

# NGS - quality control, alignment, visualisation

Sequencing technologies

# Major applications

- Transcriptome characterization
  - e.g. RNA-seq
- Epigenome characterization:
  - e.g. ATAC-seq
- DNA-protein interactions:
  - e.g. ChIP-seq
- Whole genome (assembly)
- Variant detection
- Metagenome characterization
- Any others?

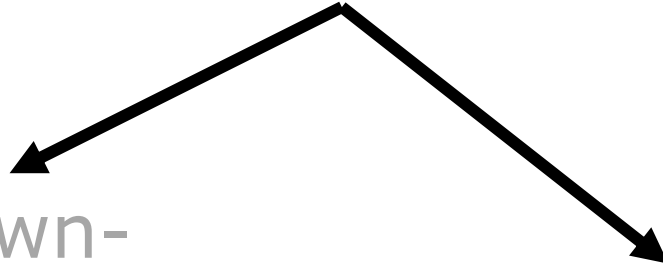
Sequencing



Quality control

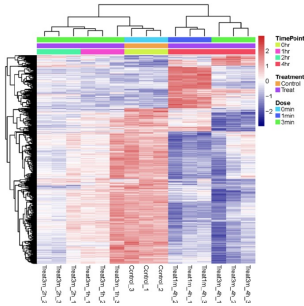
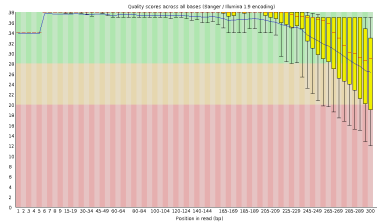


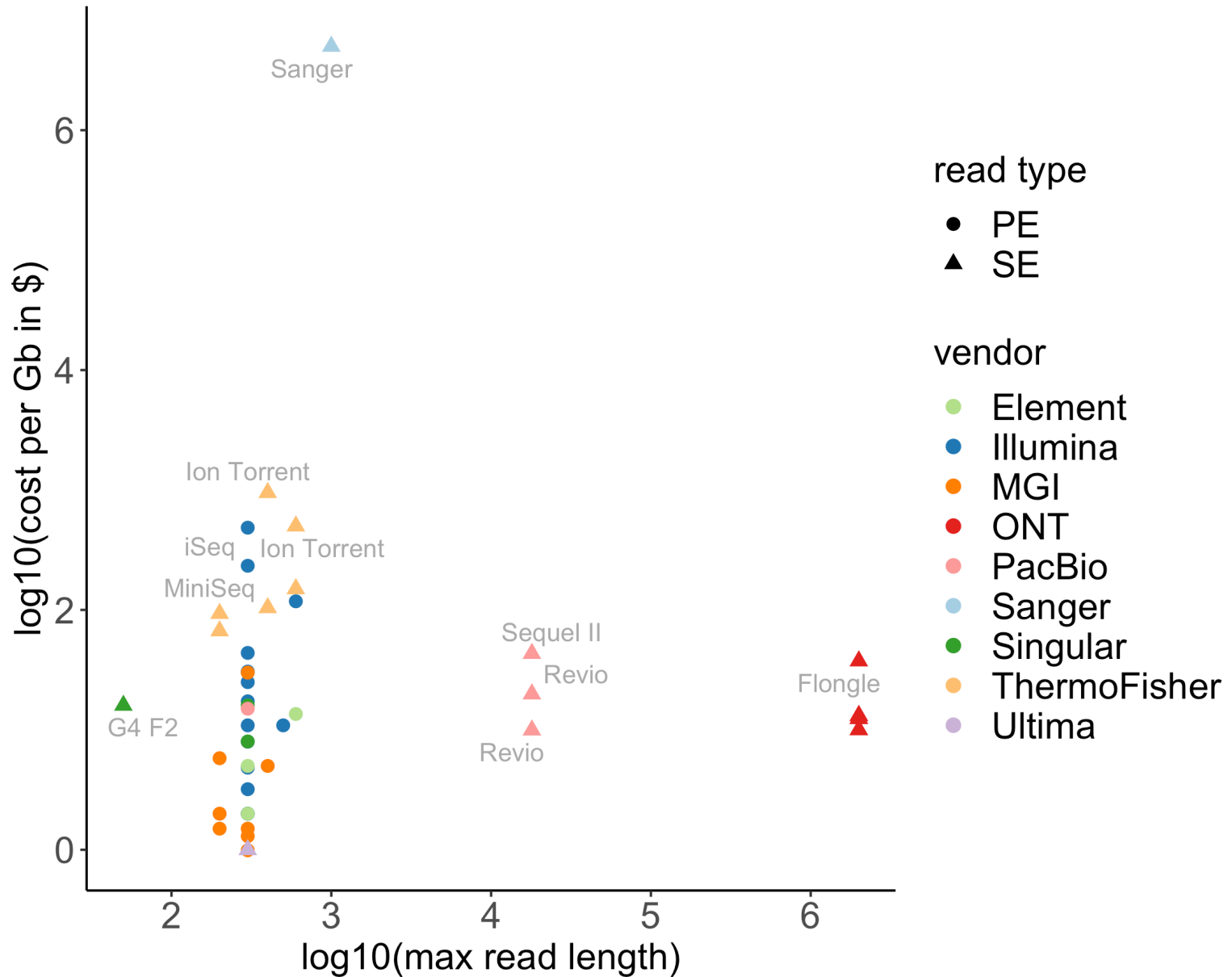
Alignment



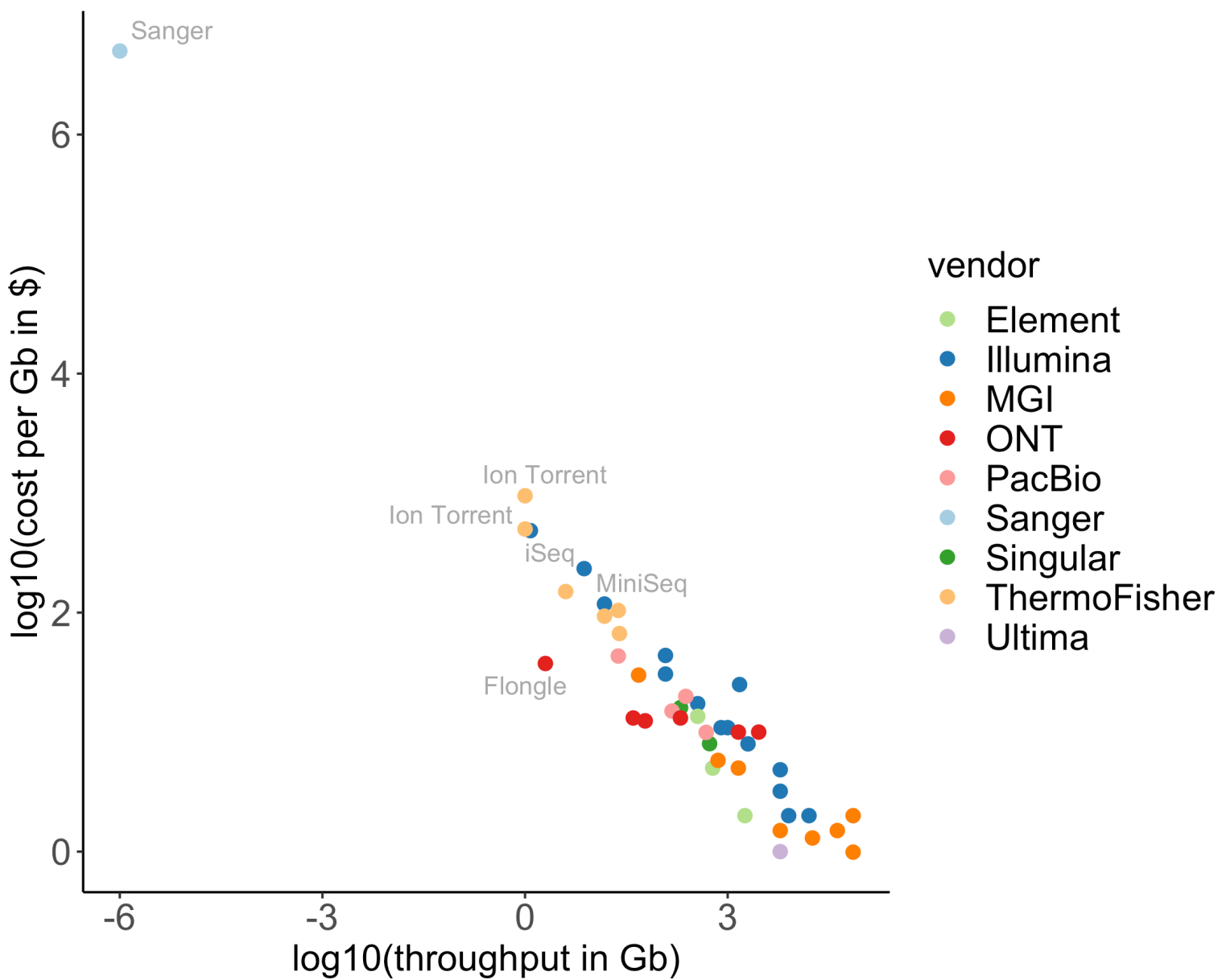
Down-  
stream  
analysis

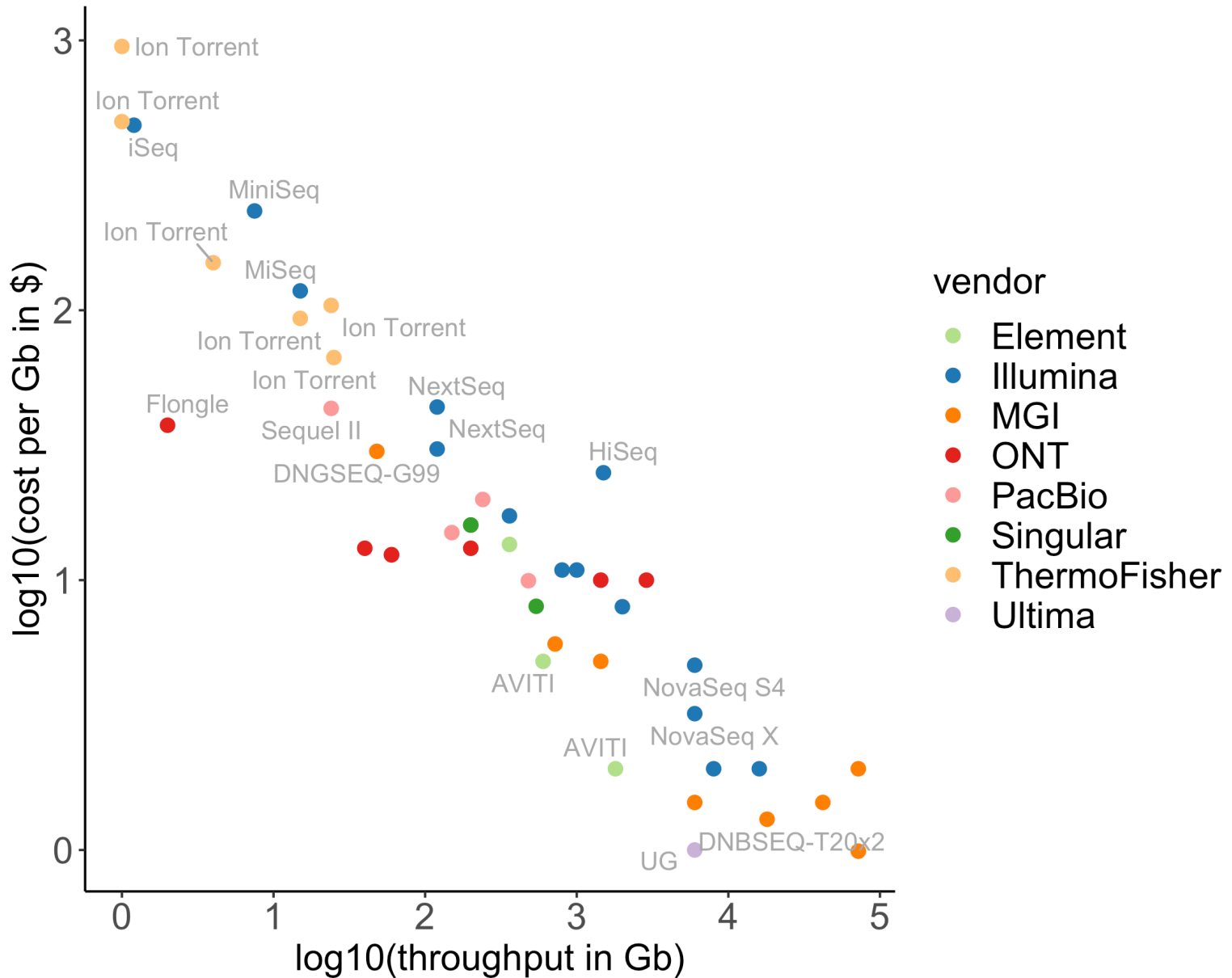
Visualisation





# Quiz Question 4





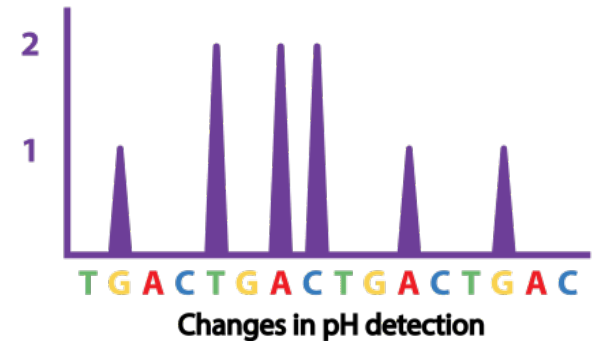
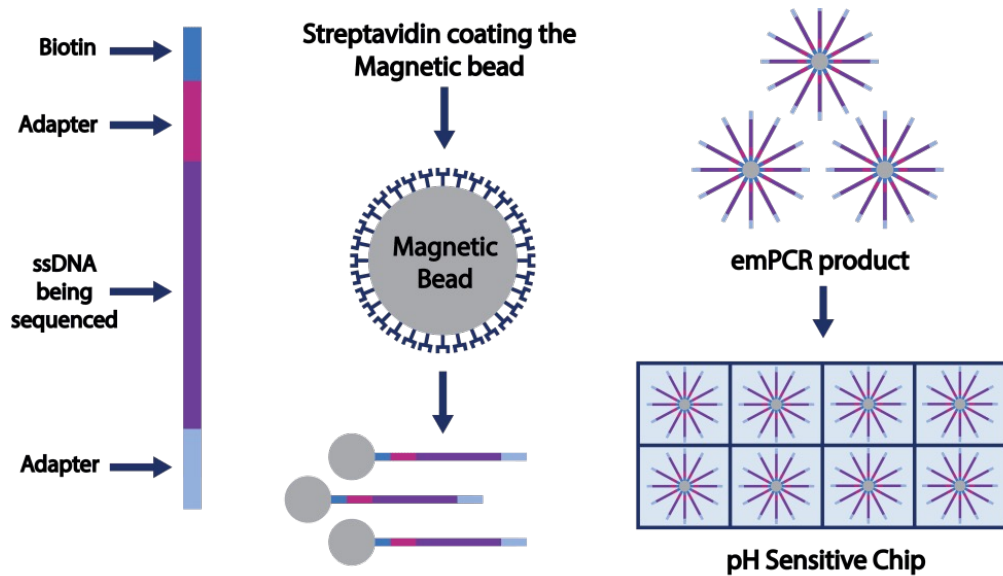
drawn from: <https://docs.google.com/spreadsheets/d/1GMMfhyLK0-q8XkIo3YxlWaZA5vVMuhU1kg41g4xLkXc/> Albert Vilella

# This course

- 2nd generation (sequencing by synthesis):
  - Ion Torrent
  - Illumina
- 3rd generation:
  - Pacific Biosciences
  - Oxford Nanopore Technology



# Ion Torrent sequencing



# Ion Torrent sequencing

- Up to  $\pm 400$  bp read length
- Scalable (but Illumina has similar size systems nowadays)
- Homopolymers (e.g. TTTT) are a challenge (impossible) to sequence

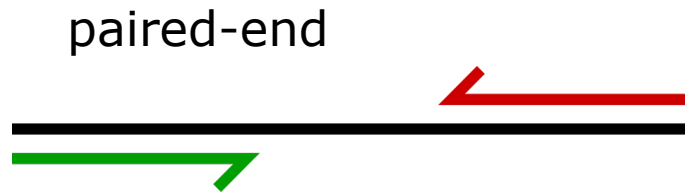


# Illumina sequencing

- Massive throughput: up to  $16 \times 10^{12}$  bases/run (NovaSeq X) =  $\sim 9,000$  whole exomes
- Most used platform today

# Illumina sequencing

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

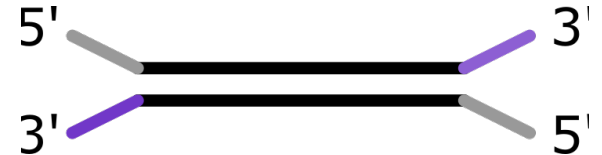


# Illumina library prep

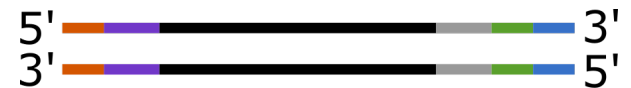
shear + size select DNA



Ligate adapters



Barcode + p5/p7 sites



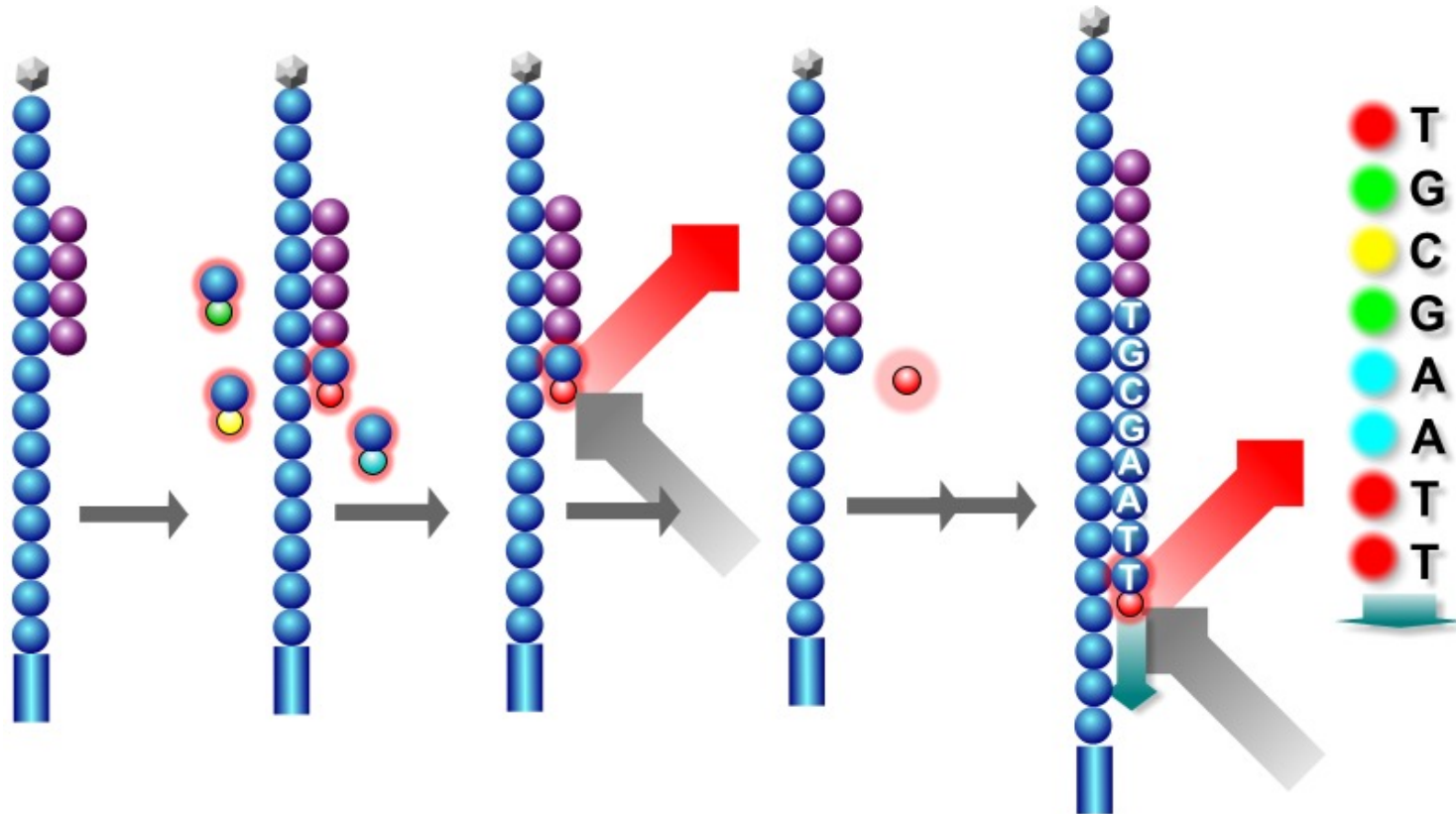
PCR: 8-16 cycles



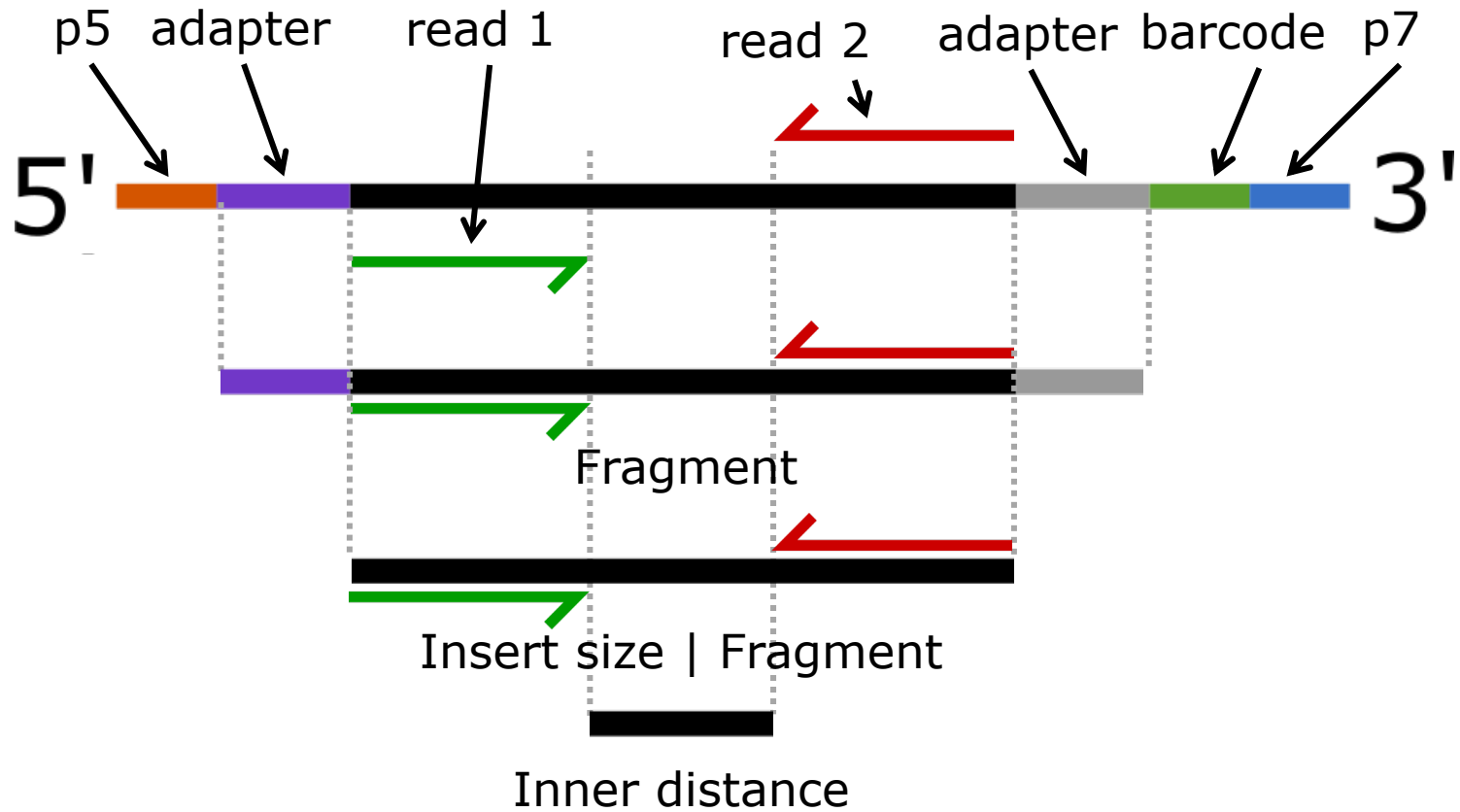
Sequencing



# Sequencing by synthesis



# Some definitions



# Some more definitions..

- **Library:** fragments from one (c)DNA sample that share a barcode
- **Sequencing run:** complete cycle of generating reads on a machine
- **Flow cell:** physical platform where sequencing reactions take place. Used once in a sequencing run.
- **Lane:** compartment within the flow cell. An Illumina flow cell often has multiple lanes (2 or 4)



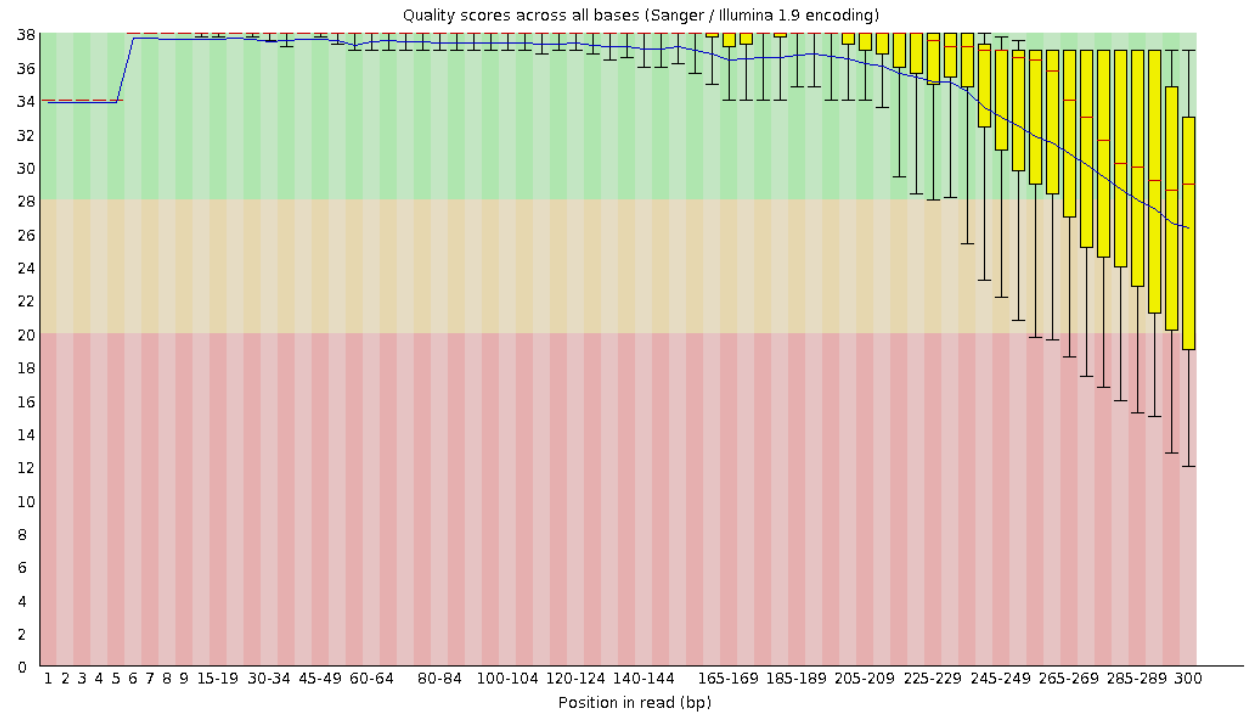
# Quiz Question 5

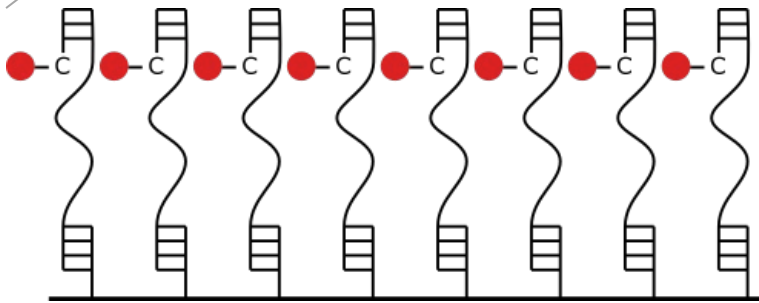
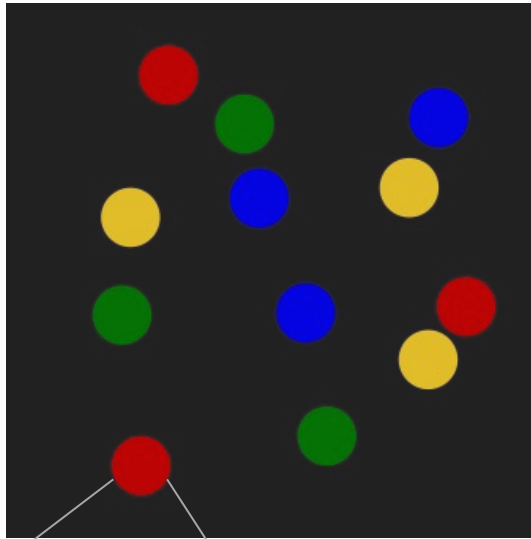
# Illumina - limitations

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

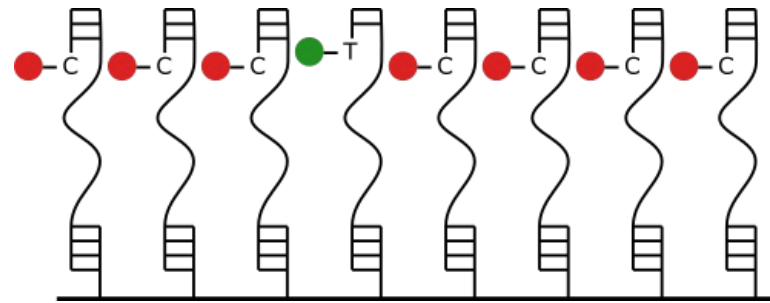
# Illumina - limitations

Sequence quality declines towards the end





in phase



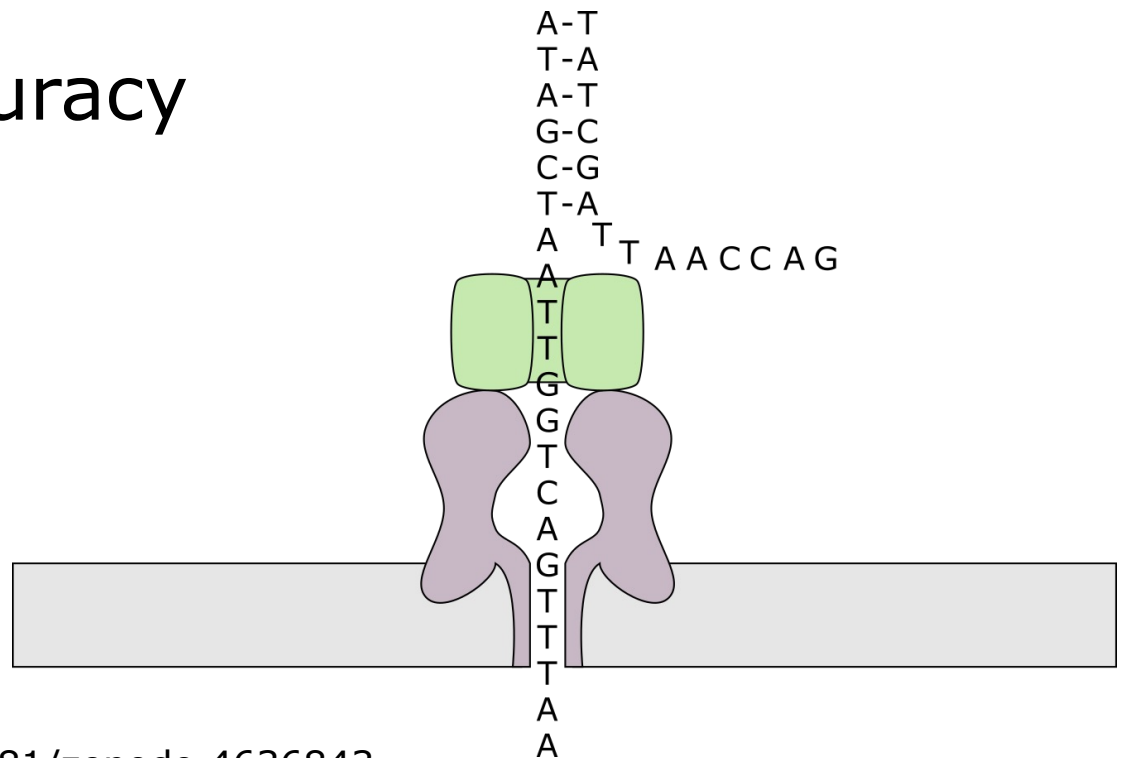
out of phase

# Long reads (3rd generation)

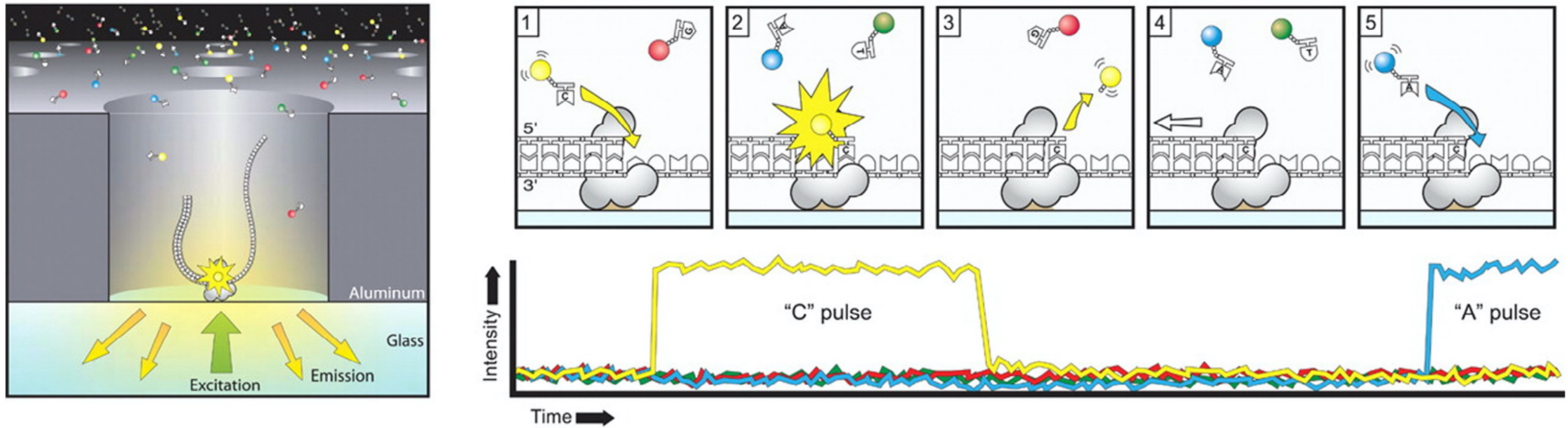
- Crux: maximizing signal from a single-molecule base read-out
- Single molecule, so no out-of-phase signal
- Two frequently used platforms:
  - PacBio SMRT sequencing
  - Oxford Nanopore Technology

# Oxford Nanopore technology

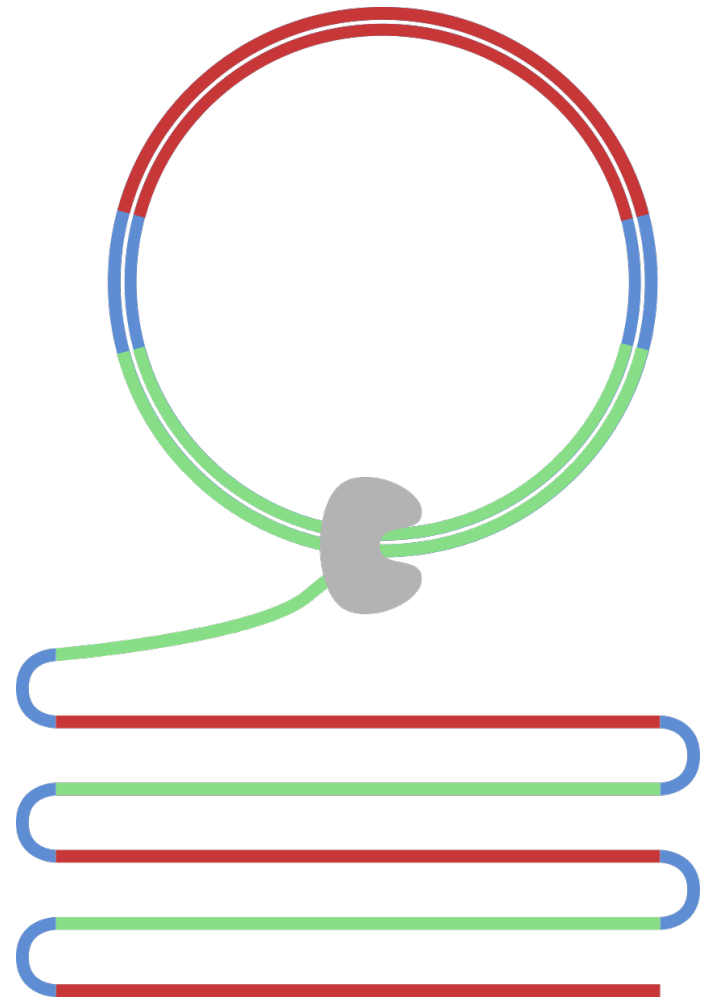
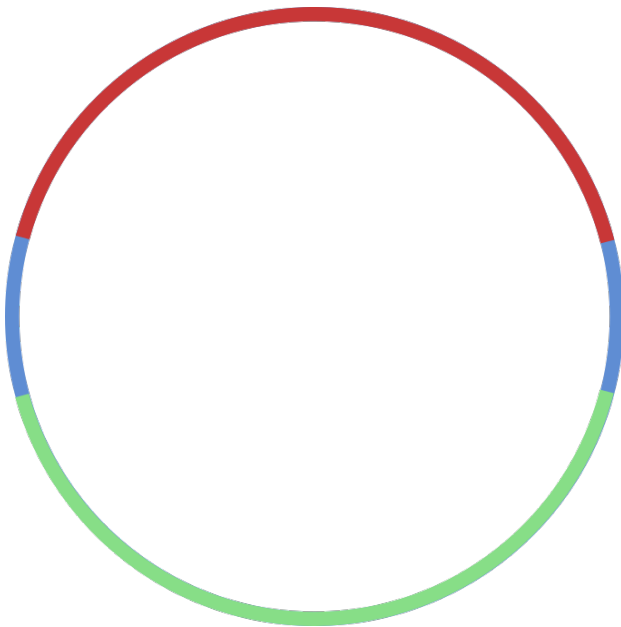
- Based on changes in electrical current
- Well-known for its scalability and portability
- ~95-97% accuracy



# PacBio sequencing



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



Hi-Fi read



# Quiz Question 6 and 7