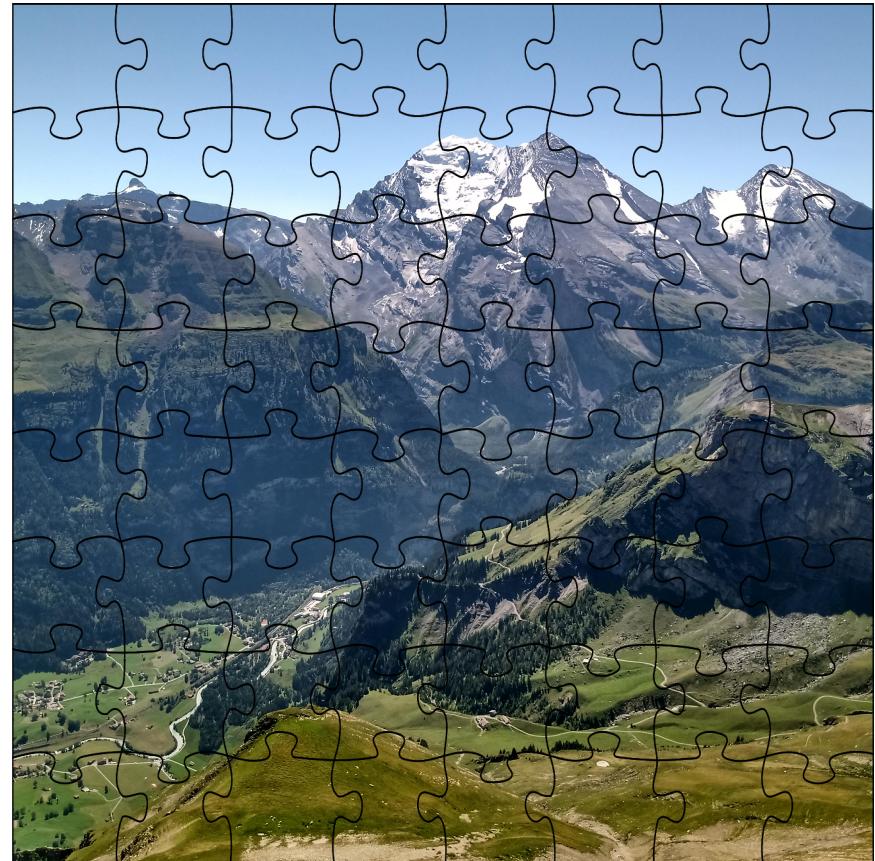
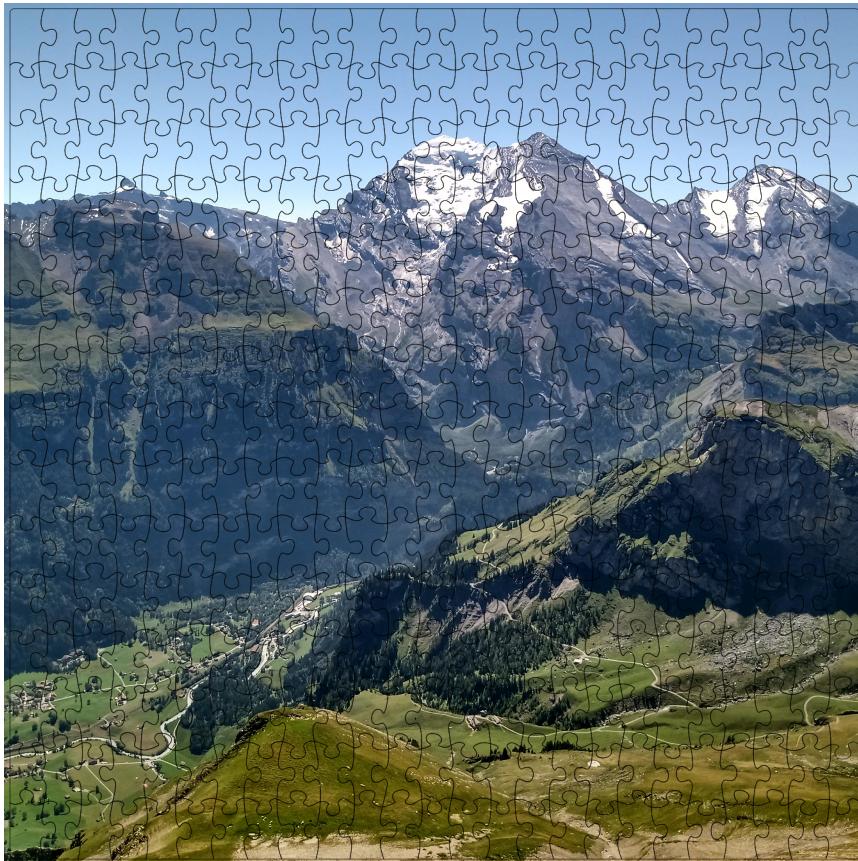


Long-read RNA-seq

What is a long read?

- Short read: 50-300 bp, often paired-end
- Long read: > 1kb, up to 20 Mb

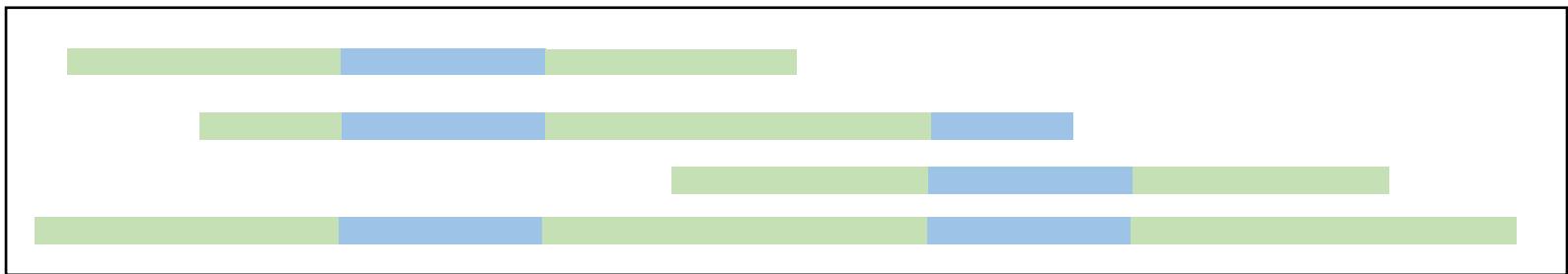
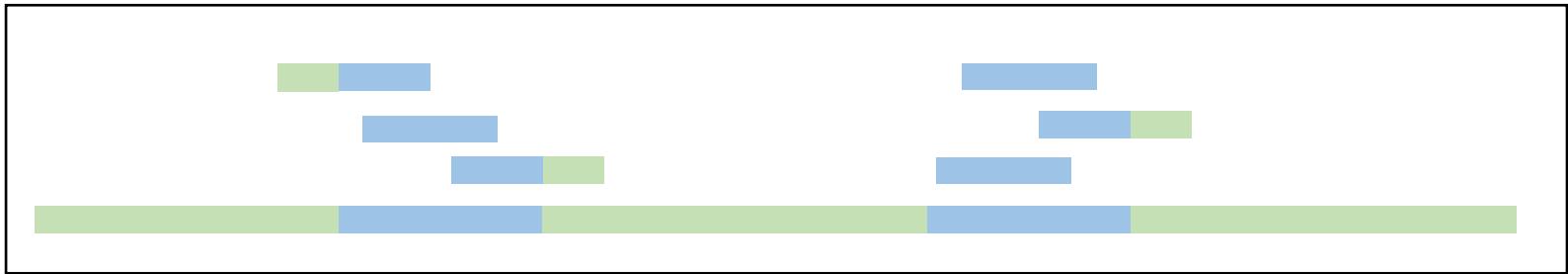
Why long reads?



Why long reads?

- Many 2nd generation data analyses have an ‘assembly’ or ‘uniqueness’ question:
 - Mapping quality?
 - Isoform variation?
 - Structural variation?
 - Haplotypes?
 - Genome/transcriptome sequence?

Mapping quality

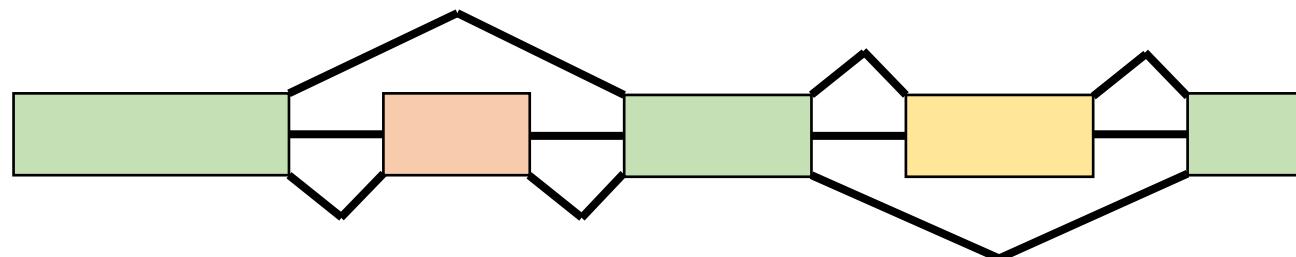


$$\begin{aligned} MAPQ \\ = -10 \log_{10} \Pr\{\text{mapping position is wrong}\} \end{aligned}$$

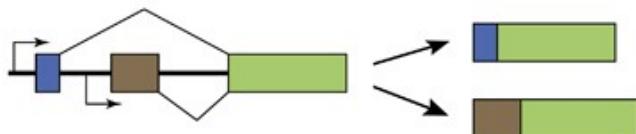
$$\begin{aligned} -10 \log_{10} \Pr\{0.01\} &= 20 \\ -10 \log_{10} \Pr\{0.5\} &= 3 \end{aligned}$$

Isoform variation

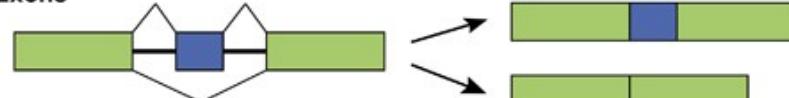
- > 95% of human genes are alternatively spliced
- Some with large phenotypic effects
- Isoforms can have different:
 - Amino acid sequence
 - Promoters
 - Polyadenylation



Alternative Promoters



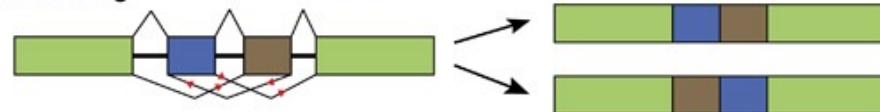
Cassette Exons



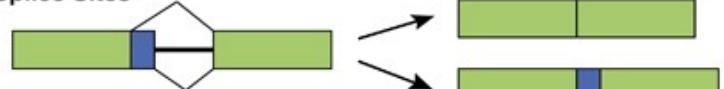
Mutually Exclusive Exons



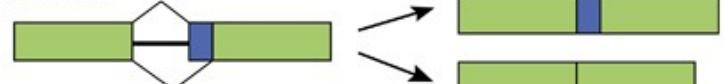
Exon Scrambling



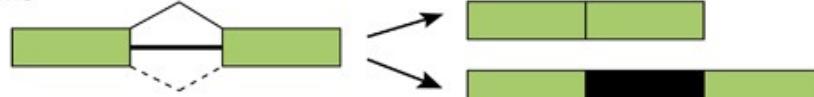
Alternative 5' Splice Sites



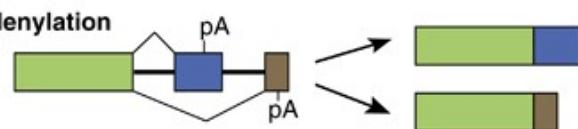
Alternative 3' Slice Sites



Retained Introns

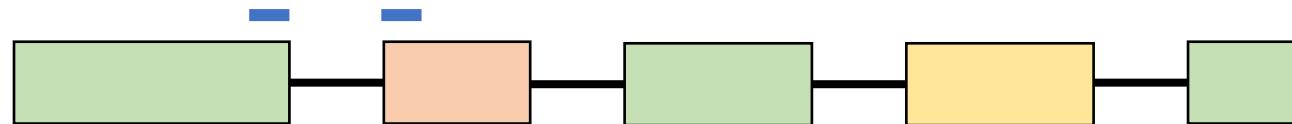
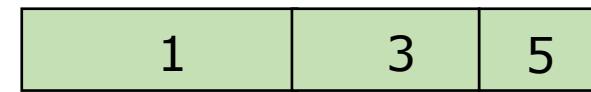
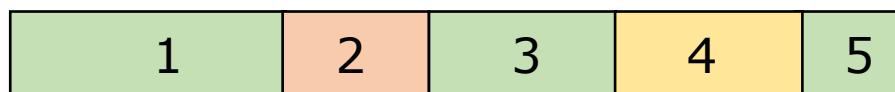
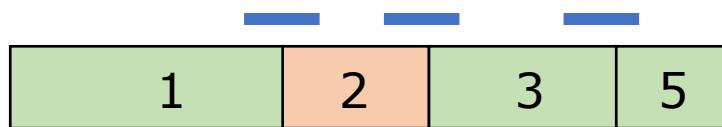
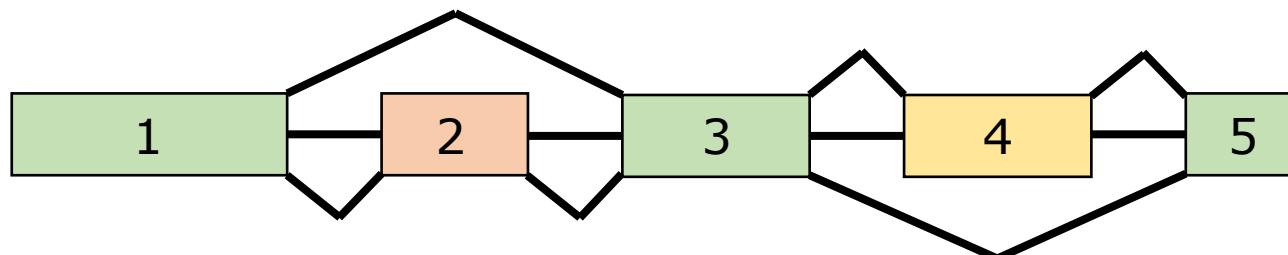


Alternative Polyadenylation



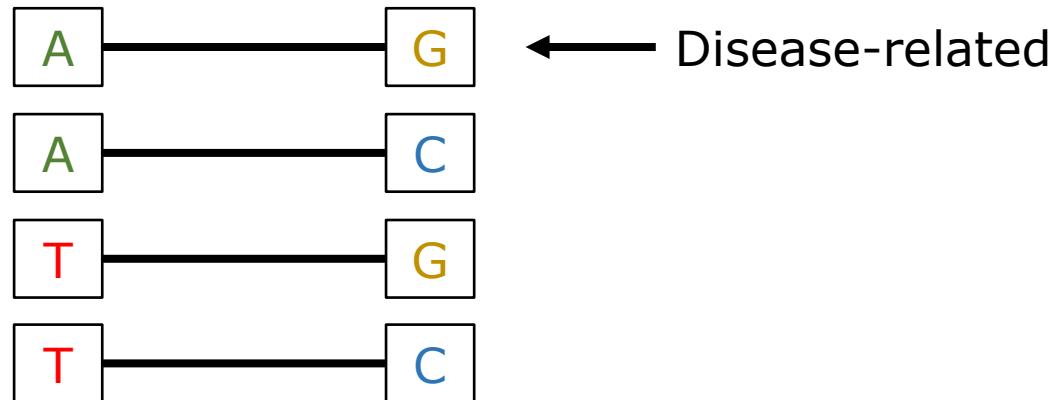
Chen, J., & Weiss, W. A. (2015).
Oncogene, 34(1), 1–14.

Isoform variation

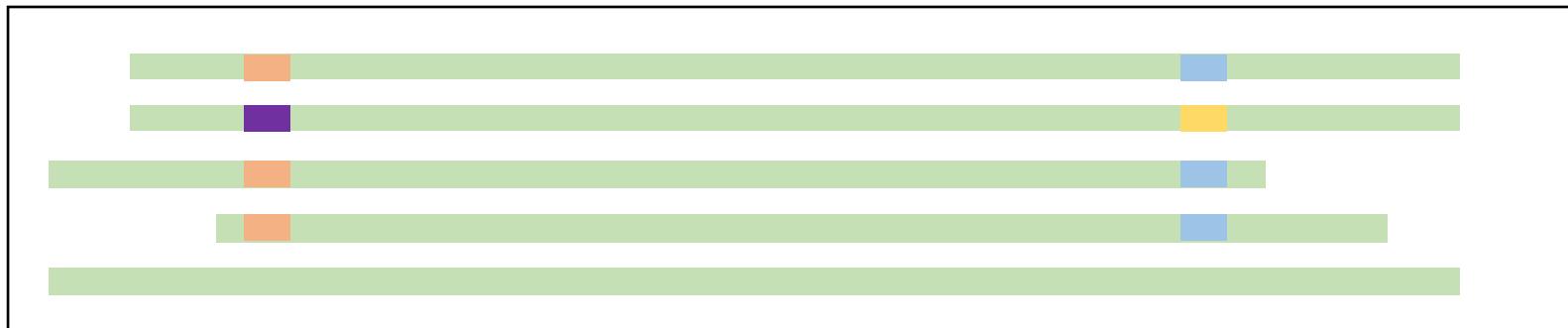
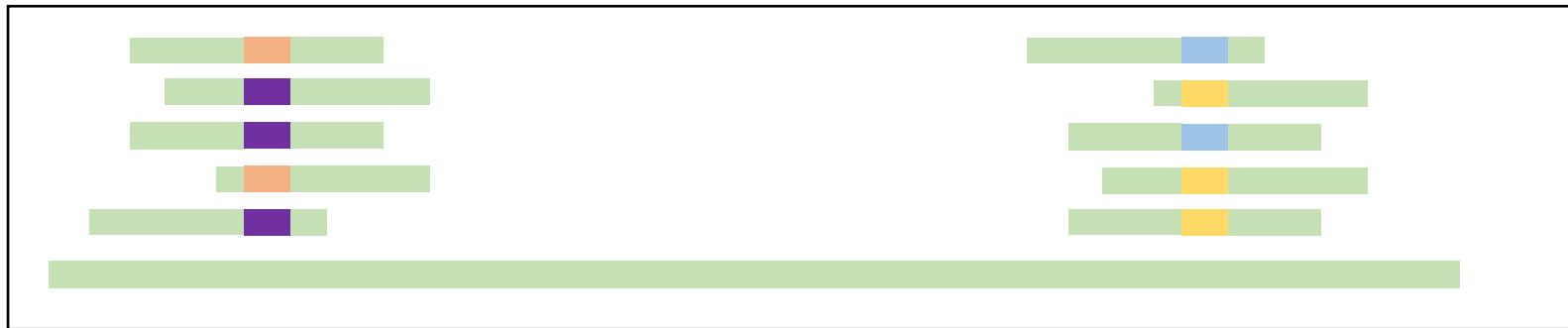


Haplotypes

- Genetic association studies often based on SNPs
- Most SNPs are bi-allelic e.g. [A/T], [G/C]
- Genetic variation is often multi-allelic



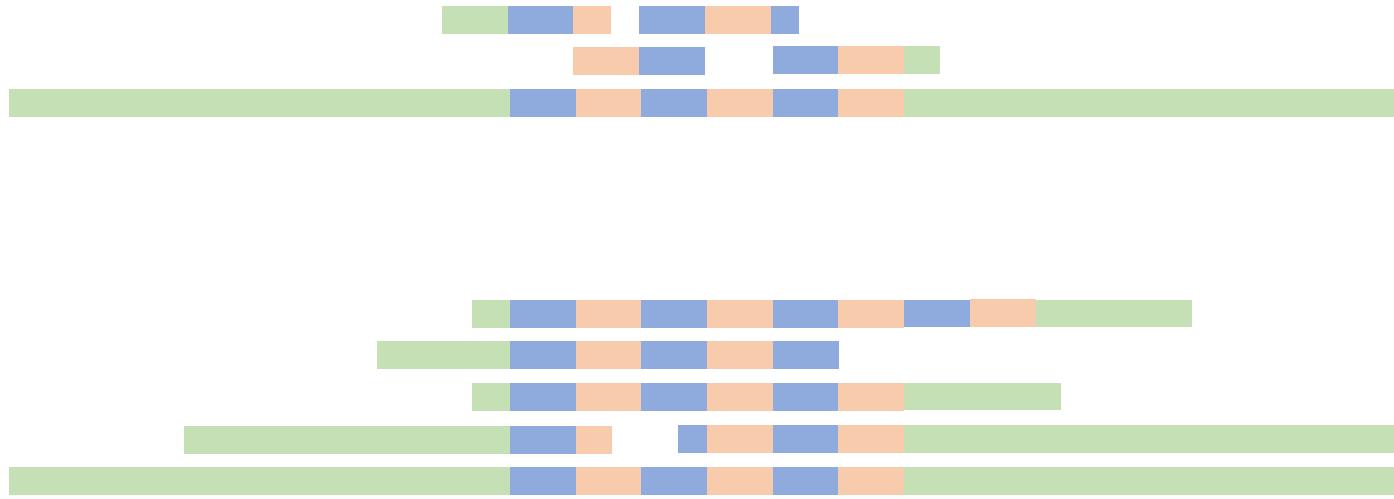
Haplotypes



Structural variation

- ~1kb – 3 Mb
- Associated with important traits
- Insertions/deletions
- Copy-number variation
- Translocations
- Inversions

Copy number variation



Why long-read RNA-seq?

- Genome annotation
- Transcriptome assembly
- Transcript & isoform discovery
- Differential isoform expression
- Poly(A)tail length (dRNAseq)
- RNA base modifications (dRNAseq)
- Transcript haplotypes
- Any other?

Sequencing (c)DNA



1st generation

Sanger

300-1000 bp

High accuracy



2nd generation

Illumina

2x 50-300 bp

High accuracy

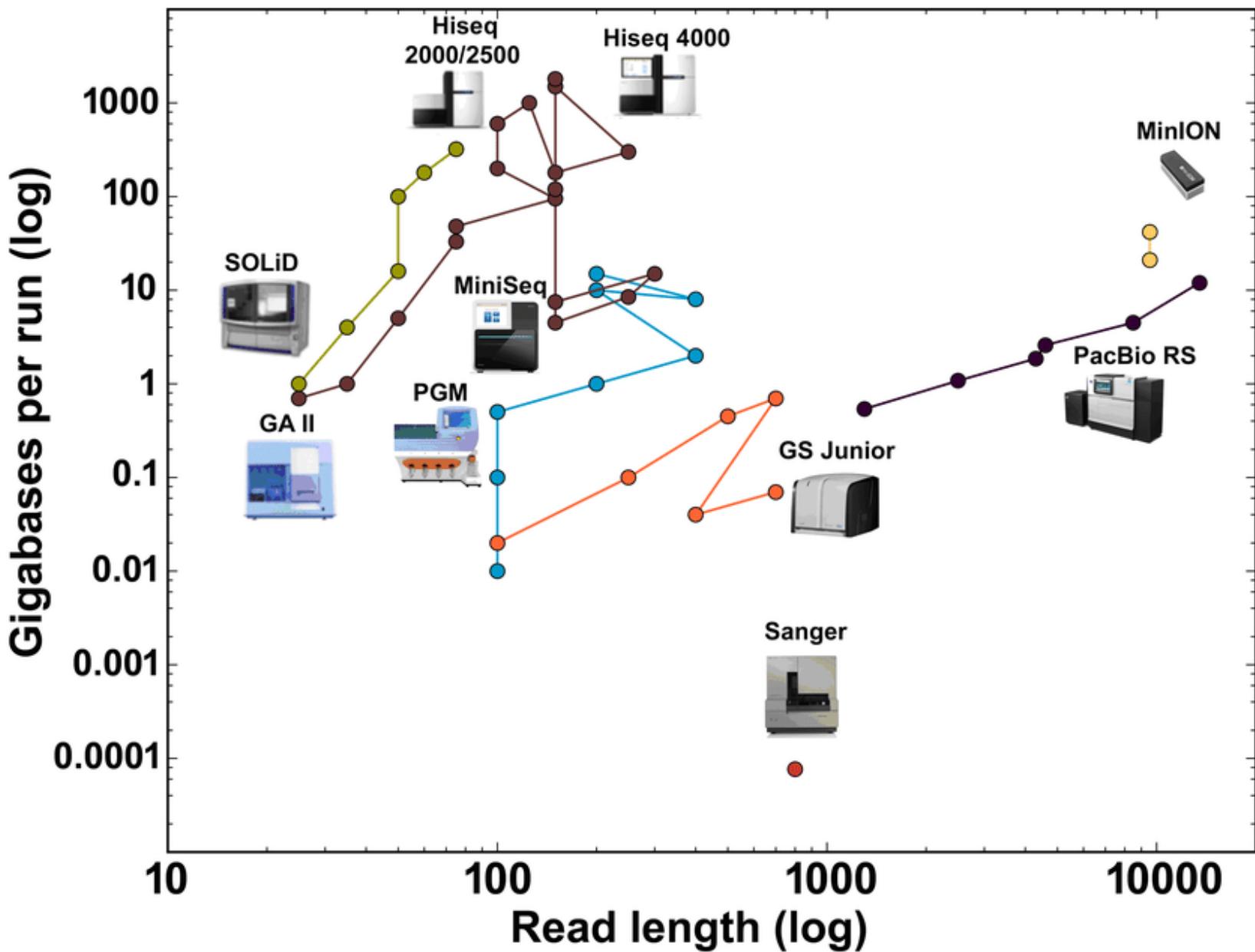


3rd generation

PacBio/ONT

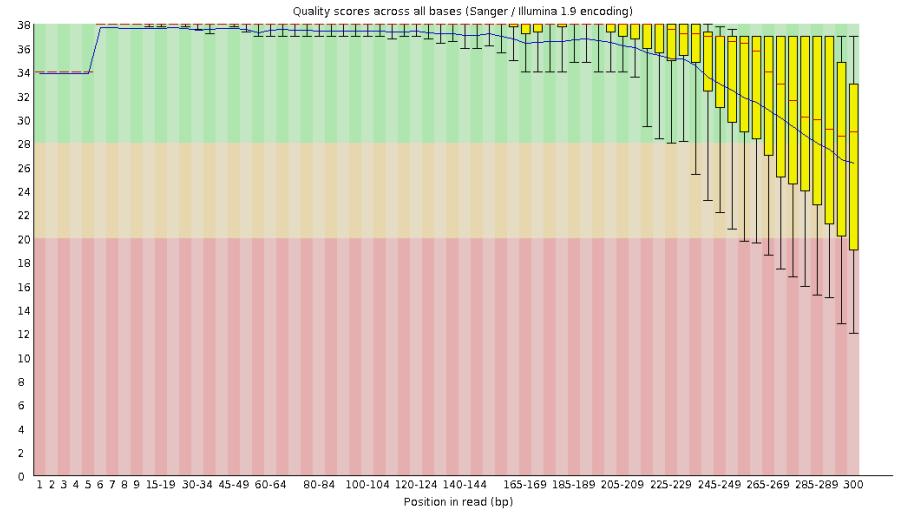
10->20 kb

Medium-High accuracy



Second generation

- Bridge amplification
- Lengths are limited by out-of-phase of signal



Third generation

- Other than 2nd generation:
 - single molecule -> No bridge amplification
- No out-of-phase of signal
- Virtually no limit in read length
- Most used platforms:
 - PacBio SMRT sequencing
 - Oxford nanopore technology



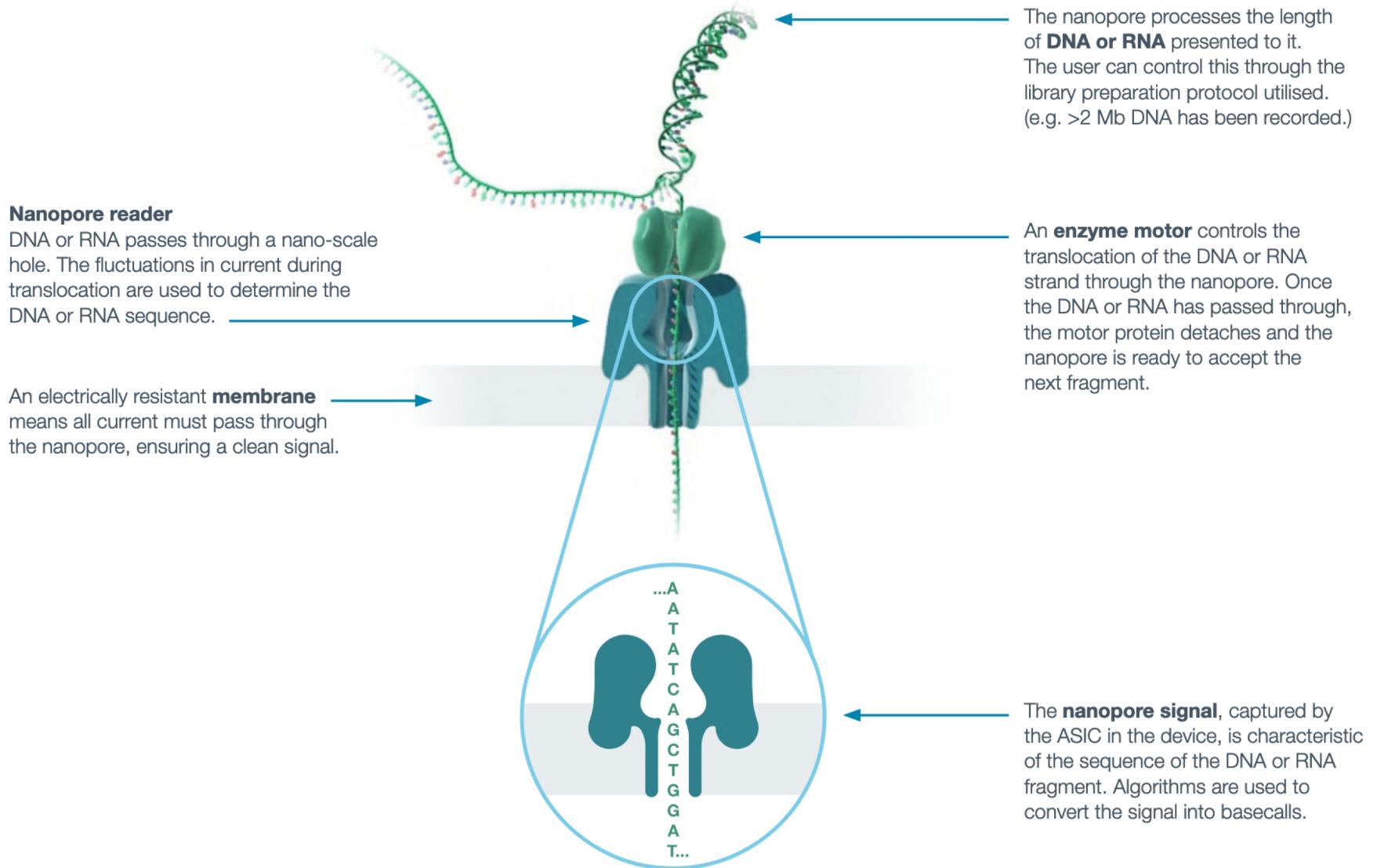
PACBIO®



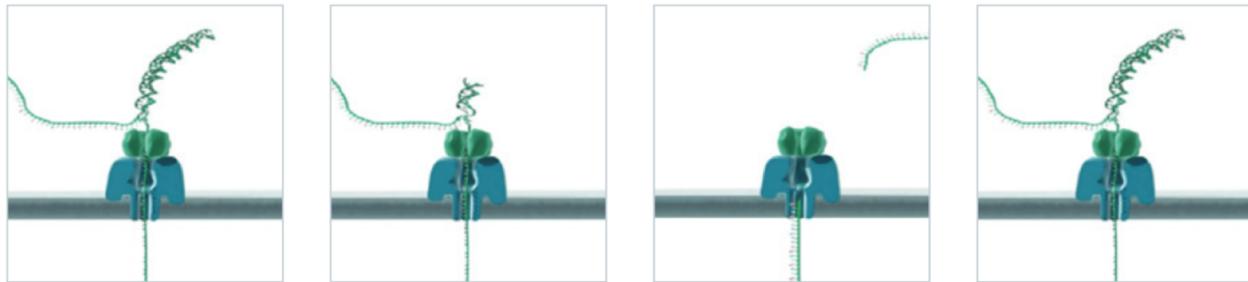
Oxford Nanopore technology

- Based on changes in electrical current
- Well-known for its scalability and portability
- ~95-97% accuracy





1D



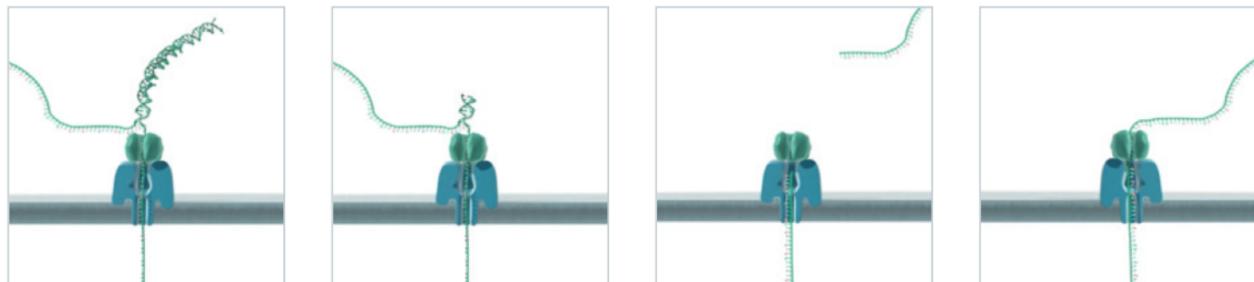
Template...

...Template...

(Exit)

Next molecule...

1D²



Template...

...Template...

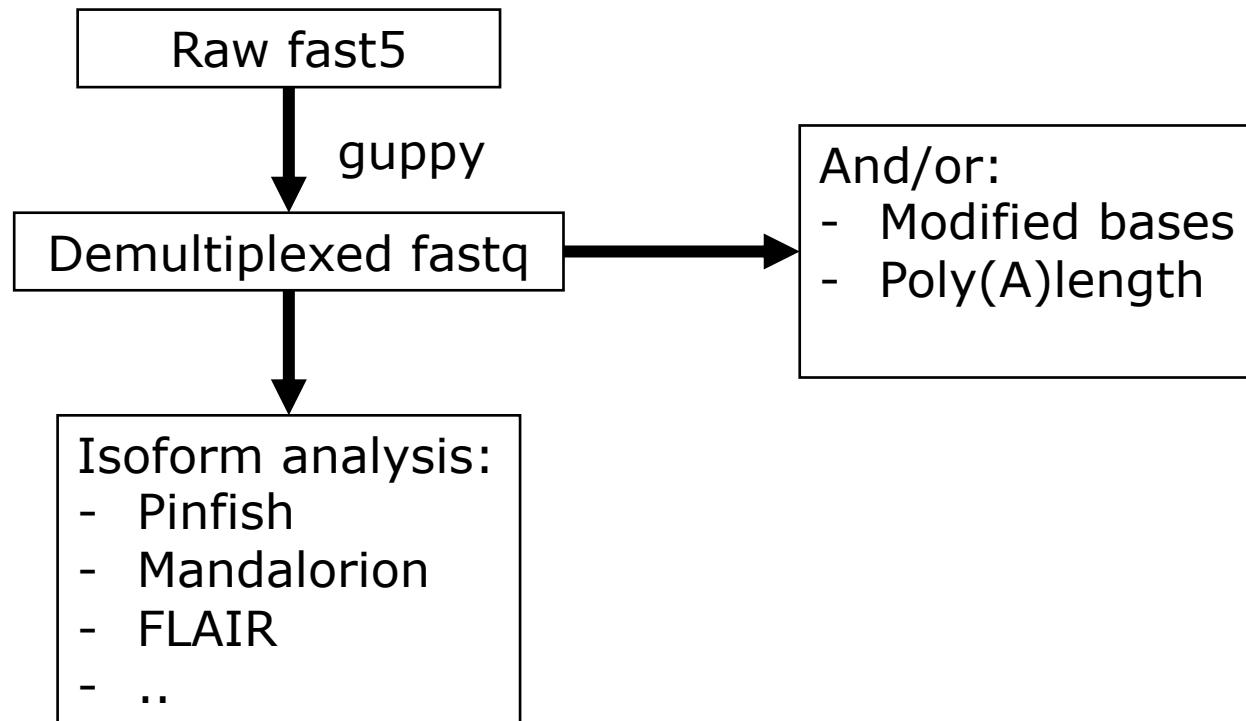
(Exit)

...Complement

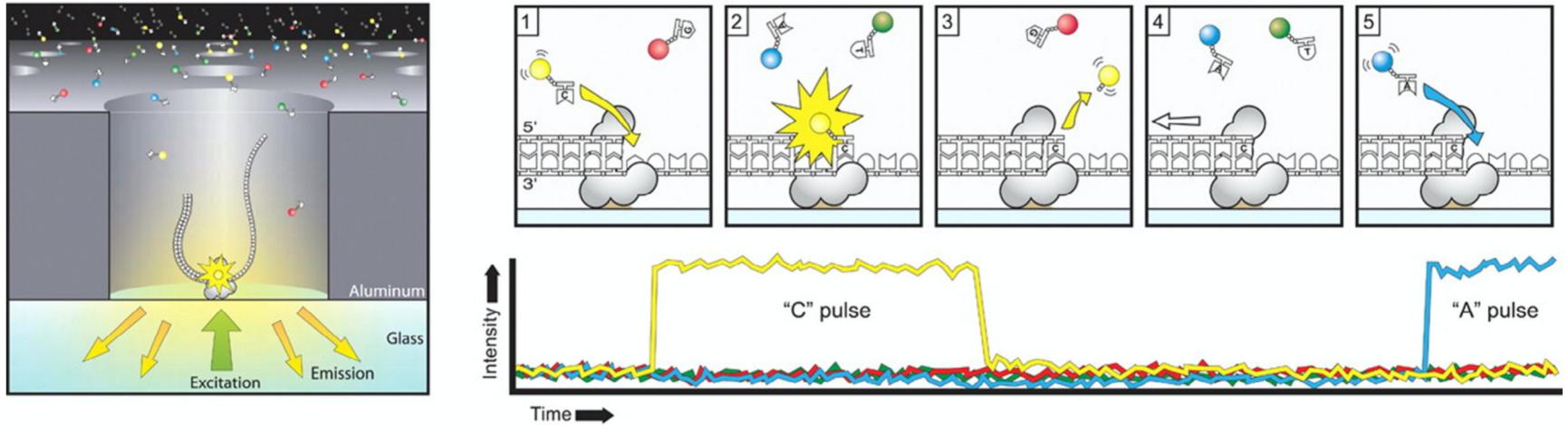
ONT RNA library prep

- cDNA PCR
 - Multiplexing
 - Low input
- Direct cDNA
 - Multiplexing
 - No PCR bias
- Direct RNA (dRNAseq)
 - No multiplexing
 - Detect base modifications
 - Polyadenylation

‘Standard’ transcript workflow



PacBio sequencing



- Polymerase bound to ZMW bottom
- Circular molecules
- Single read out ~90% accuracy
- CCS (HiFi): single molecule sequenced multiple times

Start with high-quality double stranded DNA



Ligate SMRTbell adapters and size select



Anneal primers and bind DNA polymerase



Circularized DNA is sequenced in a single pass



The polymerase reads are trimmed of adapters to yield subread



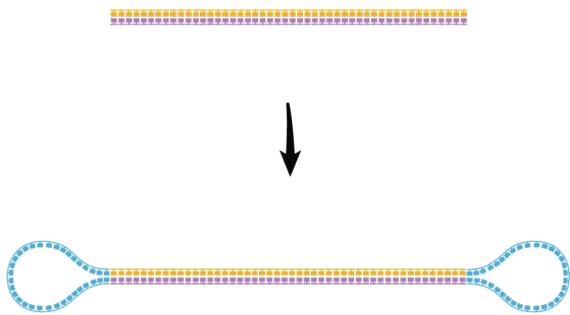
During assembly, consensus is called from multiple molecules

LONG READ

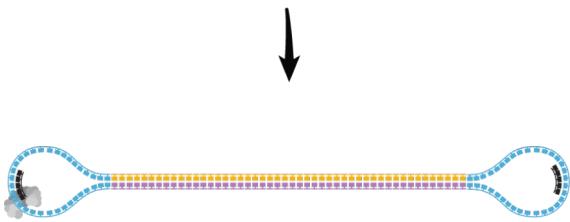
(Half of Reads >50 kb)



Start with high-quality double stranded DNA



Ligate SMRTbell adapters and size select

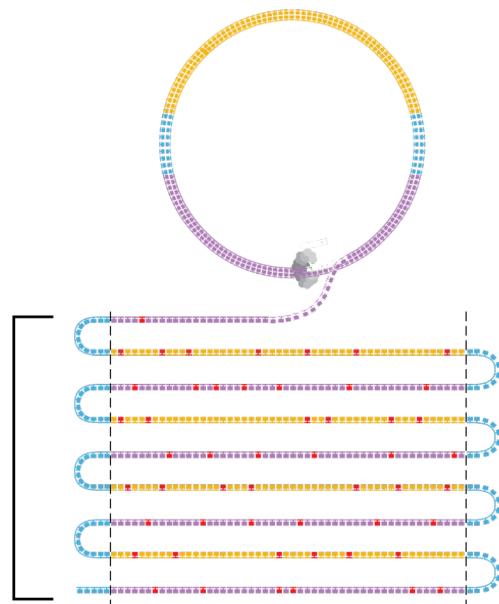


Anneal primers and bind DNA polymerase

Circularized DNA is sequenced in repeated passes

The polymerase reads are trimmed of adapters to yield subreads

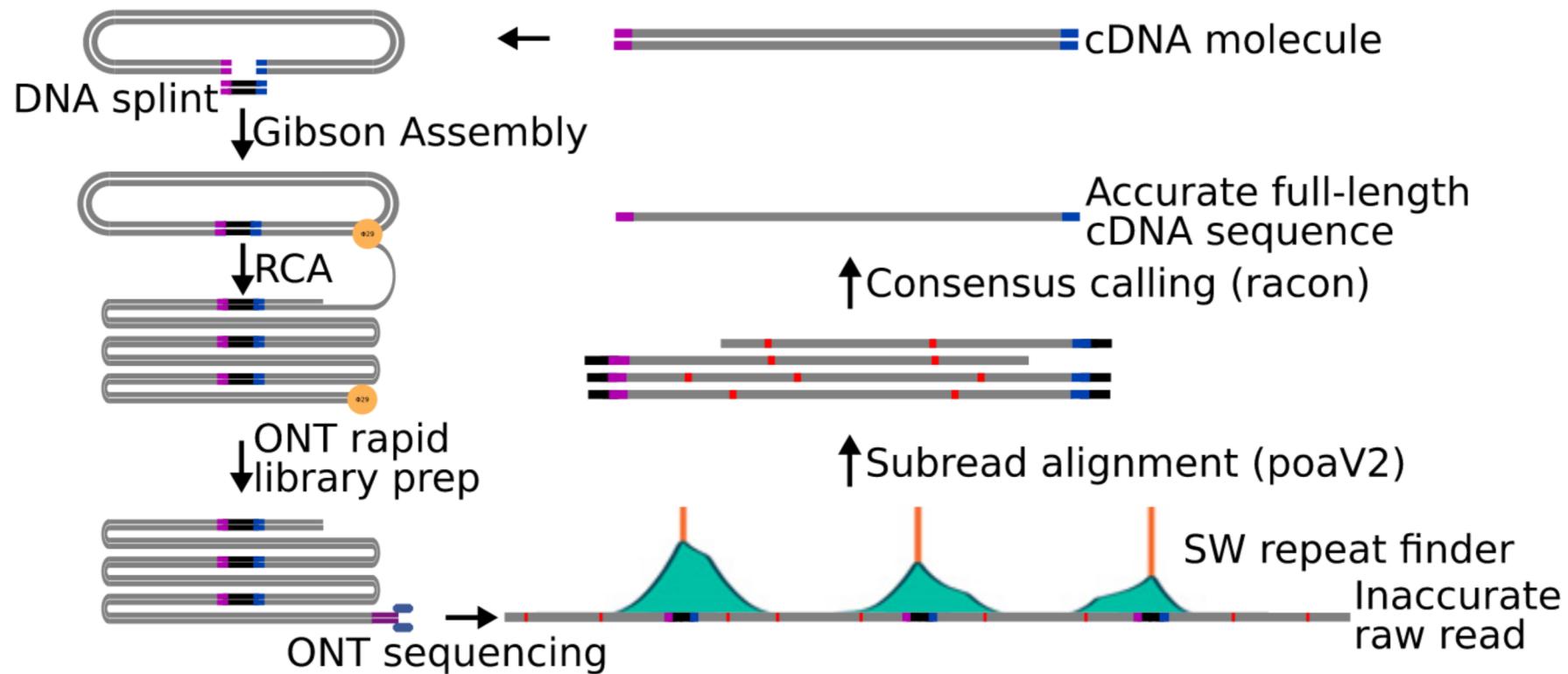
Consensus is called from subreads



CCS

(>99% accuracy)

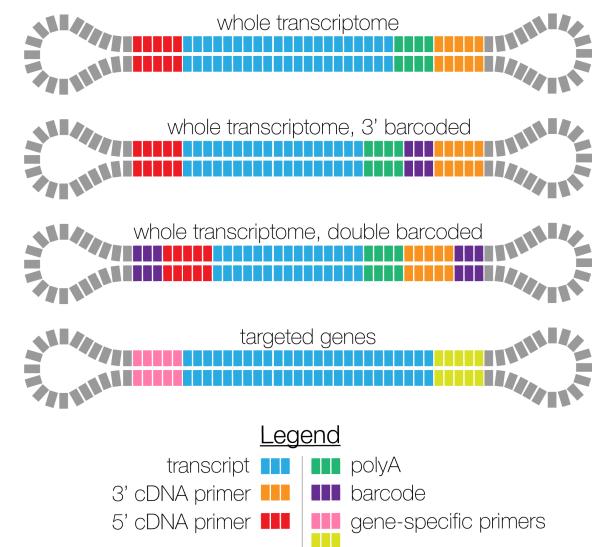
CCS for ONT? R2C2

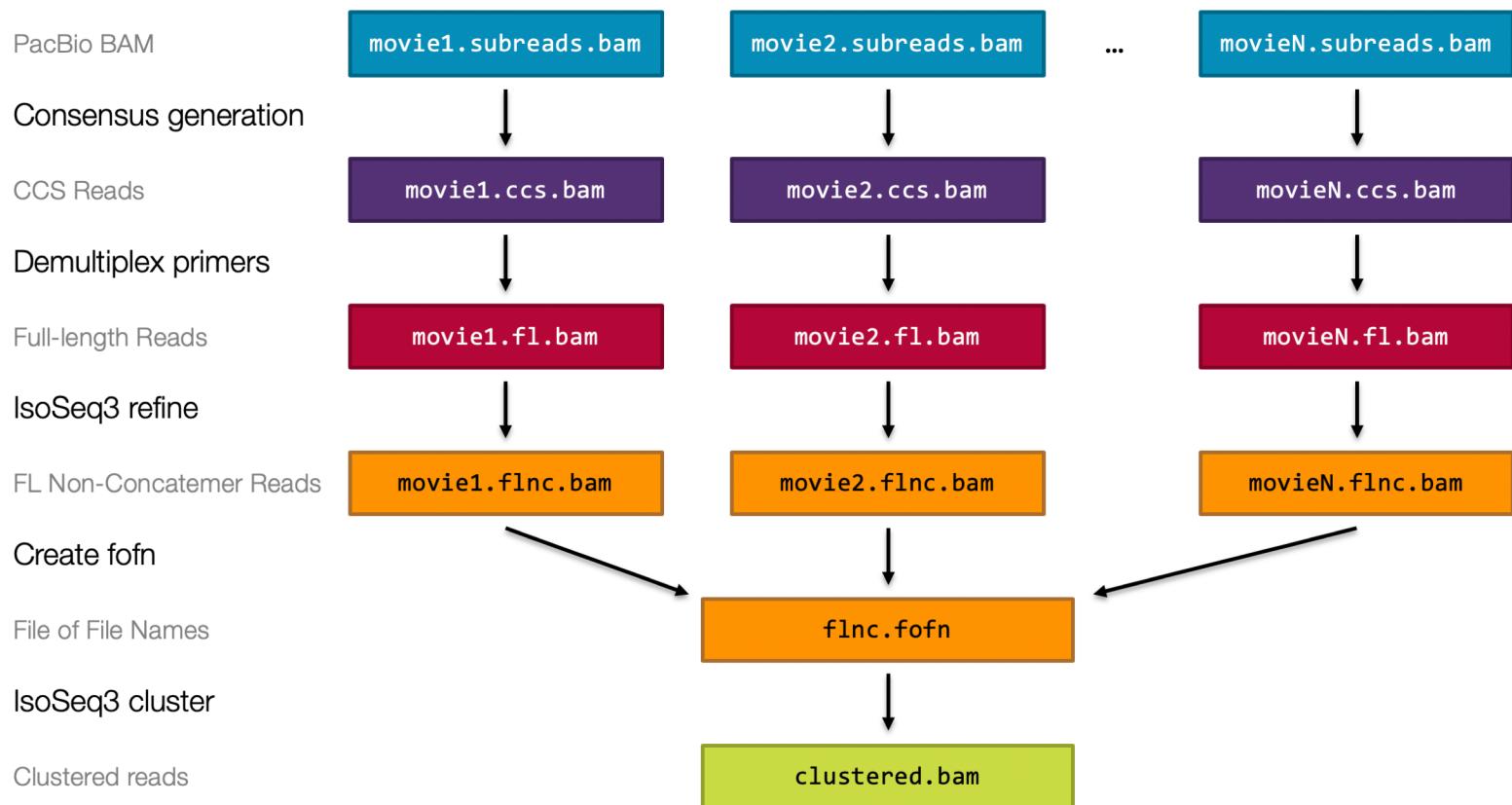


Volden, R. et al (2018). Improving nanopore read accuracy with the R2C2 method enables the sequencing of highly multiplexed full-length single-cell cDNA. *PNAS*, 115(39)

PacBio RNA library prep

- General workflow: IsoSeq
- cDNA & PCR amplification
- Multiplexing
- github.com/PacificBiosciences/IsoSeq





	ONT	PacBio
Read accuracy	~97%	~90% (>99% CCS)
Read length	Transcript length	Transcript length
Transcript quantification	Yes, PCR & cDNA free possible	Yes, with PCR
RNA base modifications	Yes (m6A) ¹	No
DNA base modifications	Yes (m5C, m6A) ²	Yes (m5C, m6A, hm5C) ³
Throughput (BIF)	~250M reads/run	~30M CCS reads/run

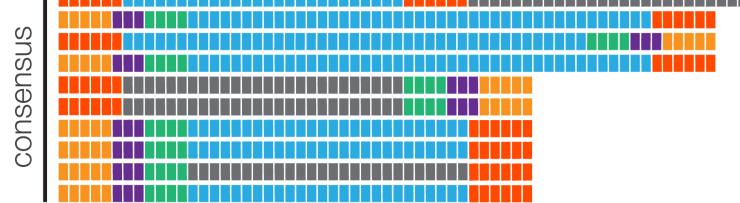
1. Liu, H., et al (2019). Accurate detection of m6A RNA modifications in native RNA sequences. *Nature Communications*, 10(1), 1–9
2. Liu, Q., et al (2019). Detection of DNA base modifications by deep recurrent neural network on Oxford Nanopore sequencing data. *Nature Communications*, 10(1).
3. Flusberg, B. A., et al (2010). Direct detection of DNA methylation during single-molecule, real-time sequencing. *Nature Methods*, 7(6), 461–465

Iso-Seq Clustering

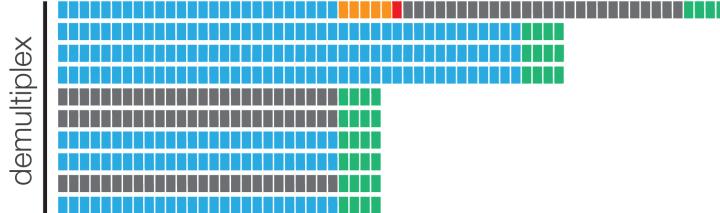
transcript gene A
transcript gene B
optional polyA

3' cDNA primer
5' cDNA primer
sample barcode

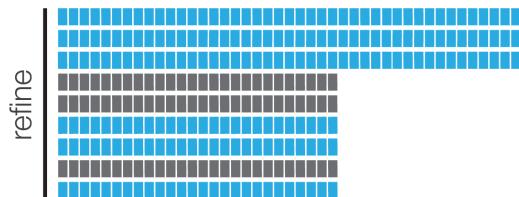
Version 8, Dr. Armin Töpfer



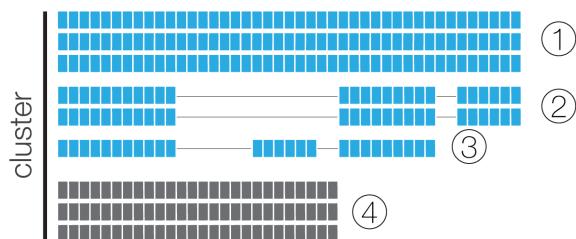
- Use **polished** CCS reads
- Only full-pass ZMWs



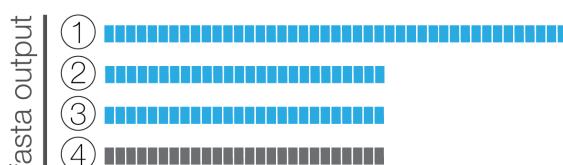
- Barcoded and unbarcoded cDNA primer removal
- Orientation
- Unwanted primer combination removal



- PolyA tail trimming
- Concatemer removal



- Hierarchical, $n^*\log(n)$ clustering,
alignment of shorter to longer sequences
- Iterative cluster merging
- Generate consensus for each read cluster
using QV guided PoA



- Fasta output is split into HQ and LQ reads
- One consensus per read cluster

Isoform analysis software

- Pinfish
 - Developed by nanoporetech
 - github.com/nanoporetech/pinfish
- Mandalorion
 - C3PO/R2C2 only
 - github.com/rvolden/Mandalorion-Episode-II
- FLAIR
 - Optional polishing with Illumina reads
 - Also PacBio
 - github.com/BrooksLabUCSC/flair