

Long-read sequence analysis

Sequencing technologies

Question 4

What is a long read?

- Short read: 50-300 bp, often paired-end (Illumina sequencing)
- Long read: > 1kb, up to 20 Mb:
 - single molecule sequencing or
 - 3rd generation sequencing

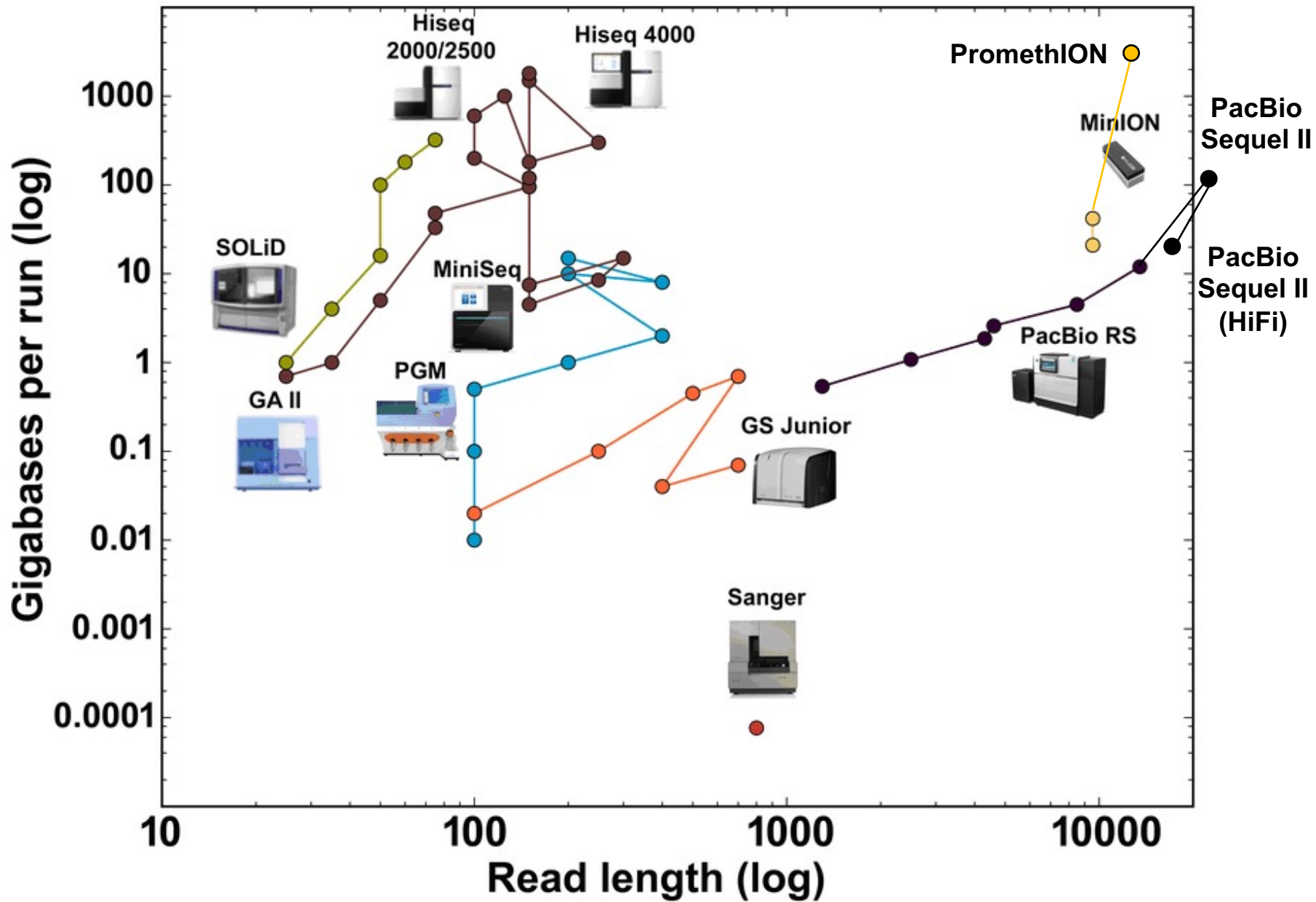


Image from: G. Silva (2016)

Illumina sequencing

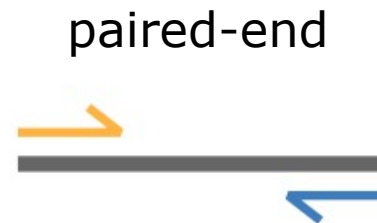
- Sequencing-by-synthesis: 2nd generation sequencing
- Massive throughput: up to 500×10^9 bases/run
- Most used platform today

illumina[®]



Illumina sequencing

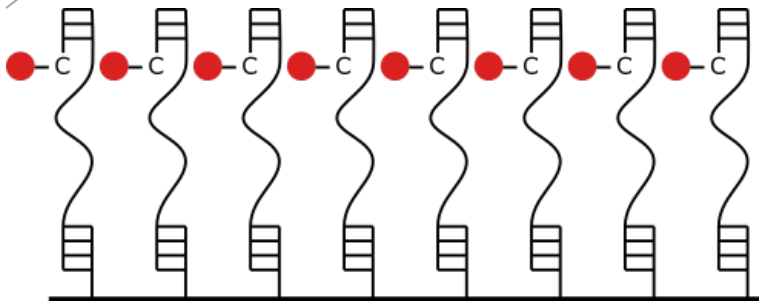
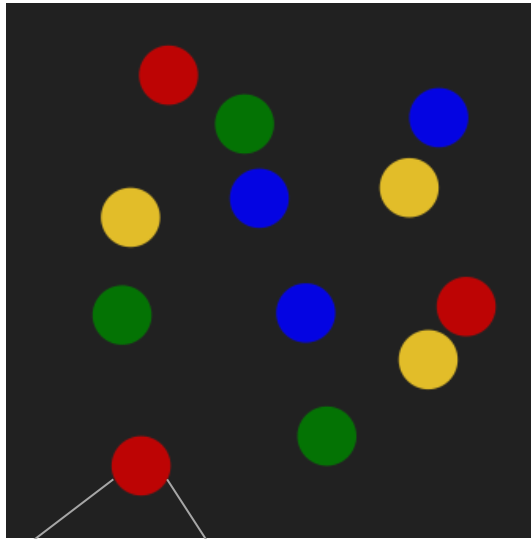
- 50 – 300 bp
- Paired-end (or single-end)



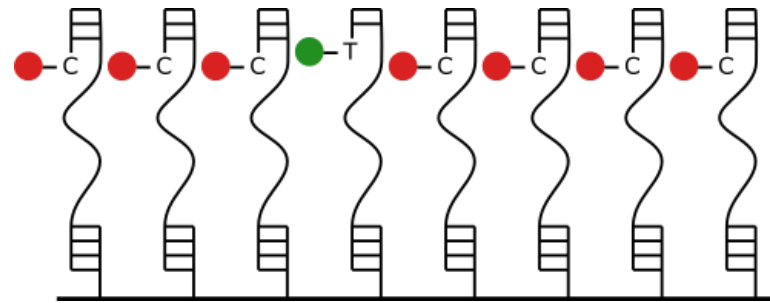
Question 5

Illumina - limitations

- Maximum read length: 300 bp
- How to reconstruct:
 - Repeats?
 - Isoforms?
 - Structural variation?
 - Haplotypes?
 - Genomes?
- Why not longer read lengths with Illumina?



in phase



out of phase

Long reads (3rd generation)

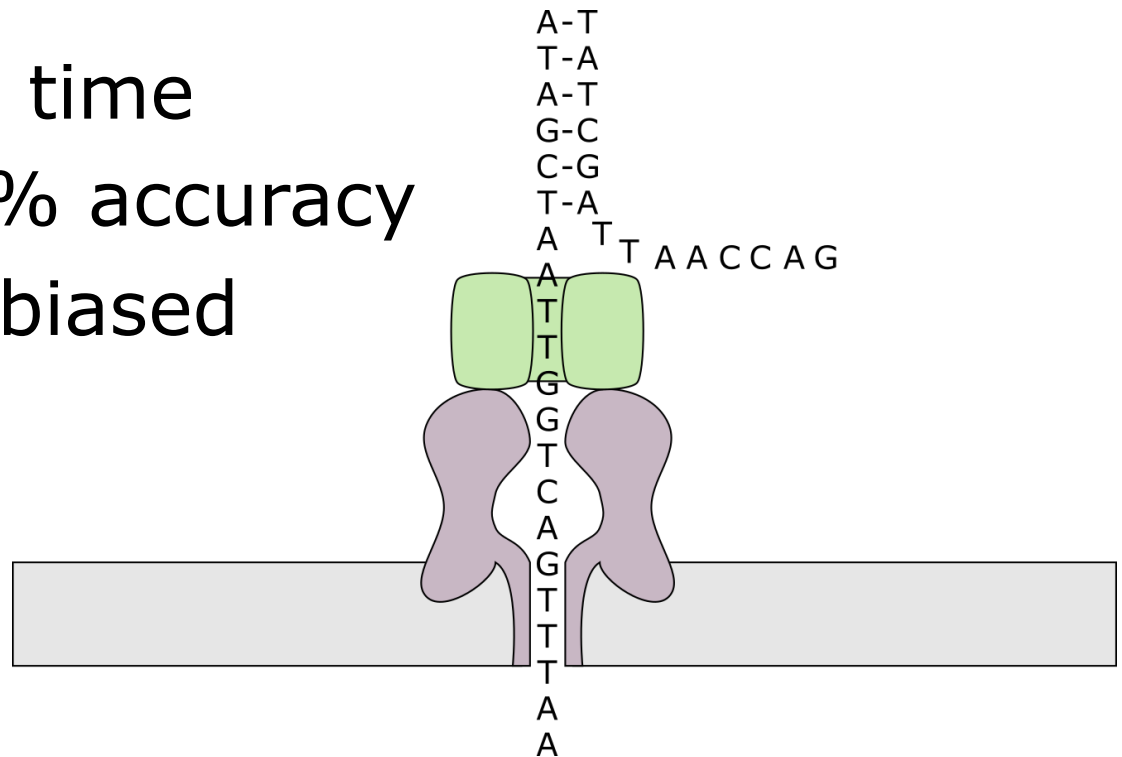
- Crux: maximizing signal from a single-molecule base read-out
- Single molecule, so no out-of-phase signal
- Two frequently used platforms:
 - PacBio SMRT sequencing
 - Oxford Nanopore Technology



Question 6

Oxford Nanopore technology

- Based on changes in electrical current
- Well-known for its scalability and portability
- 4 bp read at a time
- Up to ~95-97% accuracy
- Errors can be biased





ONT scalability

1 small
flow cell:
1 x 2.8 Gb



Flongle

1 medium
flow cell:
1 x 50 Gb



MinION

5 medium
flow cells:
5 x 50 Gb



GridION

24-48 big
flow cells:
48 x 290 Gb

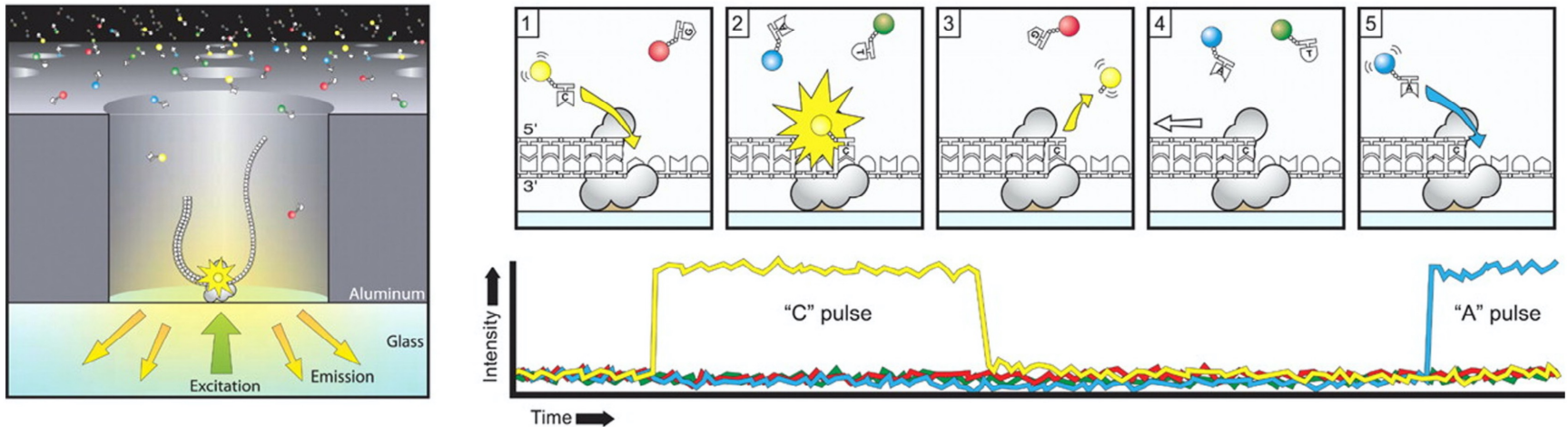


PromethION

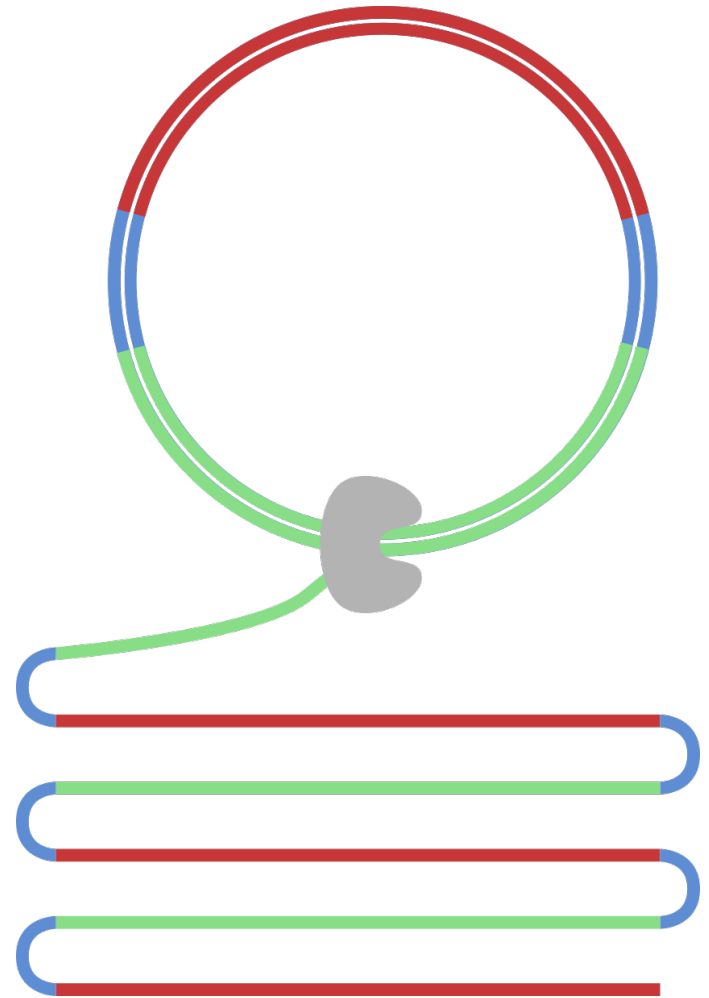
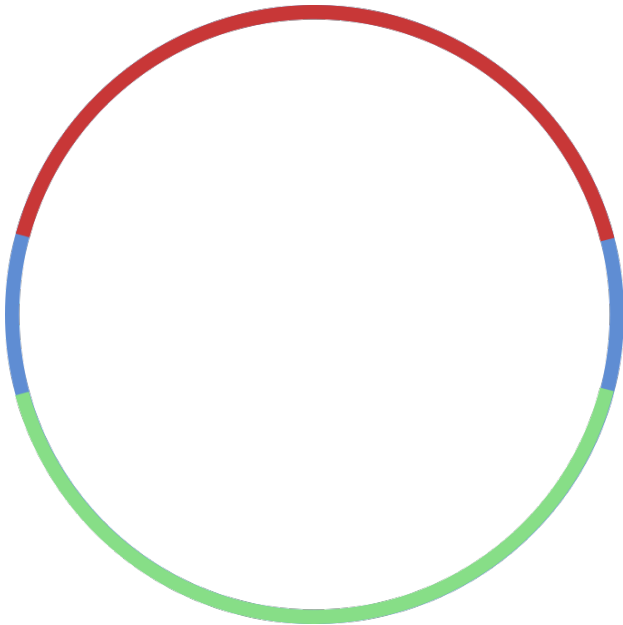
ONT library prep

- Standard kit:
 - $>1 \mu\text{g}$ HMW DNA
 - Shearing + size selection is optional
 - Multiplexing requires PCR step

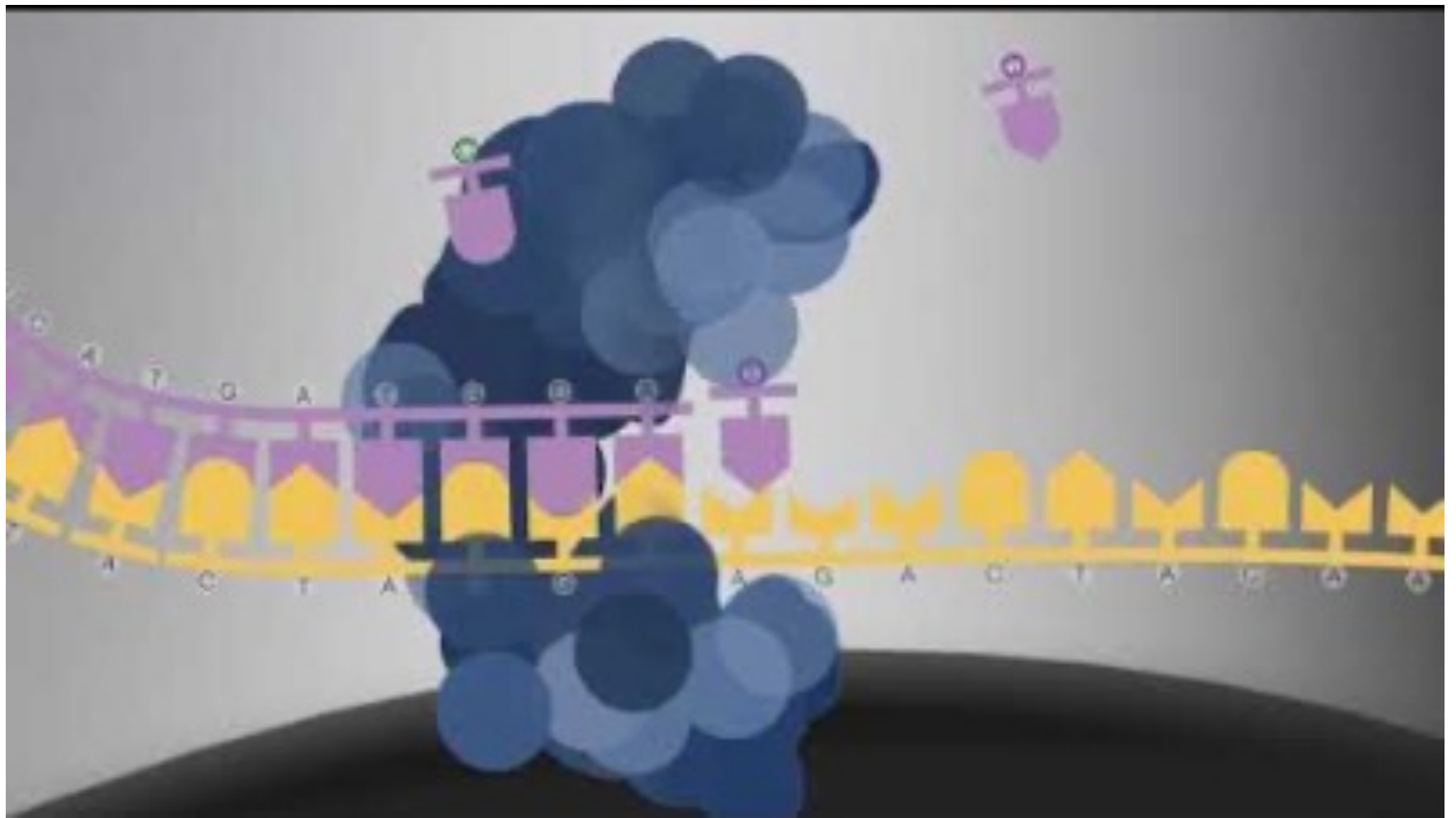
PacBio sequencing



- Polymerase bound to ZMW bottom
- Circular molecules
- Single read out $\sim 90\%$ accuracy
- HiFi: single molecule sequenced multiple times



Hi-Fi read



PacBio Sequel II



- Up to 8M CLR reads/SMRT cell
- 4M HiFi reads/SMRT cell
- start with $>5 \mu\text{g}$ HMW DNA
- Requires shearing + size selection
- Multiplexing requires PCR

	ONT	PacBio
Read accuracy	~90-95%	~90% (>99% HiFi)
Read length	up to 2 Mb	up to 30-40 kb (HiFi) up to 200 kb (CLR)
RNA base modifications	Yes (m6A) ¹	No
DNA base modifications	Yes (m5C, m6A) ²	Yes (m5C, m6A, hm5C) ³
Throughput (BIF)	~500M reads/run ⁴	~4M HiFi reads/run ~8M CLR 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

4. 48 flow cells on a PromethION

Question 7&8