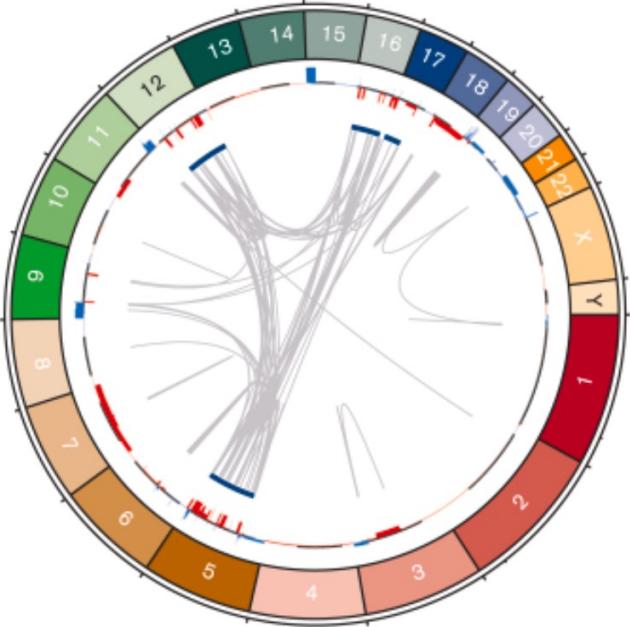
(protein-coding sequence)

Haley Abel

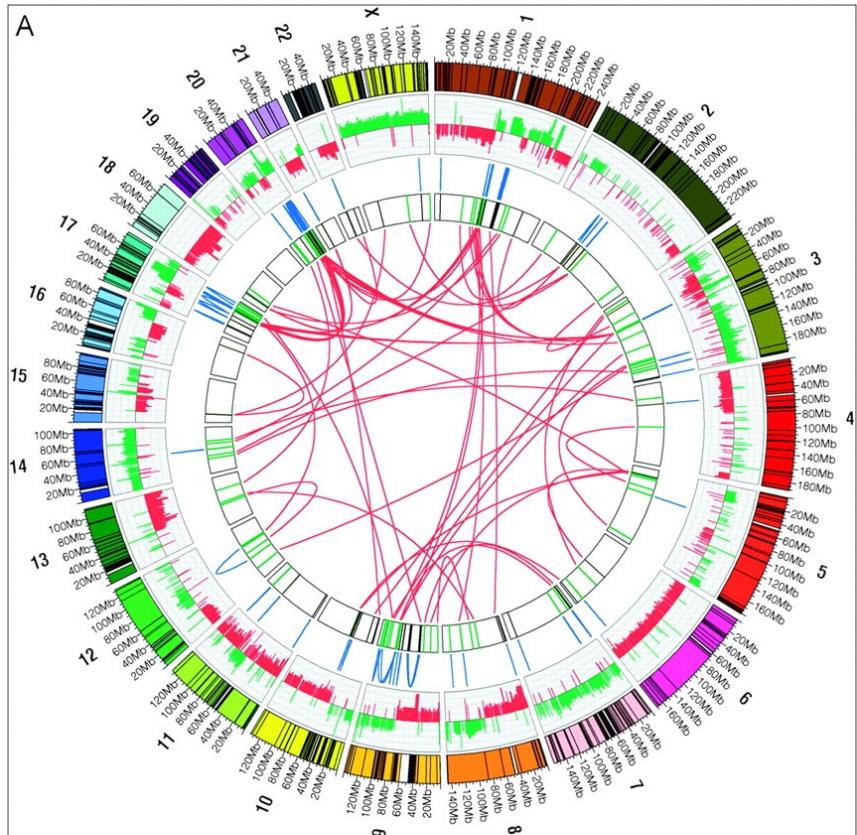
Structural variant resolution

MCF7 Breast Cancer Cell line

TP53-mutated AML



doi: 10.1182/bloodadvances.2023010156

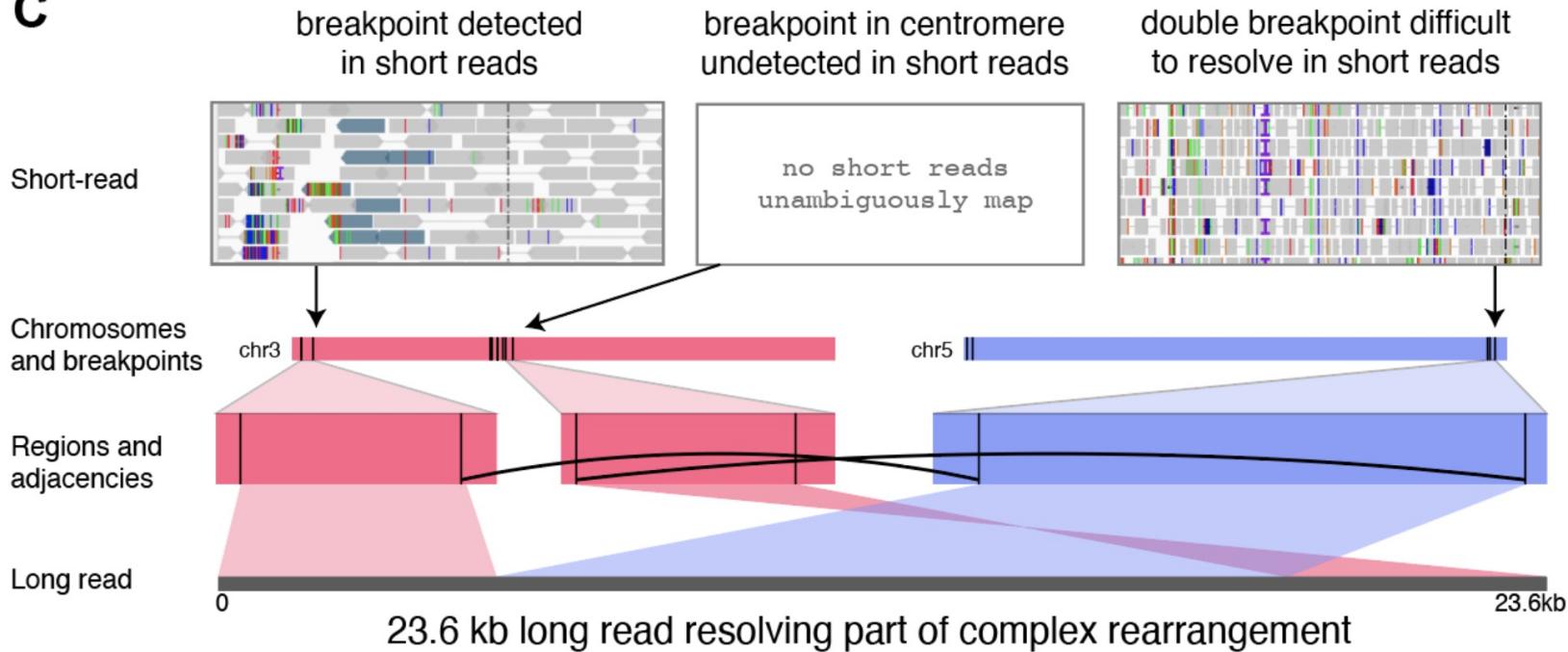


Hampton, et al. doi: 10.1101/gr.080259.108

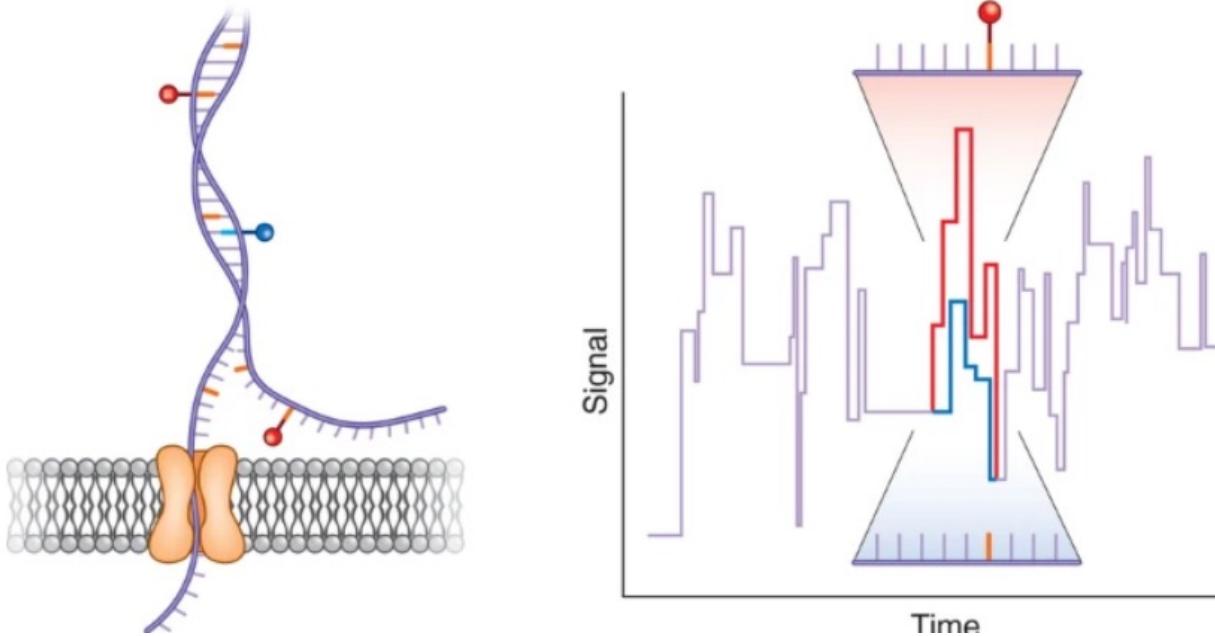


Structural variant resolution

C



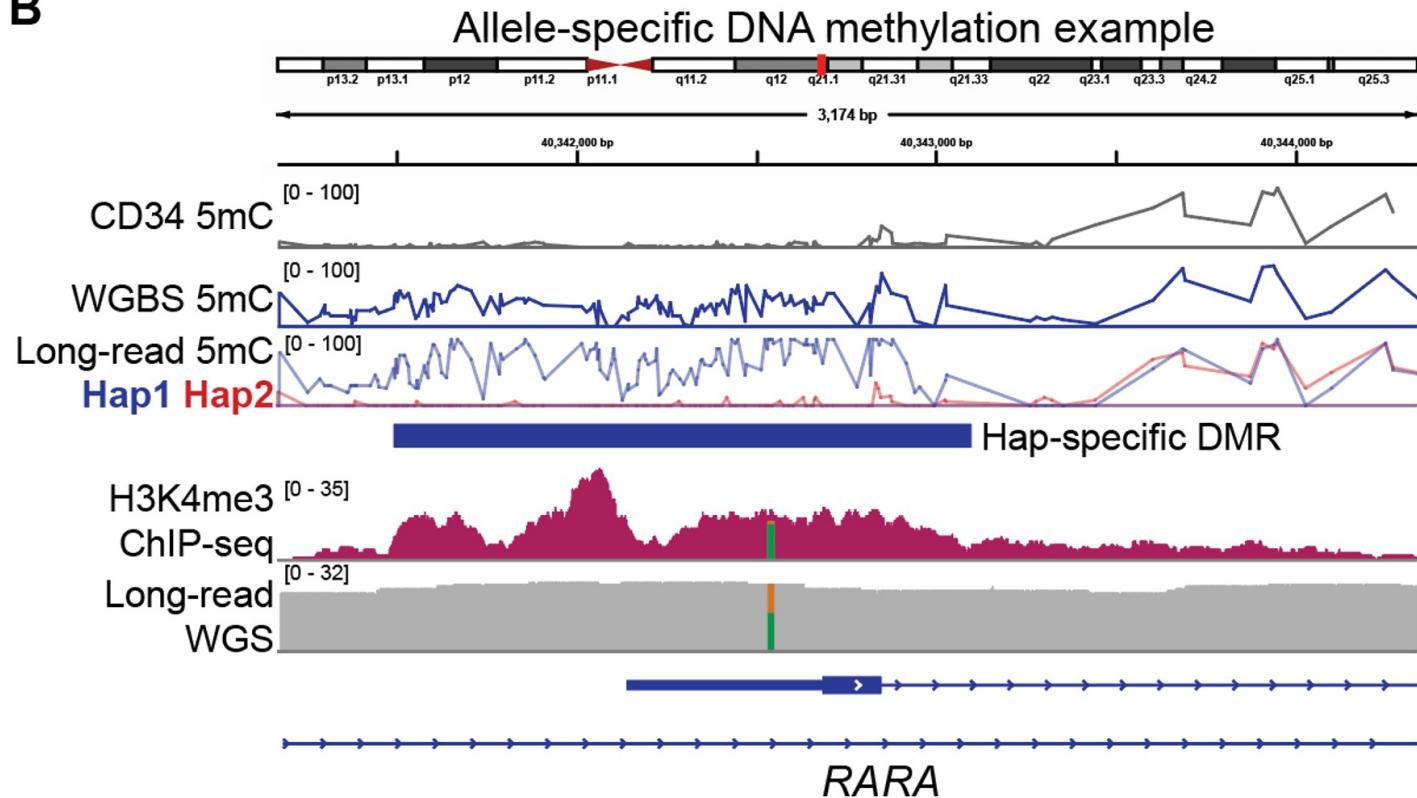
Base modification detection



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

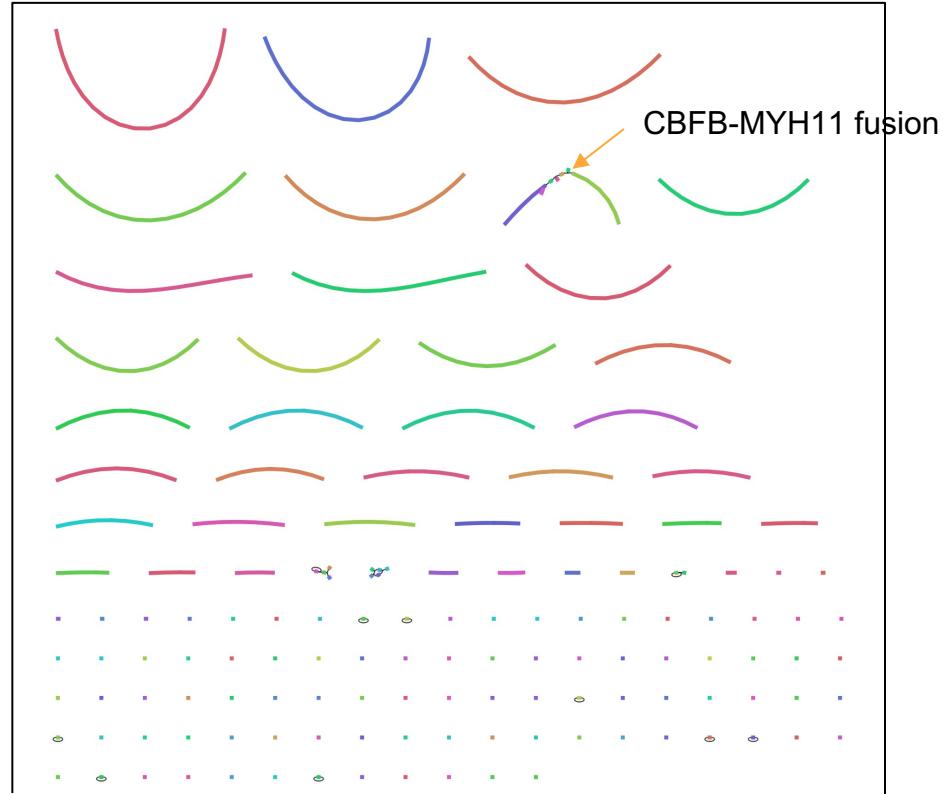
Phasing of reads/modifications

B



Genome assembly

- Assembly of personal genomes

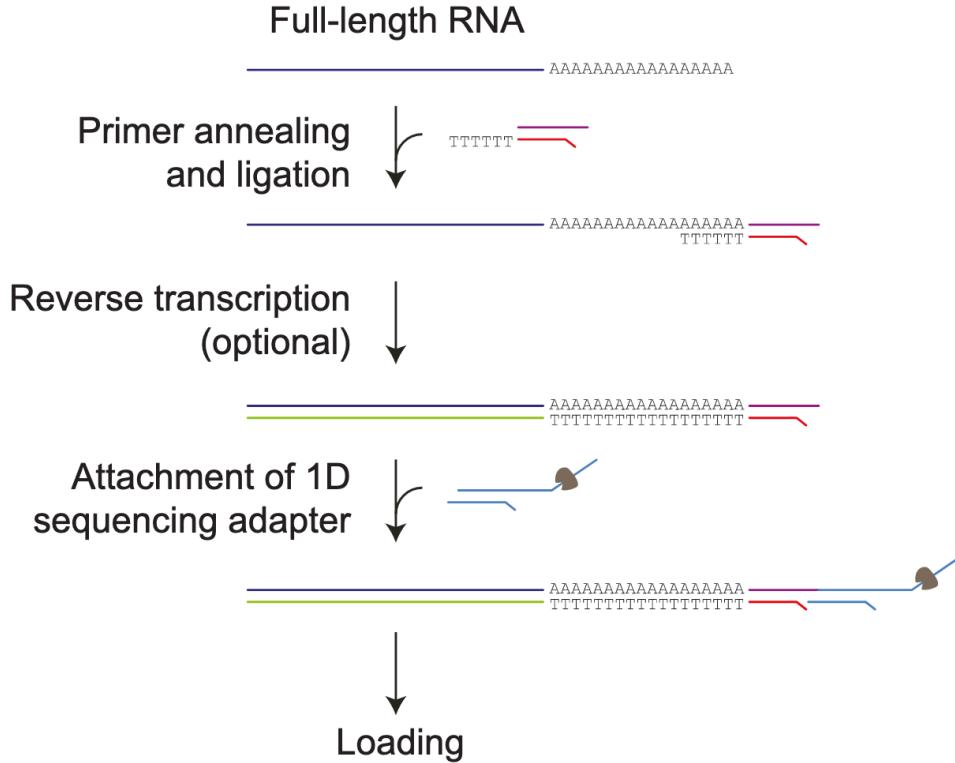


sample	hap	metric	value
Presentation	hap1	Number_of_contigs	402
Presentation	hap1	N50	88.3 Mbp

Dave Spencer, Haley Abel

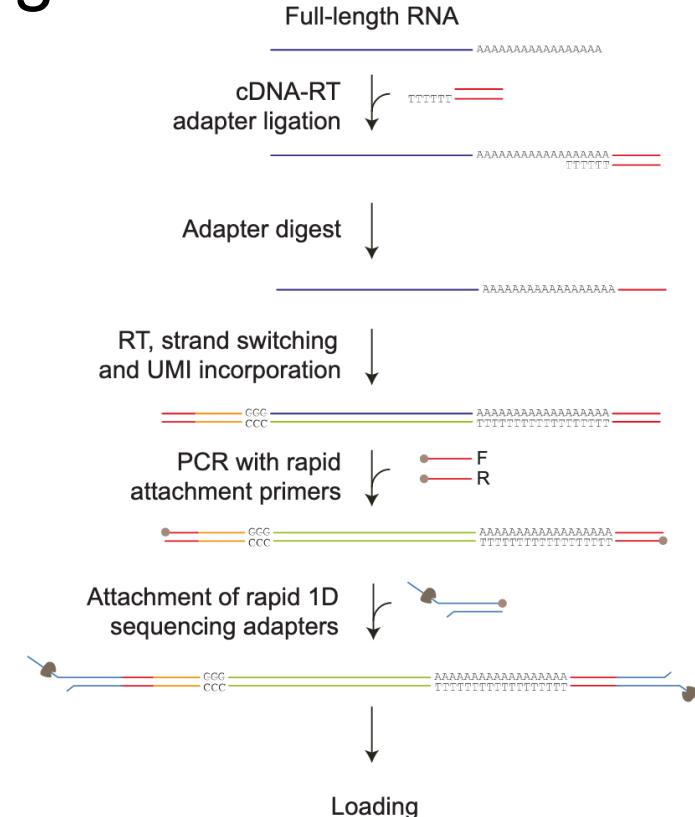
Long-read RNA sequencing

- Direct RNA
- No amplification, less bias
- Preserves base modifications (m6a, etc)



Long-read RNA sequencing

- cDNA sequencing
- much higher yields



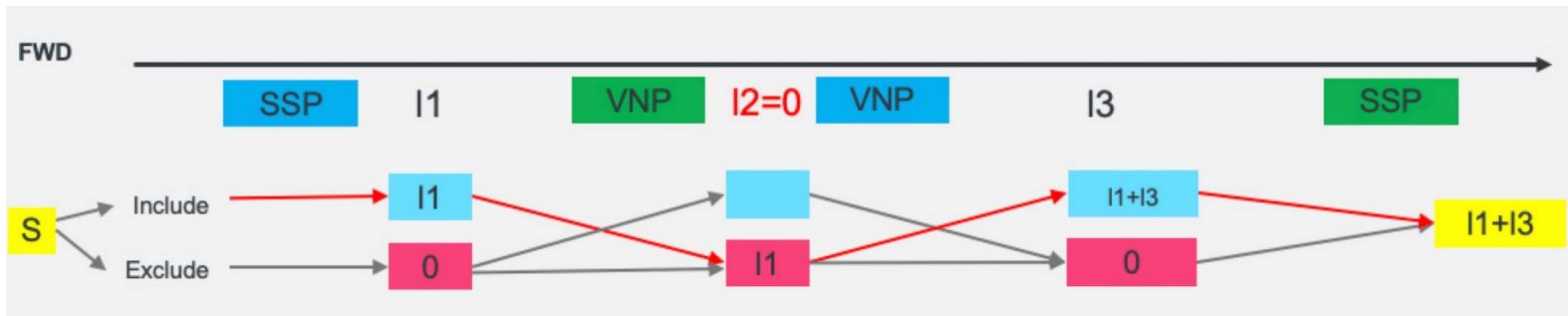
Pychopper



Pychopper



Pychopper

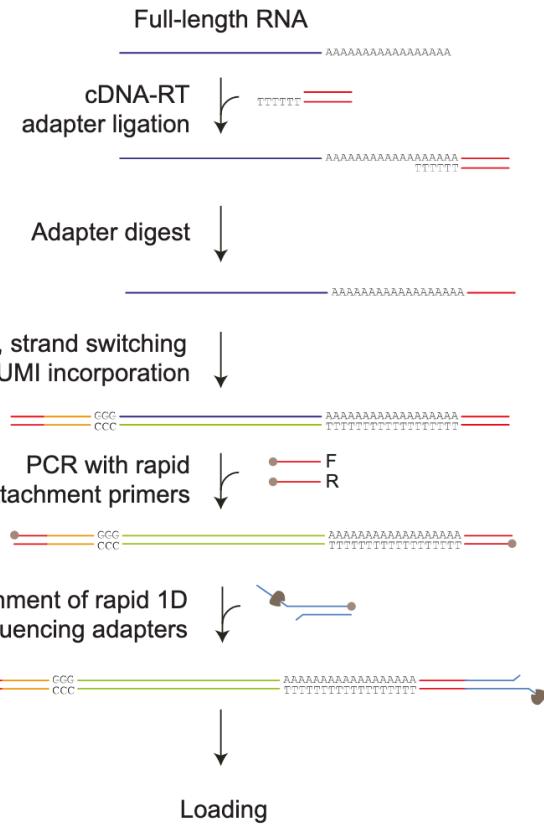


How to estimate duplication rates

- In short read data, reads at the same position are assumed to be duplicates
- How do we know if we're saturating our libraries?



UMIs



UMIs

- UMI at the 3' end of the read

TTT **GGGG**TT **GGAATT**GGCCTT**GGCAT**TT

UMIs

- UMI at the 3' end of the read

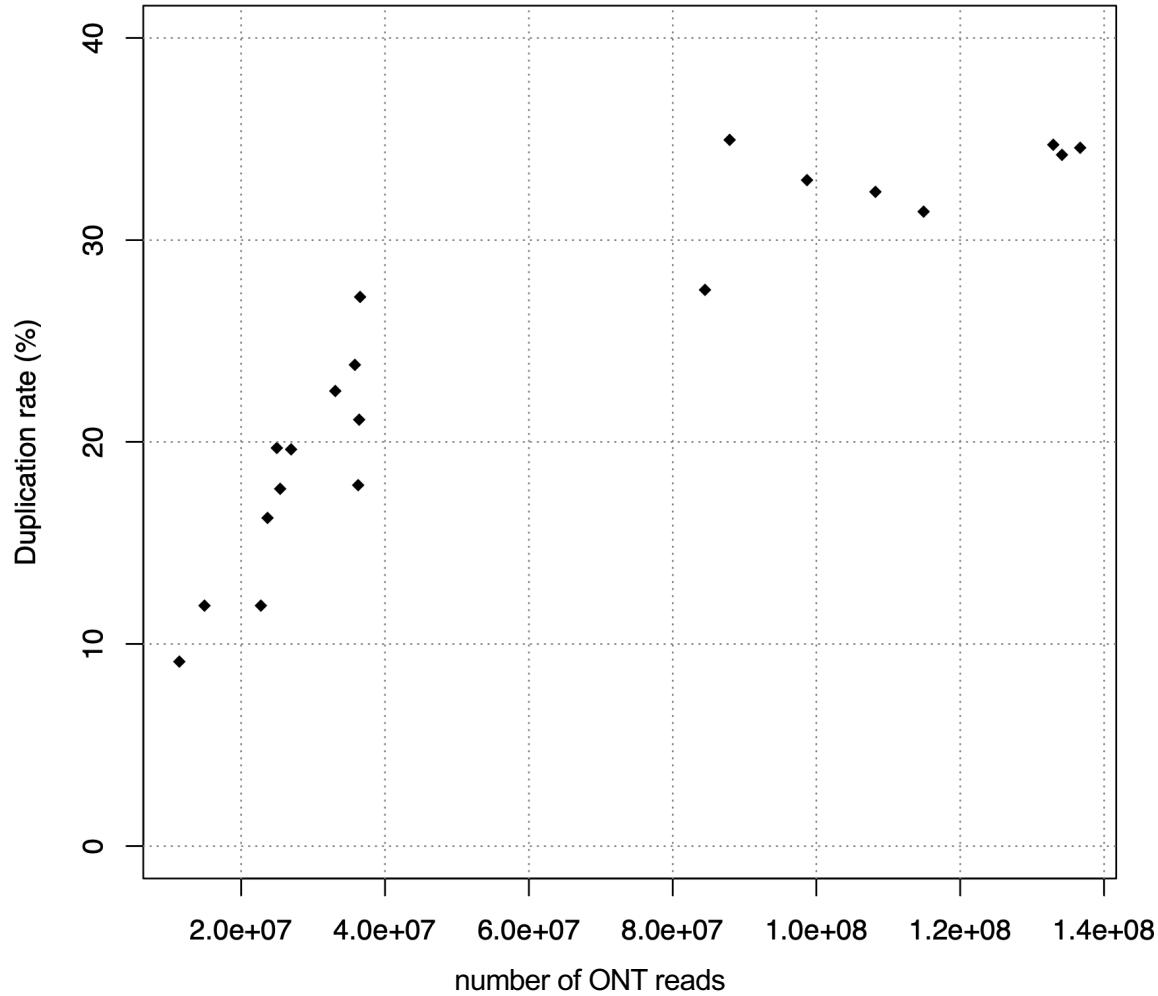
30 TTTCACCCCTCCACTTCCCCTCAGAATT
29 TTTGAAACAGCTTCACCTTGAACCTT
29 TTTCCAATAAAAAAAATTACAATTT
29 TTTCAGCAAAATAAAATTCCGGTTT
27 TTTGGAGTTGGGGTTGCGCTTGGGGTT
27 TTTGAGGTTGGAGTTGGGGTTGGCGTT
24 TTTGGAGTTGGCGTTGCGGTTGGGGTTT
23 TTT~~GGG~~TT~~G~~GAATT~~G~~C~~G~~TT~~G~~C~~A~~TTT
23 TTTGGGATTAAGATTGGCATTGCGGTT
23 TTTAGGGTTCGCGTTGGGGTTGCAGTTT
23 TTTAGGGTTAGCGTTGGAGTTGGGGTT
22 TTTGGCGTTGGGGTTGGCGTTGGCGTT
22 TTTGGCGTTGGAGTTGGGCTTGGCGTT
22 TTTACACTTGTGCTCTCCTTAGCCTTT
21 TTTGGGGTTGGAGTTGGCGTTGGCATTT
21 TTTGGCGTTGGCATTGGCGTTGGGGTTT
21 TTTGGCGTTCGGGTTGGAATTGCGTTT

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

30 TTTCACCCCTCCACTTCCCCTCAGAATT
29 TTTGAAACAGCTTCACCTTGAACCTT
29 TTTCCAATAAAAAAAATTACAATTT
29 TTTCAGCAAAATAAAATTCCGGTTT
27 TTTGGAGTTGGGGTTGCGCTTGGGGTT
27 TTTGAGGTTGGAGTTGGGGTTGGCGTT
24 TTTGGAGTTGGCGTTGCGGGTTGGGGTTT
23 TTT~~GGG~~TT~~G~~GAATT~~G~~GC~~T~~T~~G~~CA~~T~~TT
23 TTTGGGATTAAGATTGGCATTGCGGTT
23 TTTAGGGTTCGCGTTGGGGTTGCAGTTT
23 TTTAGGGTTAGCGTTGGAGTTGGGGTT
22 TTTGGCGTTGGGGTTGGCGTTGGCGTT
22 TTTGGCGTTGGAGTTGGGCTTGGCGTT
22 TTTACACTTGTGCTCTCCTTAGCCTTT
21 TTTGGGGTTGGAGTTGGCGTTGGCATTT
21 TTTGGCGTTGGCATTGGCGTTGGGGTTT
21 TTTGGCGTTGGGTTGGAATTCGCGTT

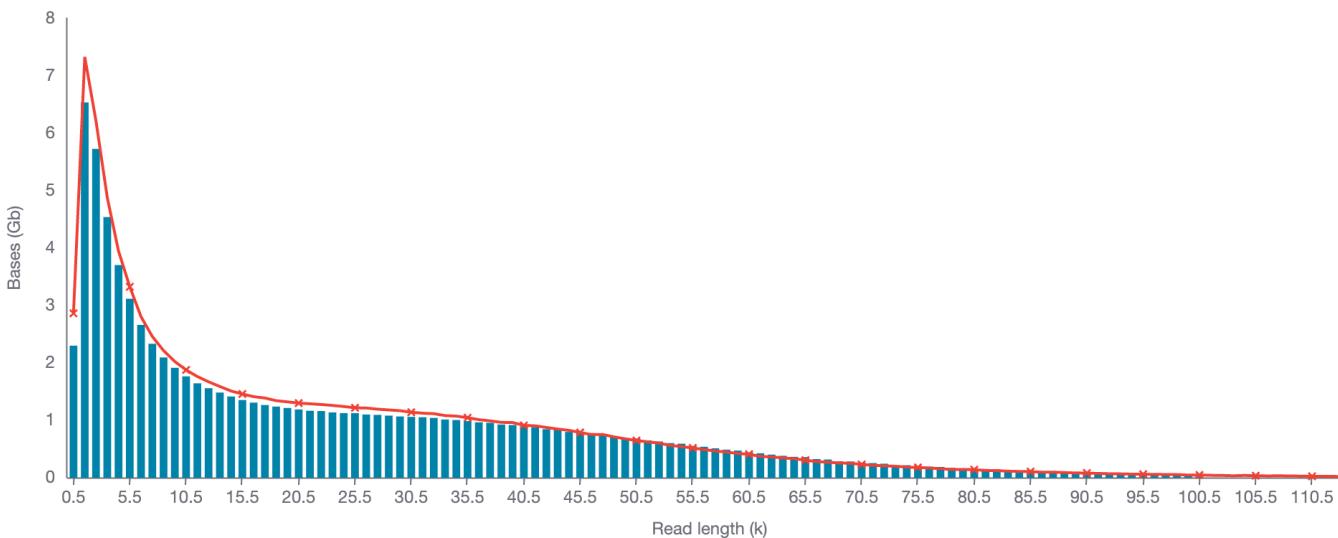
$$0.95^{28} = 0.237$$



What does the data look like?

Legend

Basecalled Estimated



Genomic DNA – standard prep

Estimated N50

17.75 kb

% Basecalled

100%

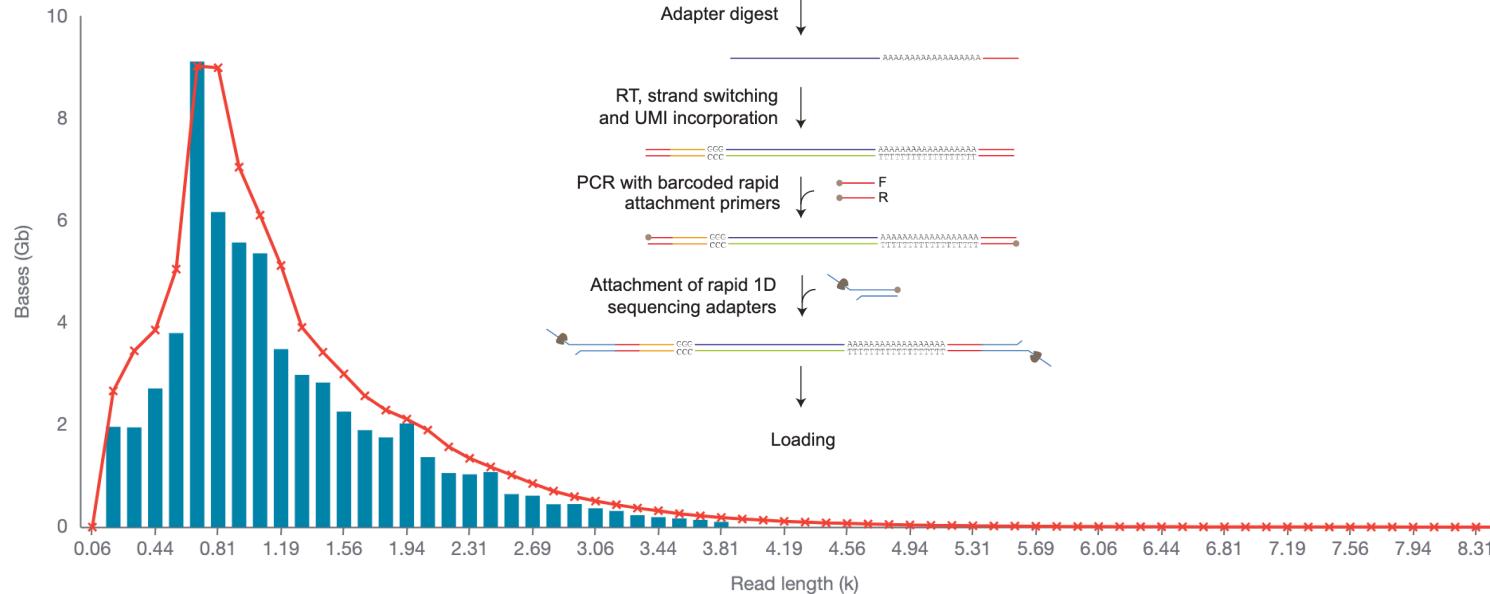
their relative amounts.

Read length (kb)	Aggregated reads (Mb)
100 - 164	886.98
164 - 228	36.06
228 - 292	4.02
292 - 344	0.35

What does the data look like?

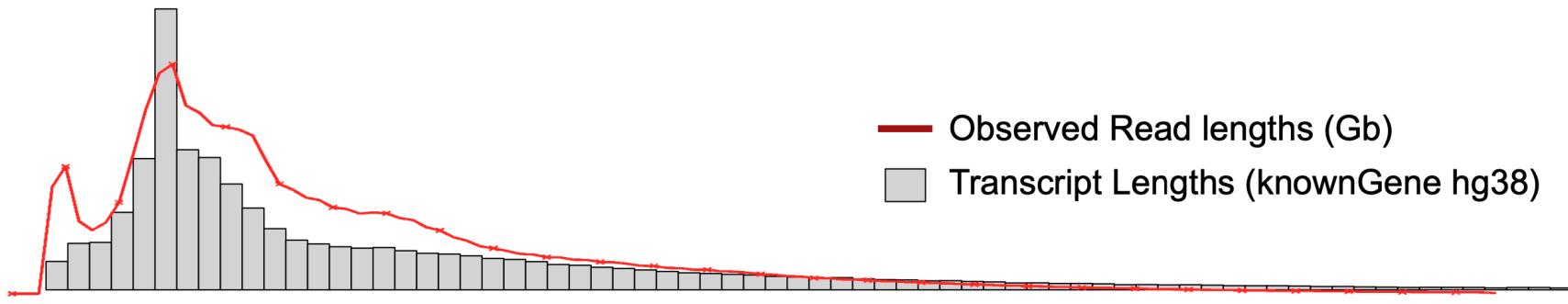
Legend

Basecalled Estimated

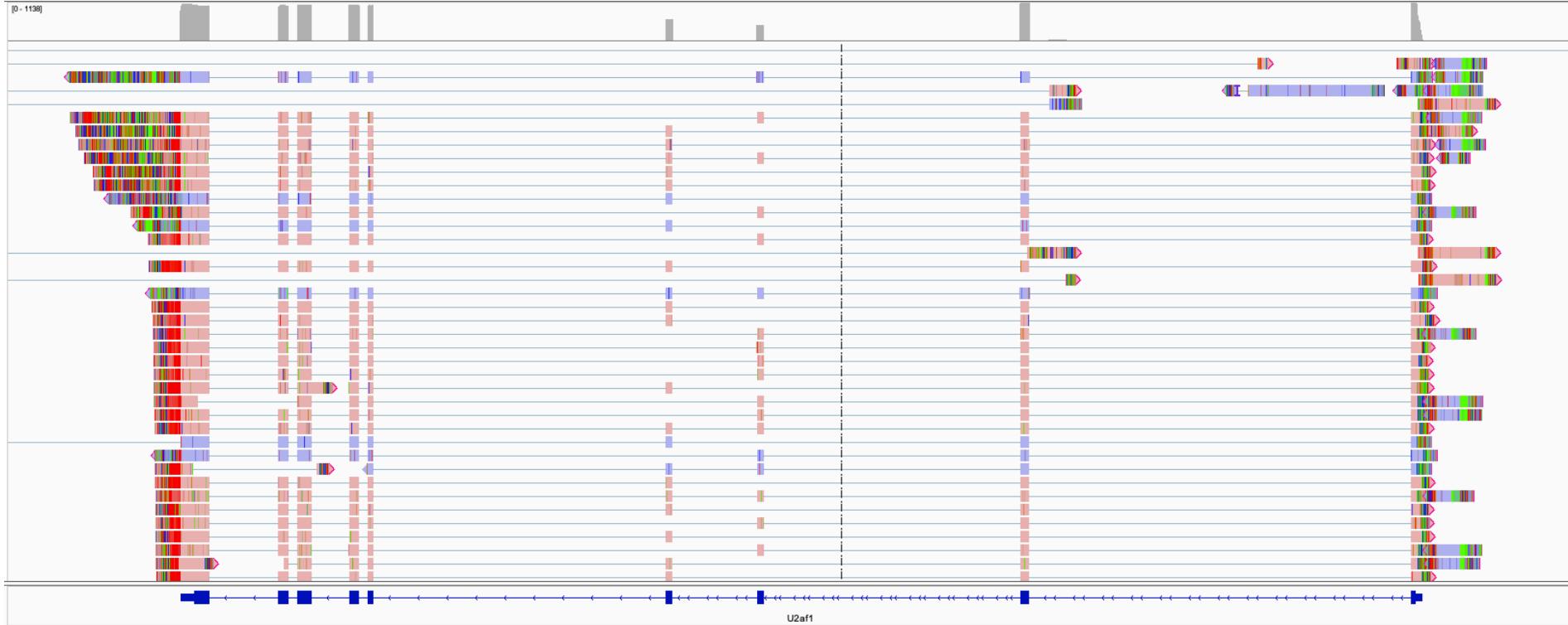


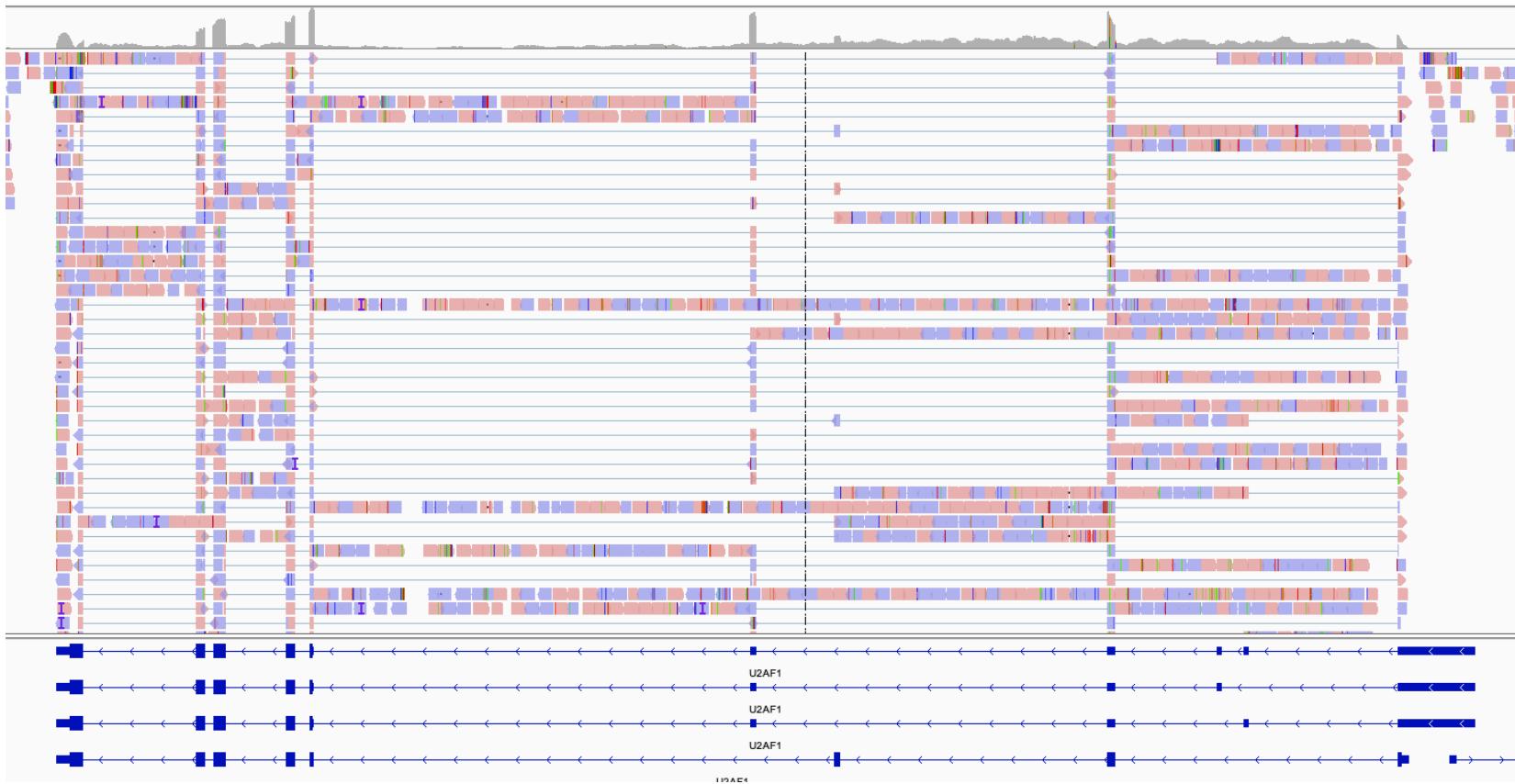
RNA/cDNA – standard prep

What does the data look like?



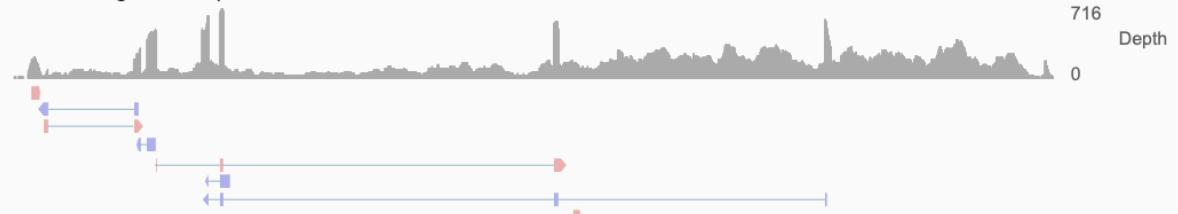
cDNA – standard prep





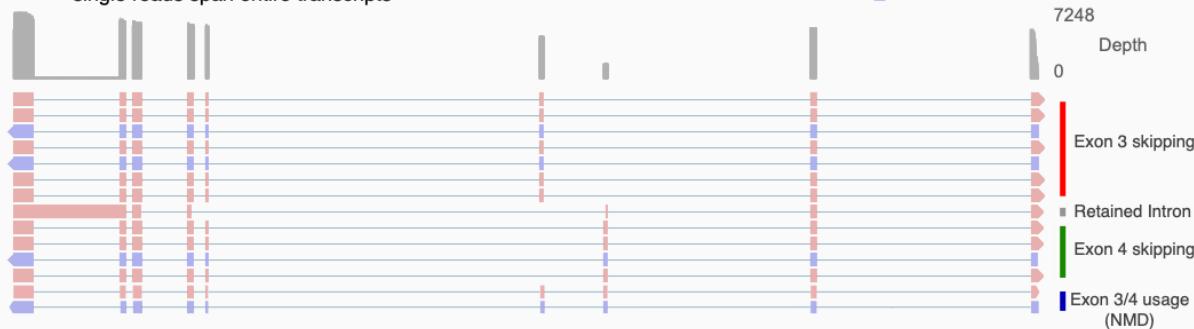
Short-read coverage (65M reads)

single reads span 1-3 exons



Long-read coverage (48M reads)

single reads span entire transcripts



U2AF1 - ENST00000291552

U2AF1 - ENST00000380276

U2AF1 - ENST00000464750

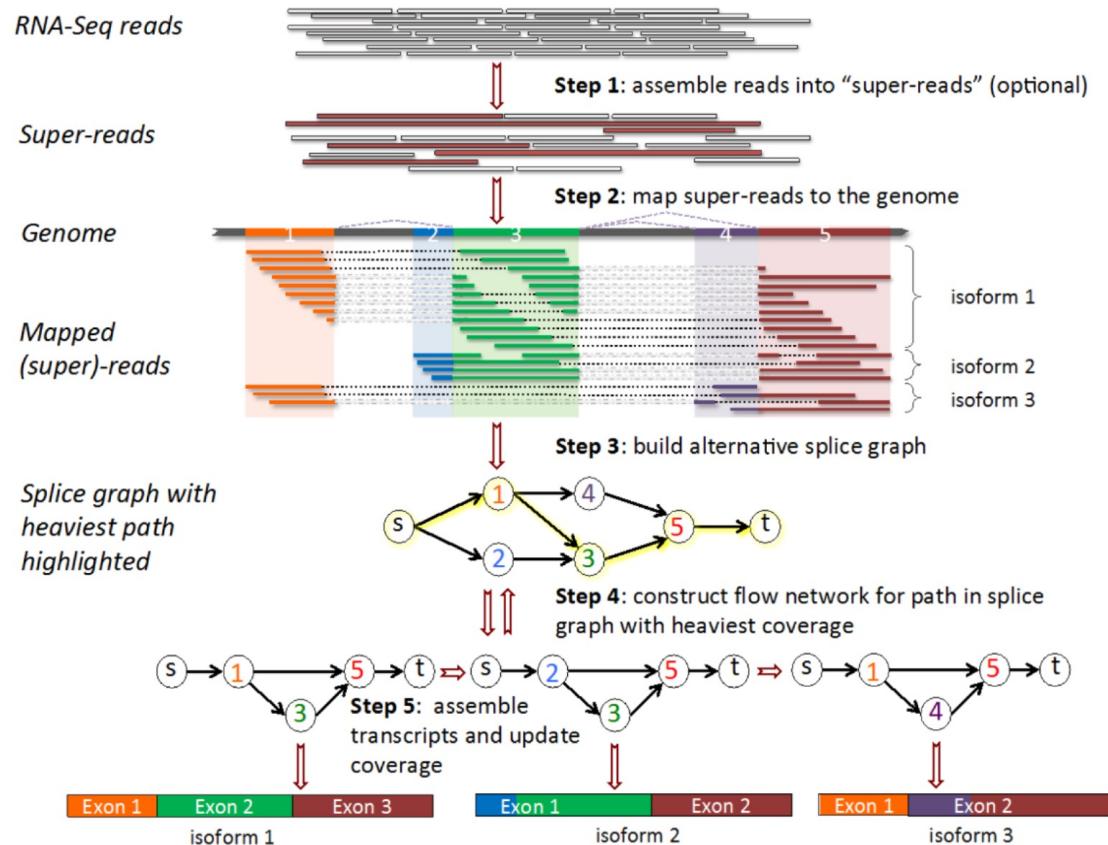
chr21

43,095,000

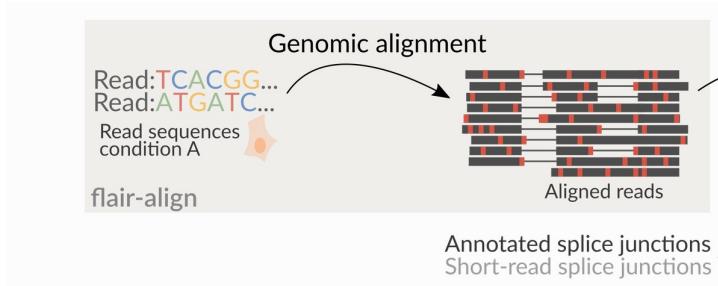
43,100,100

43,105,000

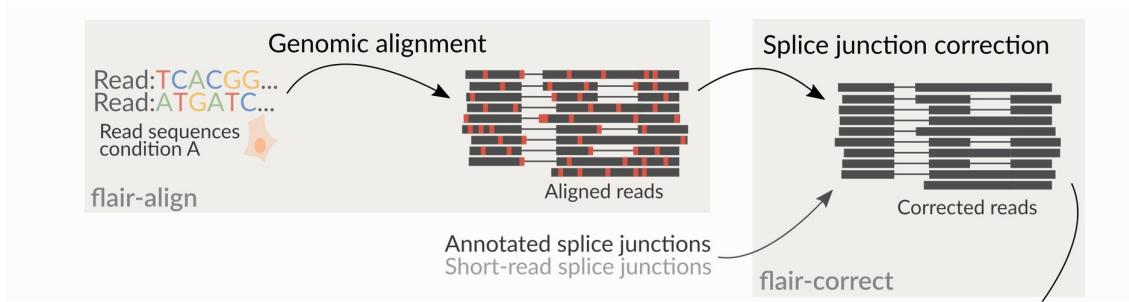
Estimating transcript abundance – short-read, Stringtie



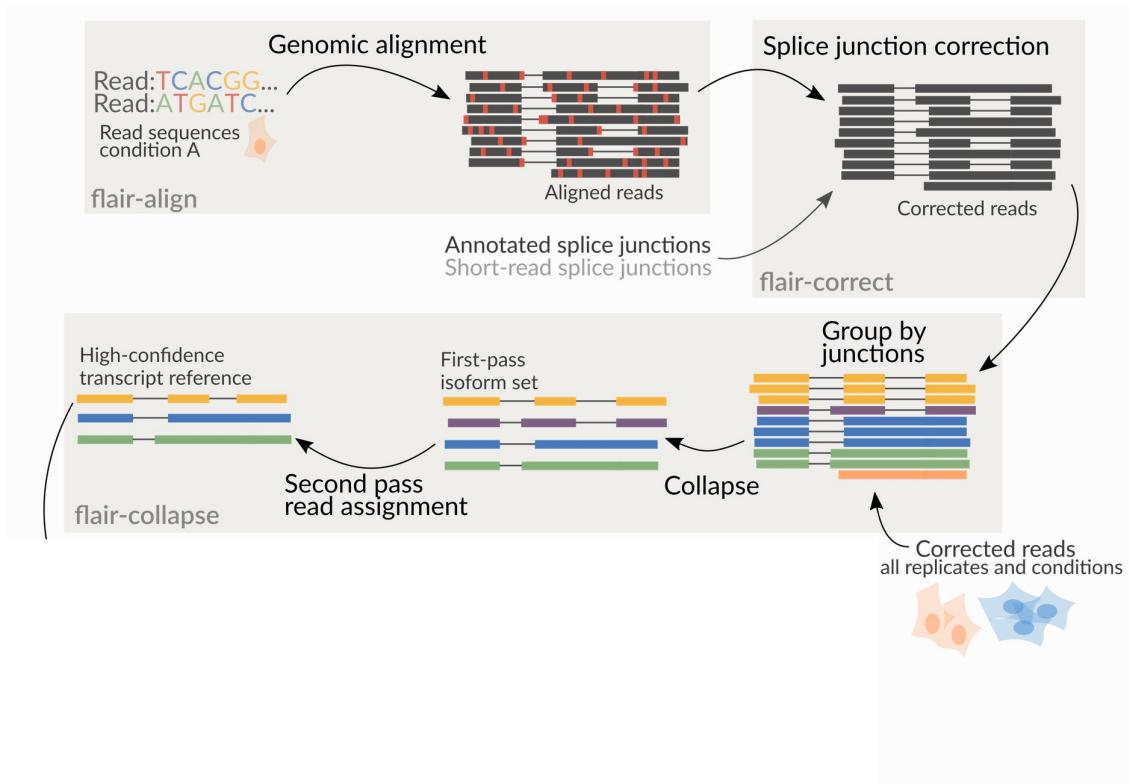
Estimating transcript abundance – long read



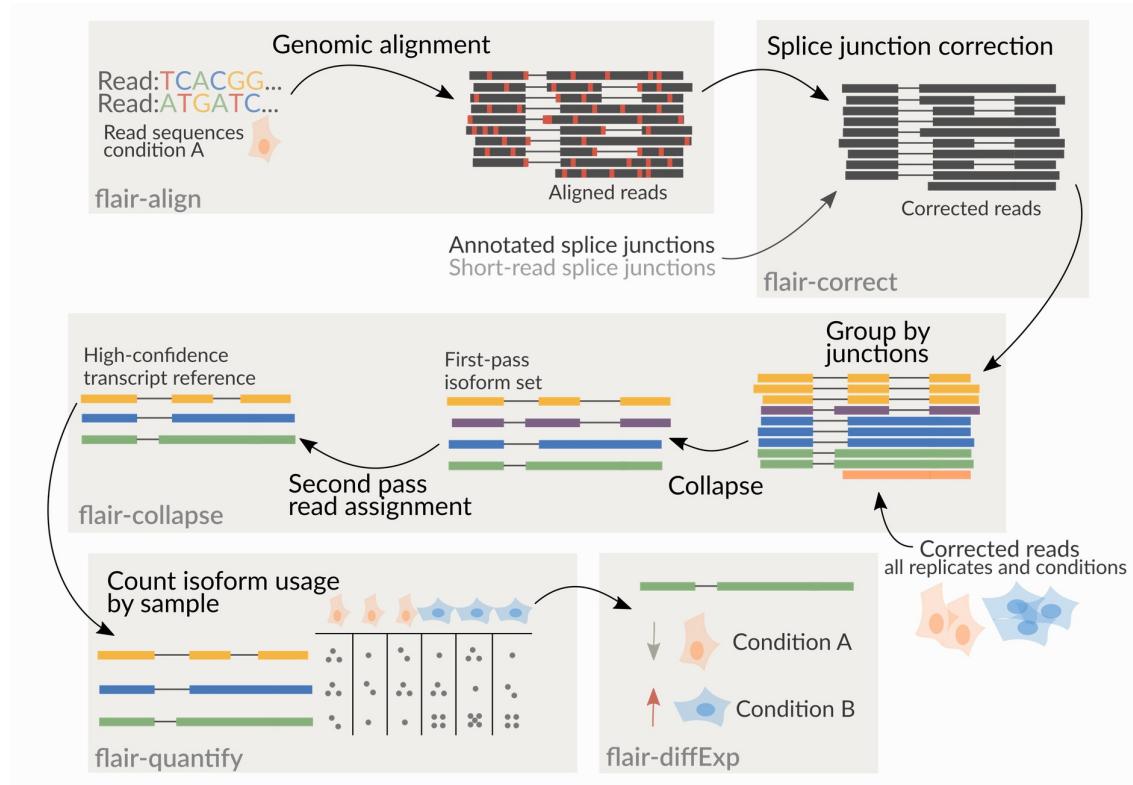
Estimating transcript abundance – long read



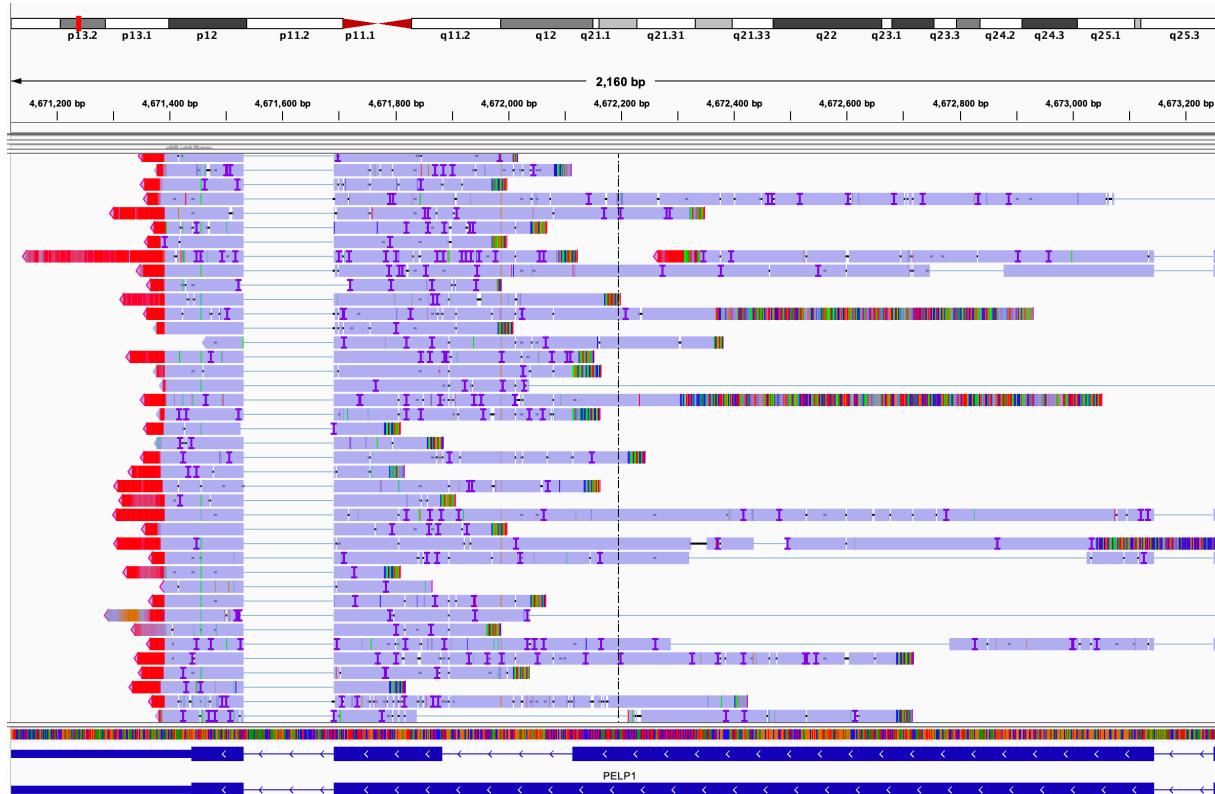
Estimating transcript abundance – long read



Estimating transcript abundance – long read



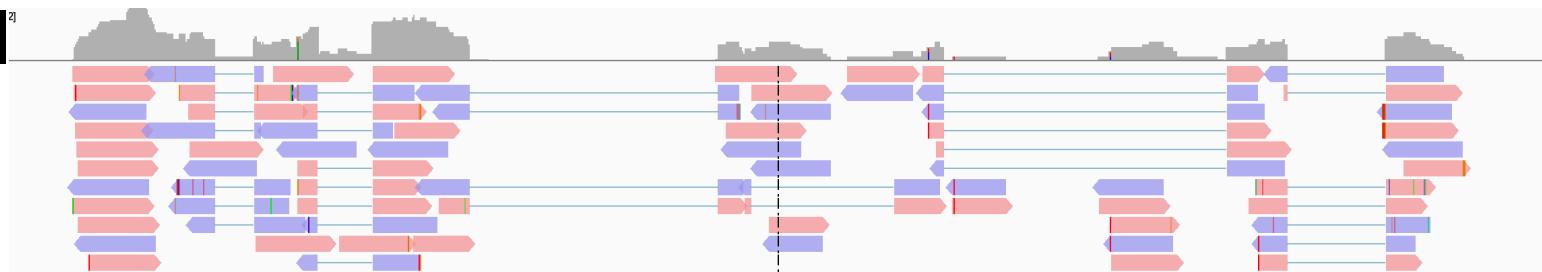
Technical artifacts



Looks like fragmentation of this RNA throughout the long exon. (This is an egregious case – most are more mild)

"Ful"

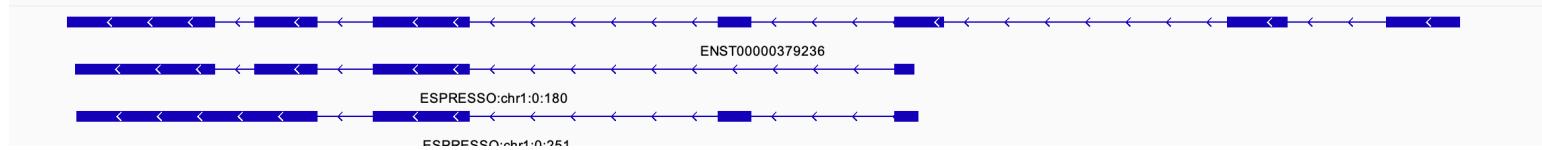
Short-read



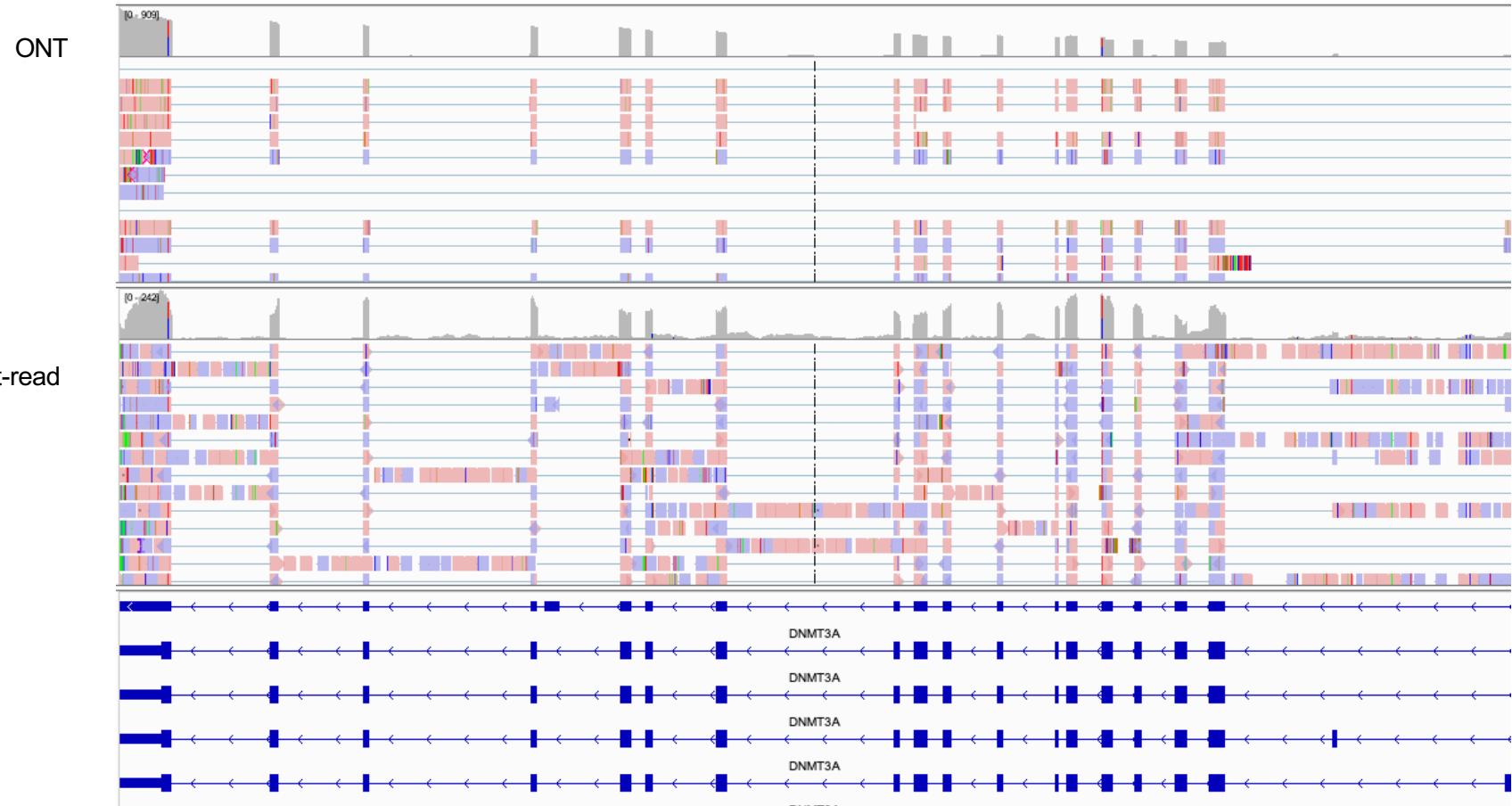
Long-read



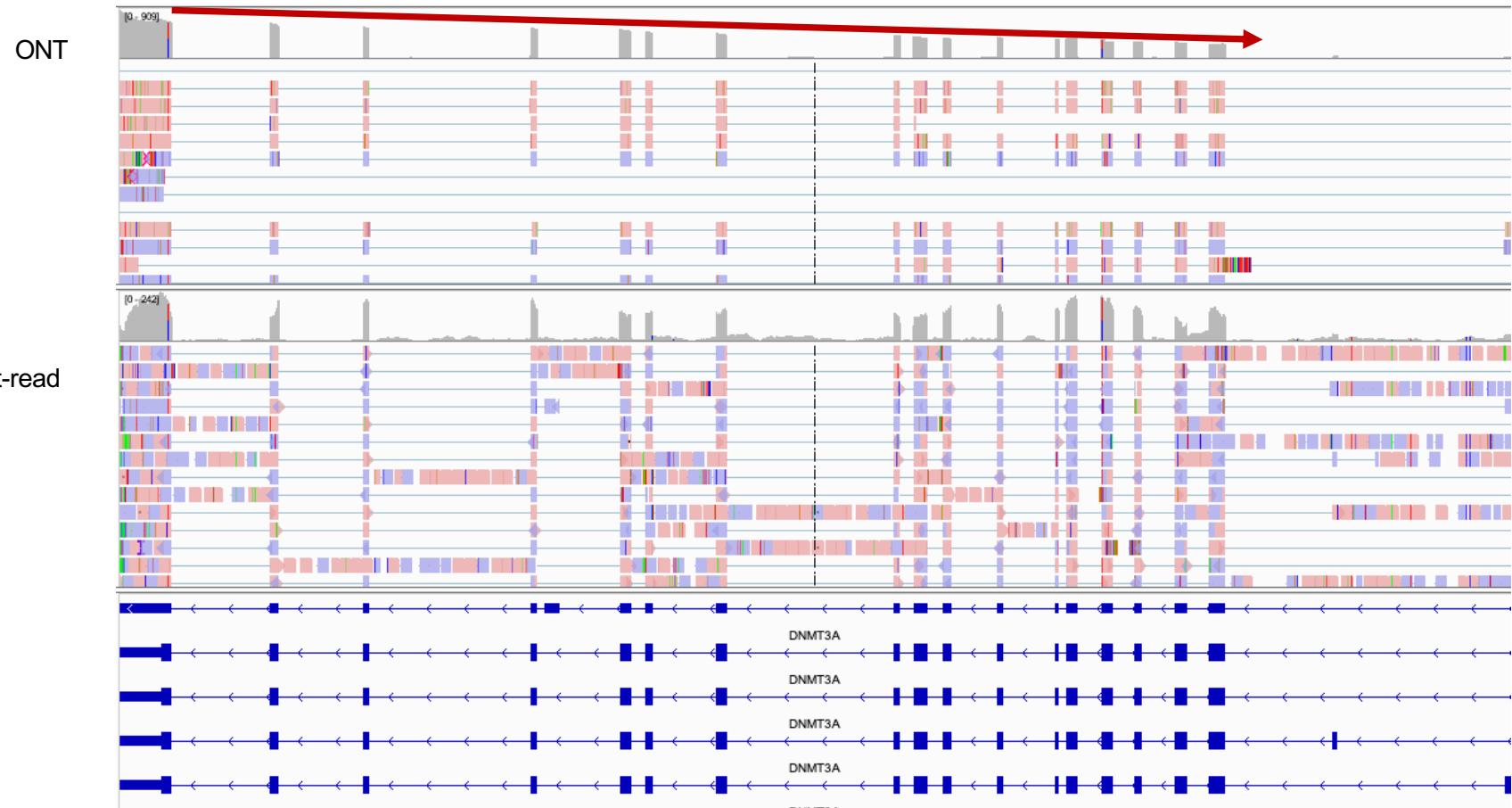
ESPRESSO
transcript
assembly



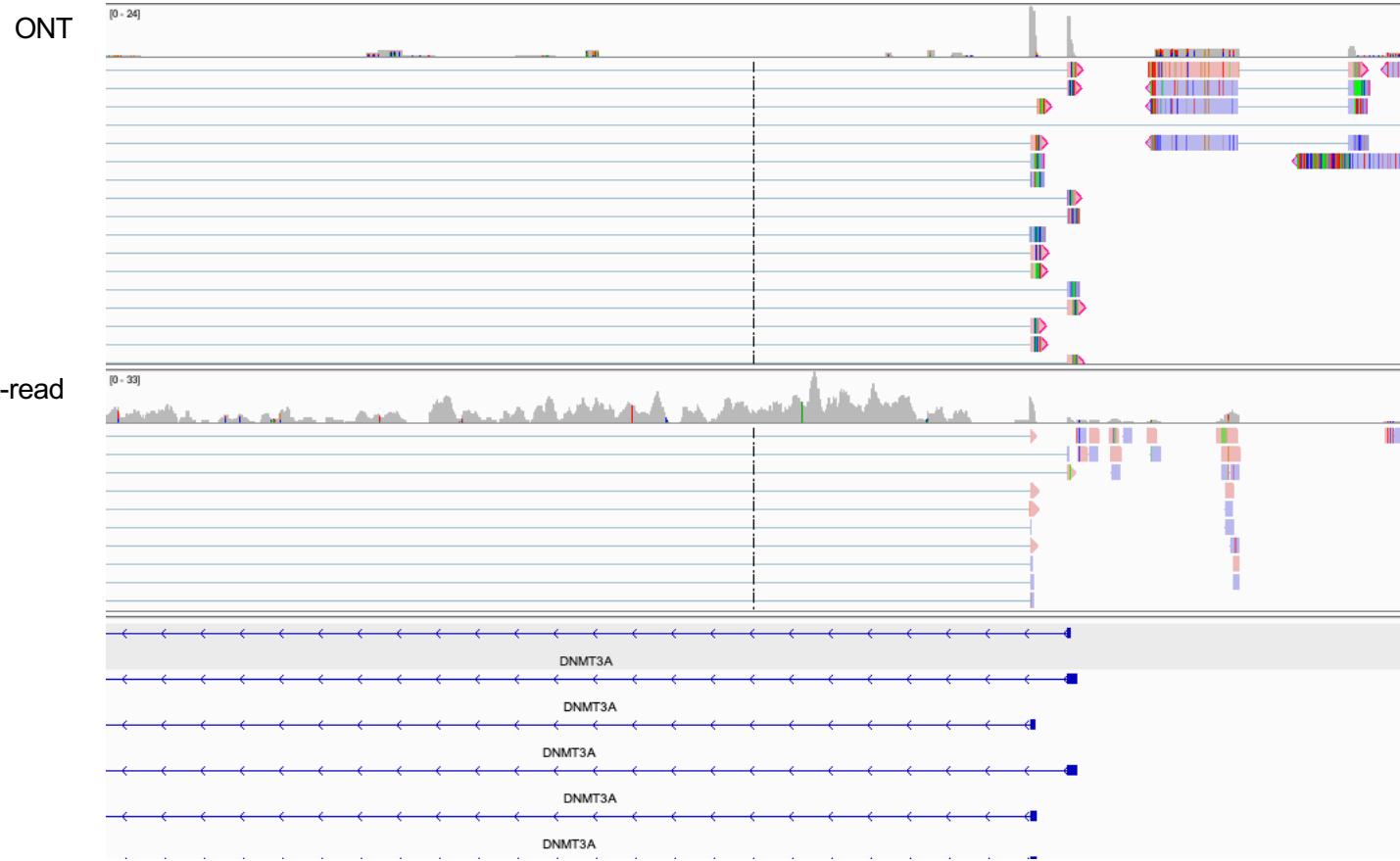
human DNMT3A



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