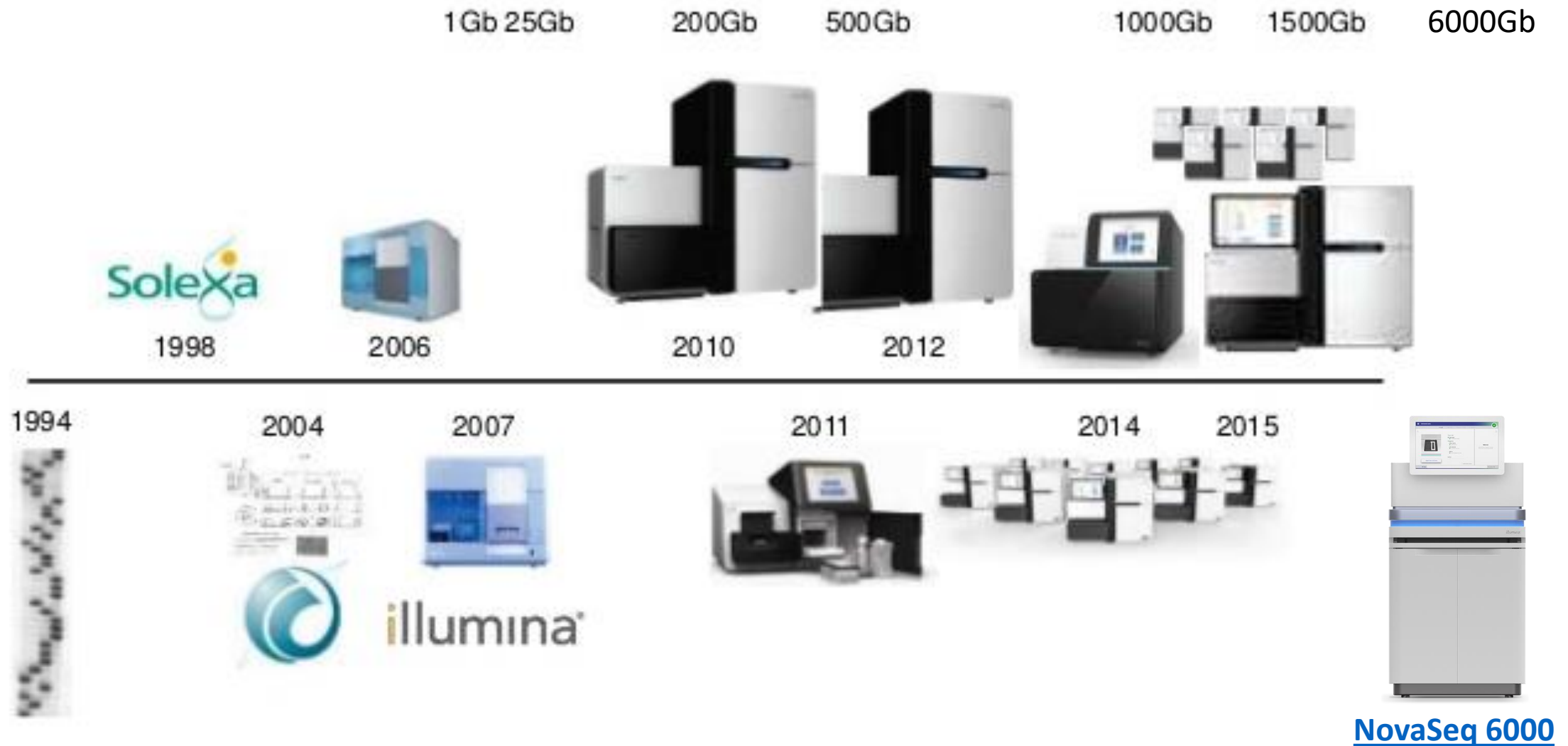


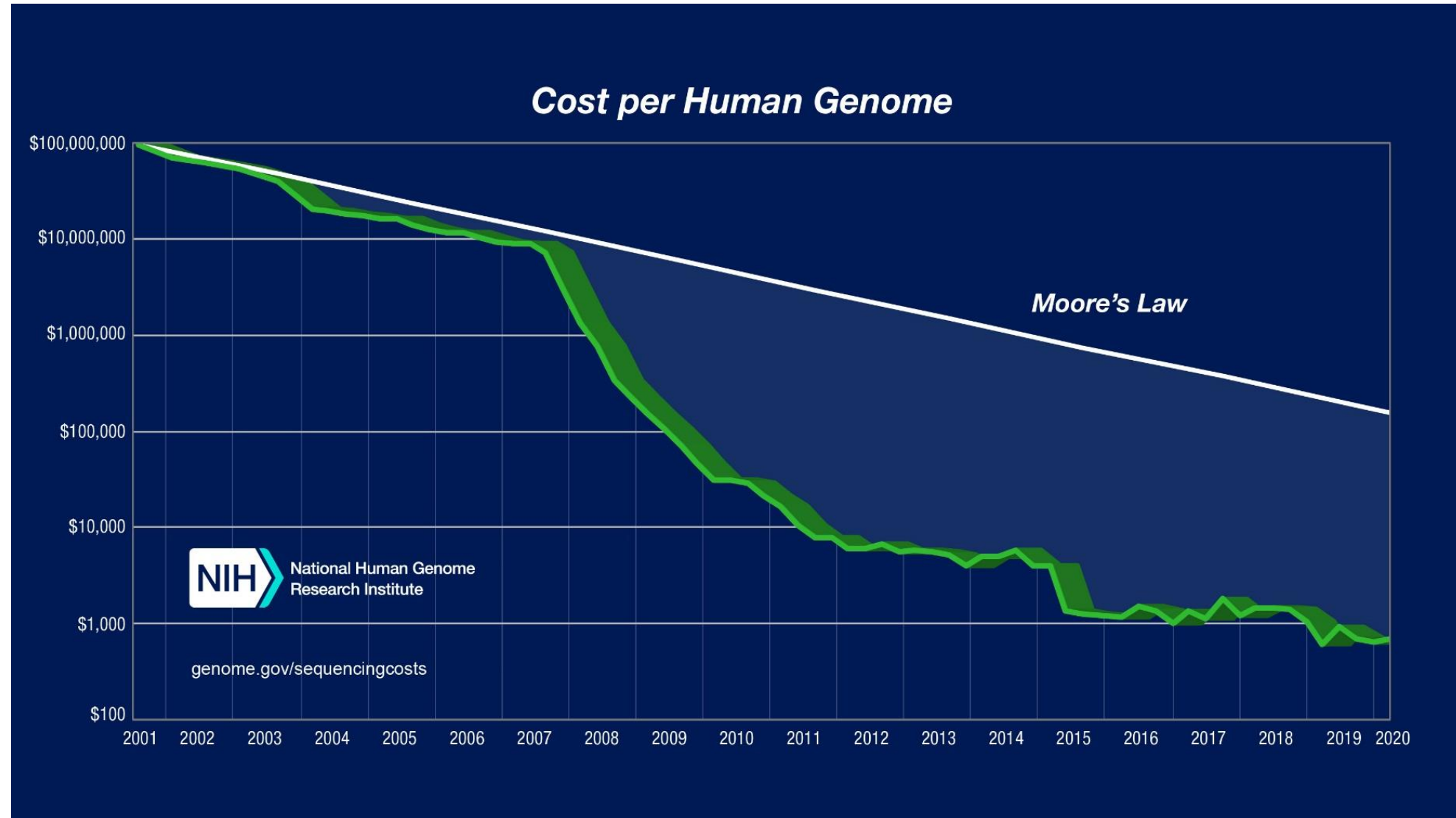
# **Next generation sequencing (NGS) introduction**

# High Throughput sequencing

# History of Illumina sequencing



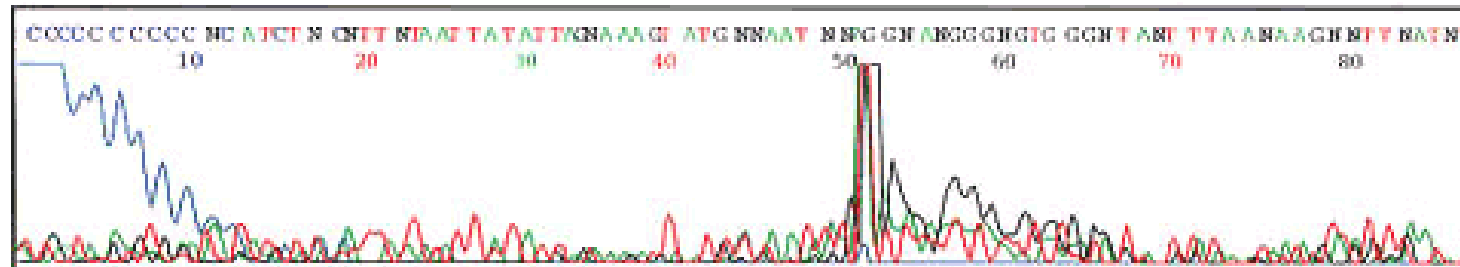
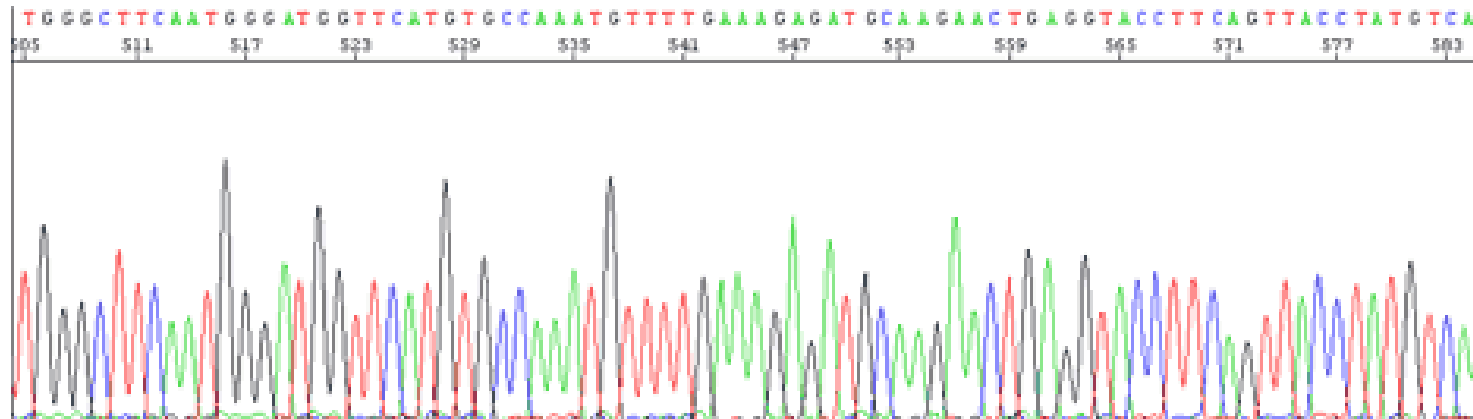
# Sequencing costs have decreased massively over time



# High Throughput Sequencing (=Next Generation Sequencing)

- **Short-read sequencing technologies** (2nd generation):
  - Sequence millions of clonally amplified molecules
  - Reads typically 150 bp long
  - Illumina
- **Long-read technologies** (3rd generation):
  - Single molecules are sequenced in real-time, fast but expensive and high error rates
  - Reads typically kb long
  - PacBio
  - Nanopore

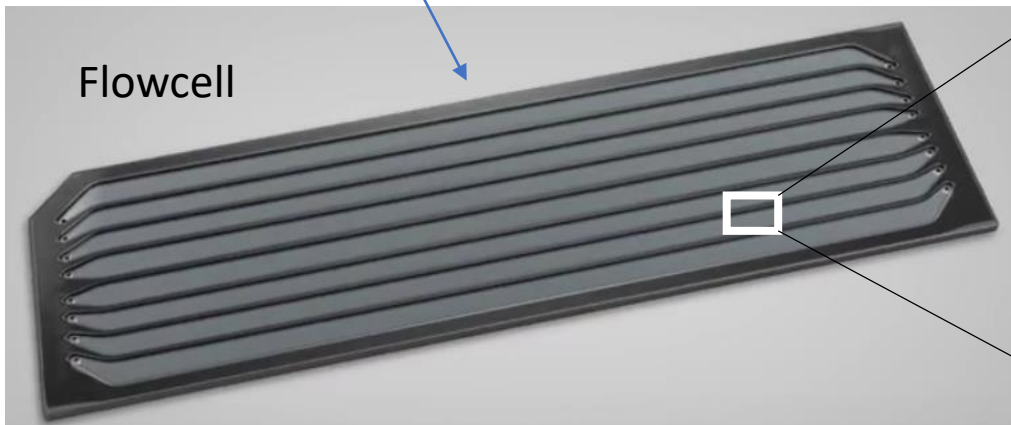
# Sanger Sequencing



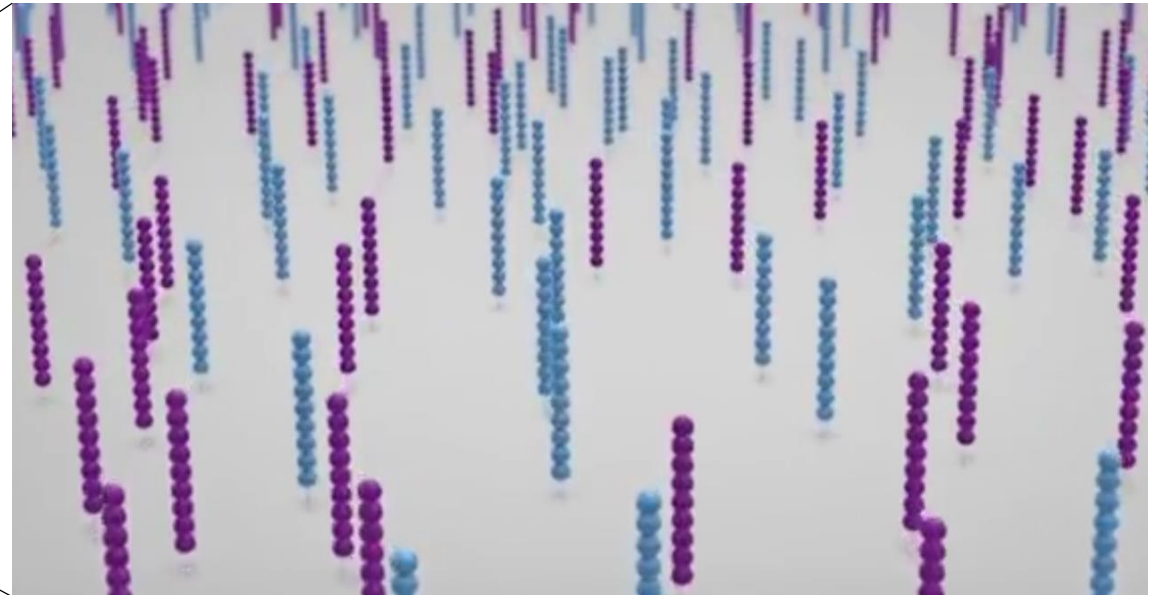
- Manually check each sequence
- Resequenced failed sequences

# Illumina flowcell: millions of DNA sequences

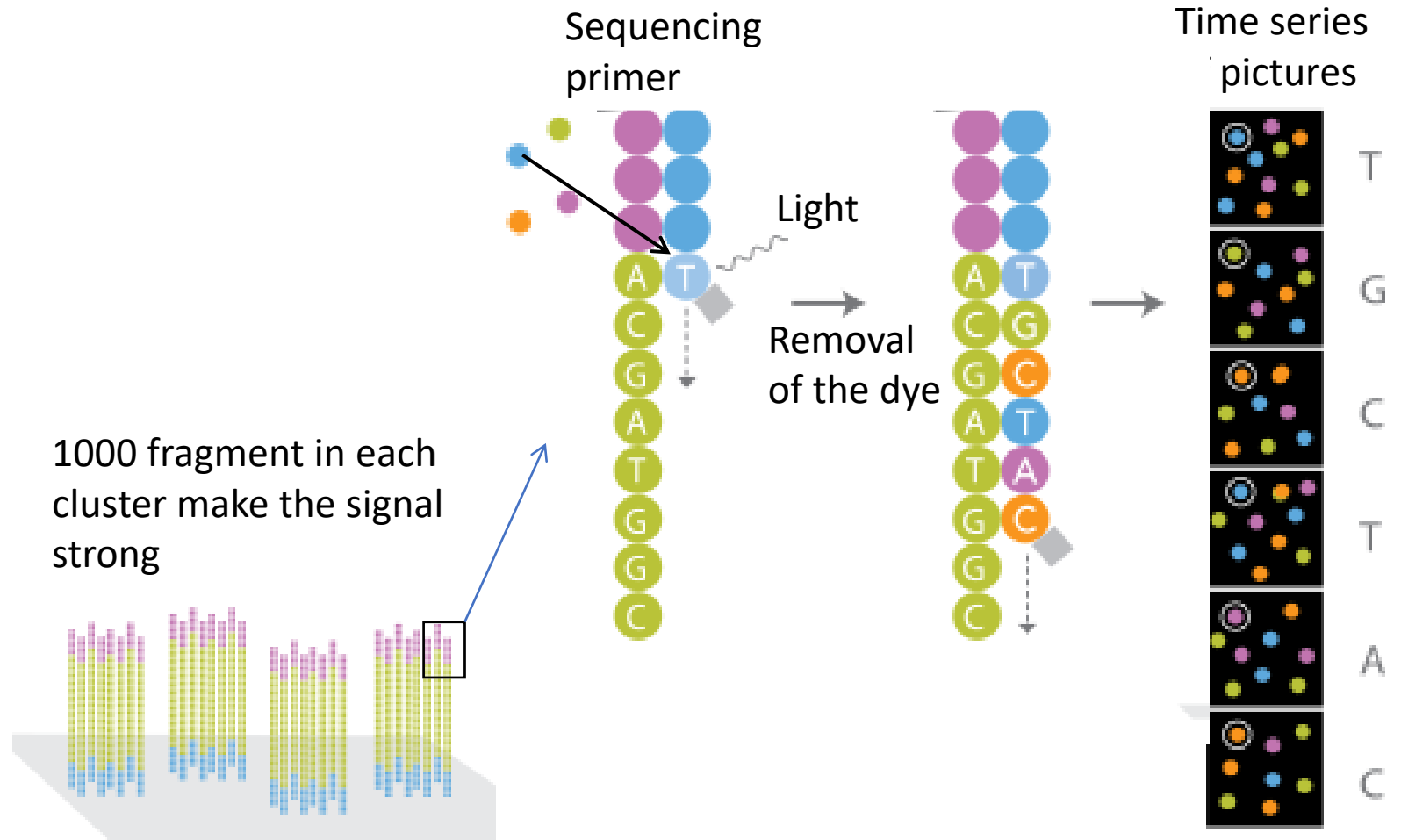
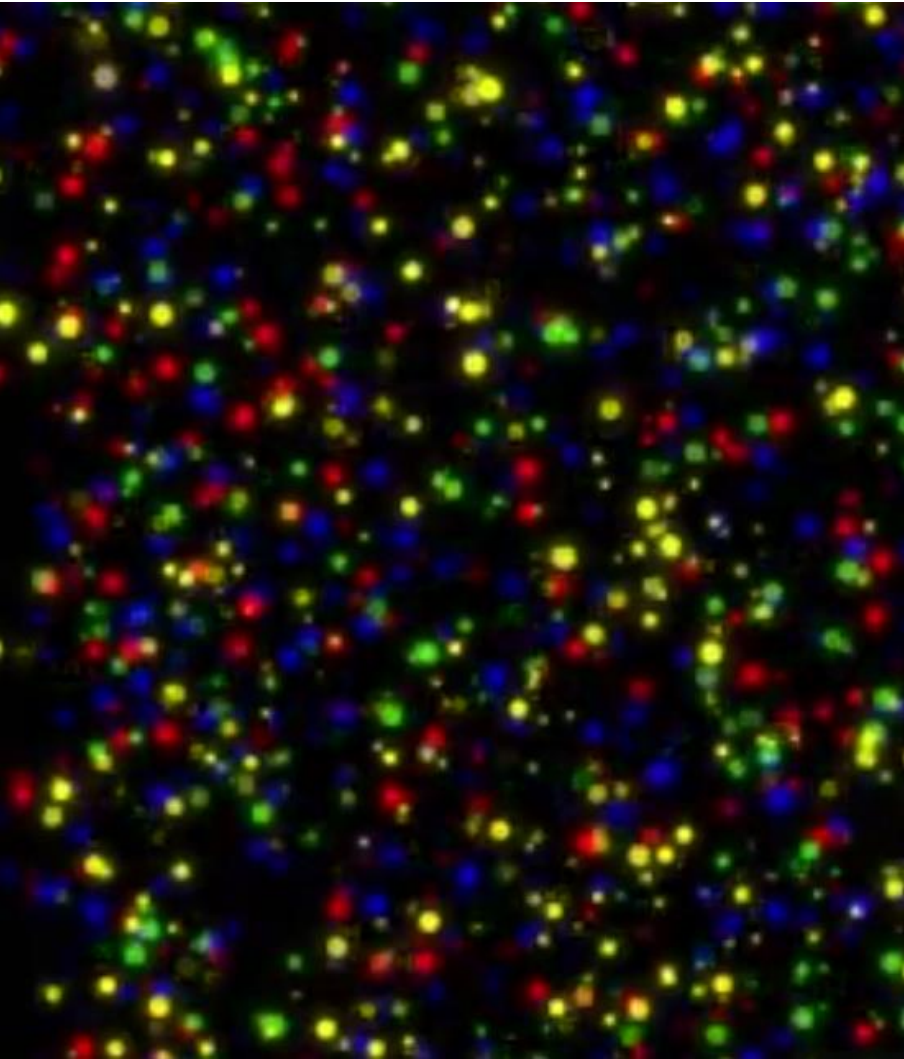
DNA fragments with Illumina adapters



Each lane contains a dense lawn of Illumina primers











# Sequencing by synthesis by Illumina












# 2-channel sequencing by synthesis

(used by these Illumina machines: Novaseq, Nextseq, MiniSeq)

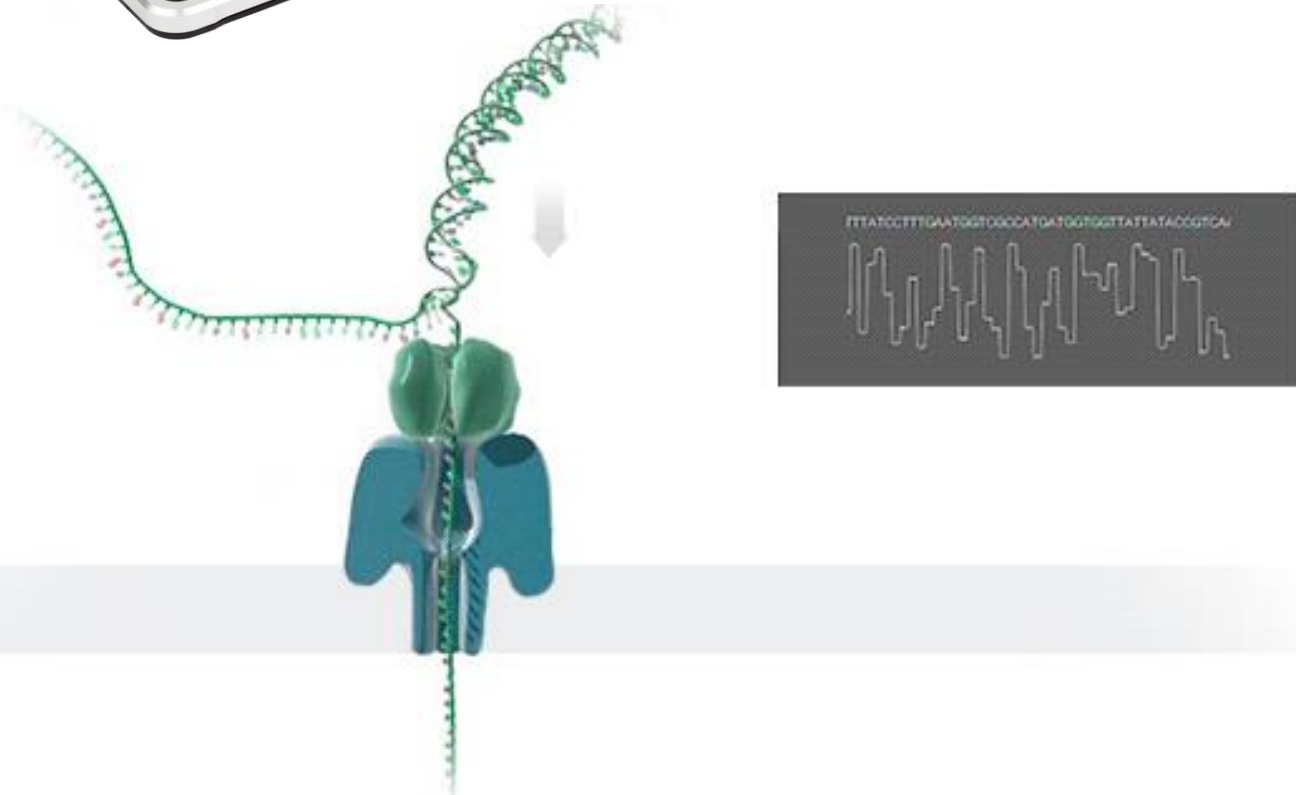
4-Channel Chemistry				
	 <b>A</b>	 <b>G</b>	 <b>T</b>	 <b>C</b>
Image 1				
Image 2				
Image 3				
Image 4				
Result	<b>A</b>	<b>G</b>	<b>T</b>	<b>C</b>

2-Channel Chemistry				
	 <b>A</b>	<b>G</b>	 <b>T</b>	 <b>C</b>
Image 1				
Image 2				
Result	<b>A</b>	<b>G</b>	<b>T</b>	<b>C</b>

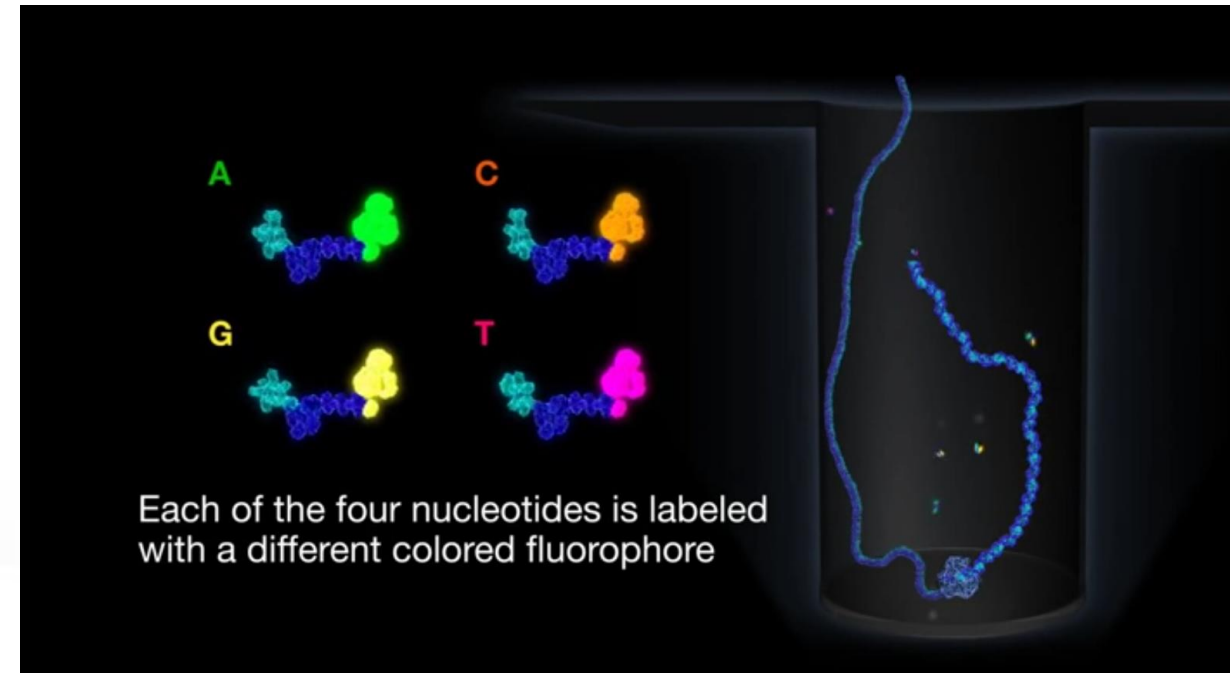
# Long read sequencing technologies



**Nanopore**

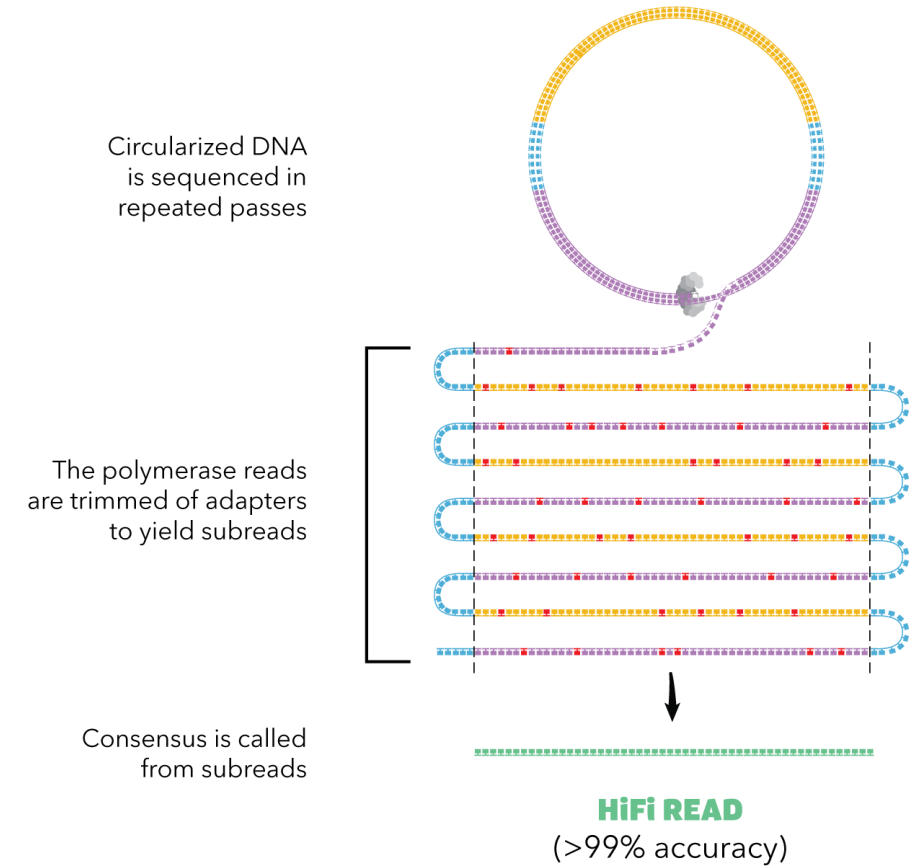
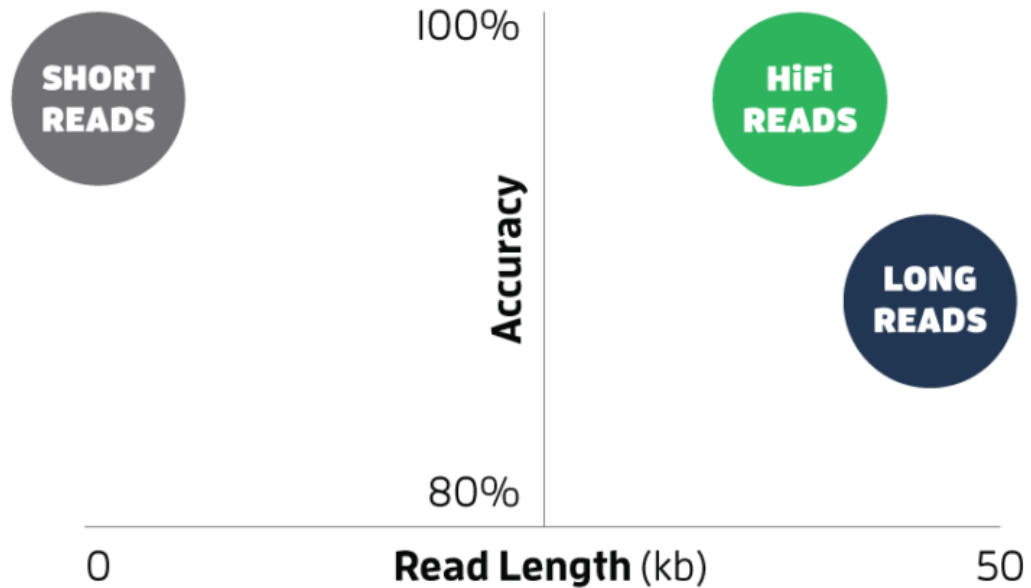


**PacBio**

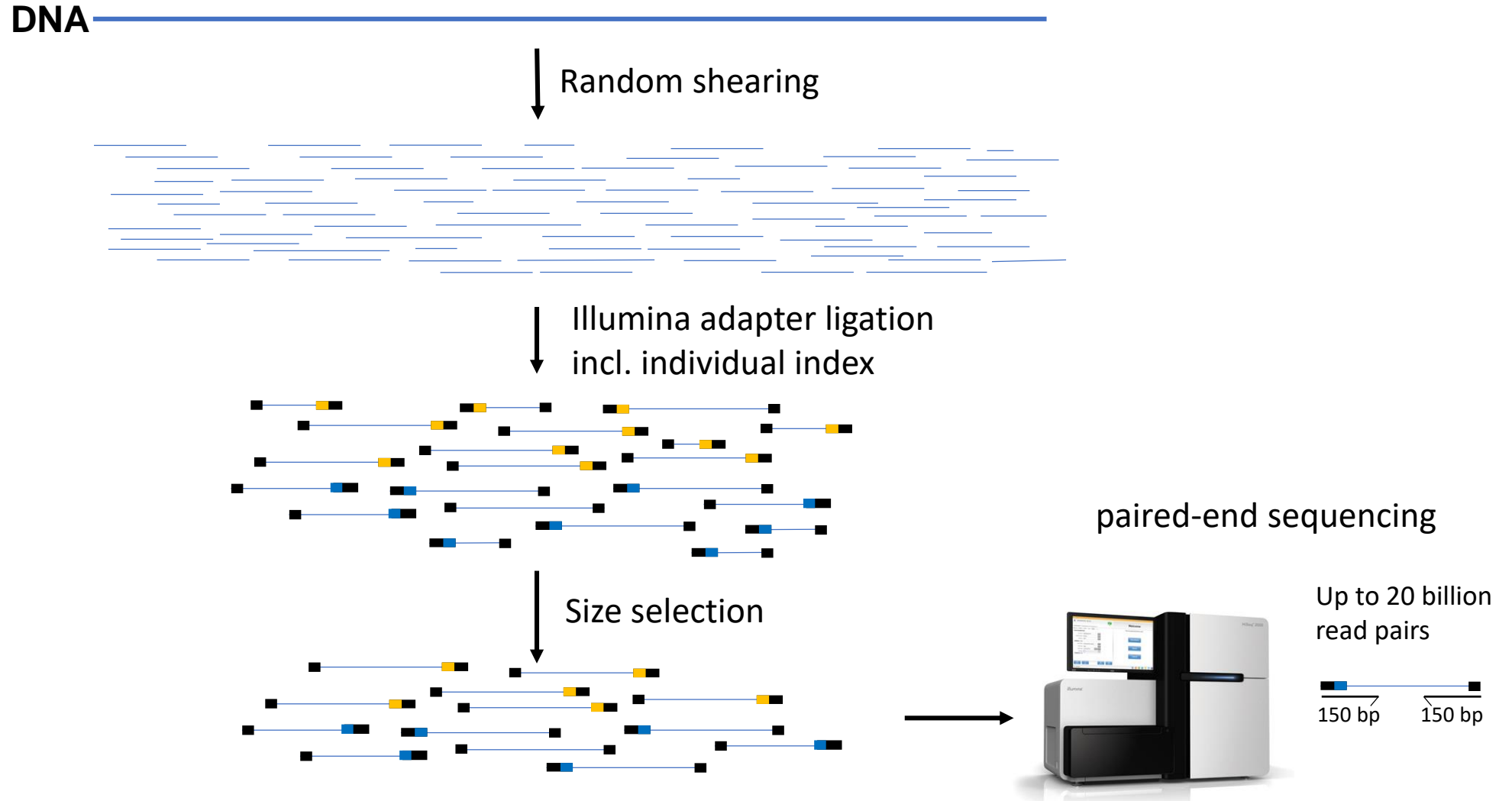


Each of the four nucleotides is labeled with a different colored fluorophore

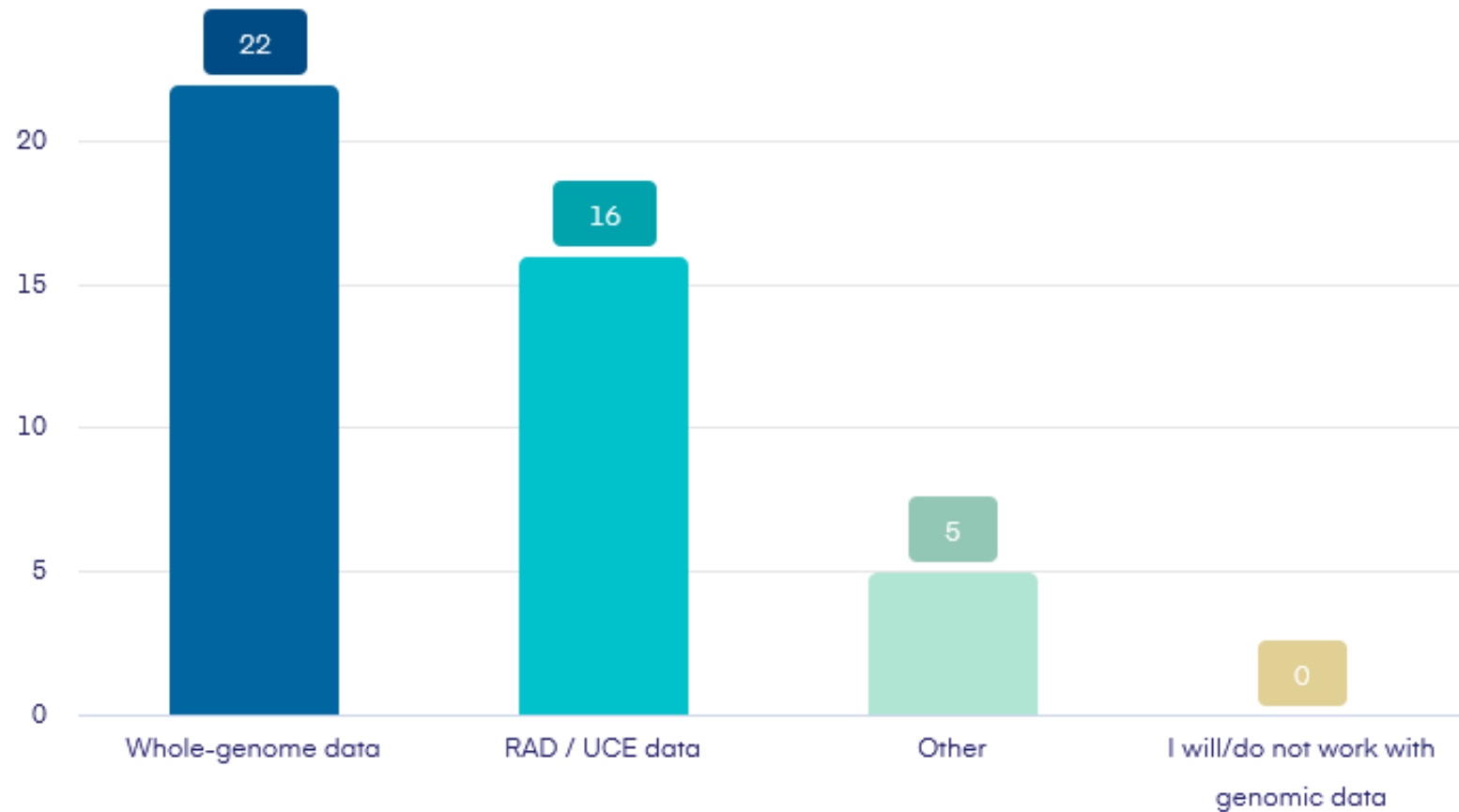
# PacBio HiFi reads



# Whole-genome sequencing (shotgun sequencing)

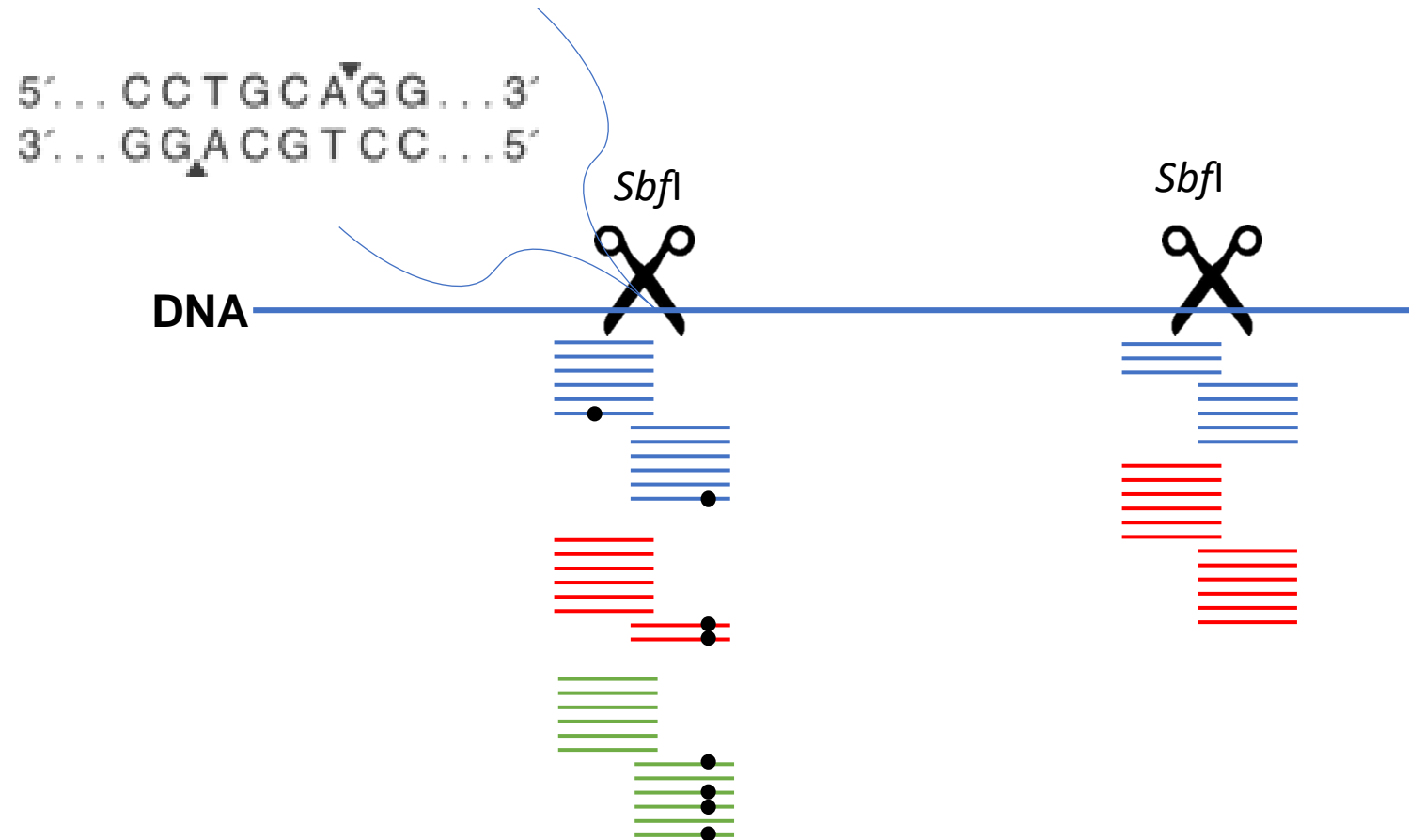


# DNA preparation methods



# RAD sequencing

Restriction **A**ssociated **D**NA sequencing



# Trade-offs: Splitting reads (i.e. costs) among:

- Number of sites to sequence
- Number of samples
- Sequencing depth
- Example: 1 HiSeq2500 flow cell  
~250 mio read pairs of 125 bp each -> 75 Gb data
  - 5 whole-genomes of a species with 1 Gb genome size at 15x coverage
  - 50 whole-genomes of a species with 500 Mb genome size at 3x coverage
  - 30 Mbp sequenced for 100 samples with a reduced-representation technique at a sequencing depth of 25

# Considerations in choosing the library preparation and sequencing techniques

- Research question and planned analyses
- Genome size
- Availability & quality of reference genome (no ref genome -> not wgs)
- Available budget
- Number of samples to sequence (tradeoff with sequencing depth)
- Amounts of DNA available
- Sequencing depth aimed at
- SNP density required
- Divergence between samples
- Heterozygosity of samples
- Phase required
- Accuracy of each single position (if high needed, avoid PCR-based methods)
- Importance of annotations
- Neutral dataset or specific regions wanted



# Fastq format

```
@HWUSI-EAS611:34:6669YAAXX:1:1:5069:1159 1:N:0:  
TCGATAATACCGTTTTTTTCCGTTTGATGTTGATACCATT  
+  
IIHIIHHIIIIIIIIIIIIIIIIIIIIIIIIHIIIHIIIIII  
@HWUSI-EAS611:34:6669YAAXX:1:1:5243:1158 1:N:0:  
TATCTGTAGATTTACACAGACTCAAATGTAAATATGCAGAG  
+  
DF=DBD<BBFGGGGGGGGBD@GGGD4@CA3CGG>DDD:D,B  
@HWUSI-EAS611:34:6669YAAXX:1:1:5266:1162 1:N:0:  
GGAGGAAGTATCACTTCCTTGCCTGCCTCCTCTGGGGCCT  
+  
:GBGGGGGGGGGGDDGGDEDGGDGGGGDHHDHGHHGBGG:GG
```

Header (must start with @)

Base calls (sequence)

Quality scores

# Quality scores

```
@EAS139:136:FC706VJ:2:2104:15343:197393 1:Y:18:ATCACG
CCGTCAATTCATTAGTTTTAACCCTTGCGGCCGTACTCCCCAGGCGGT
+
AAAAAAAAAAAAA:9@:::??@@::FFAAAAACCAA:::BB@@?A?
```

ASCII encoding

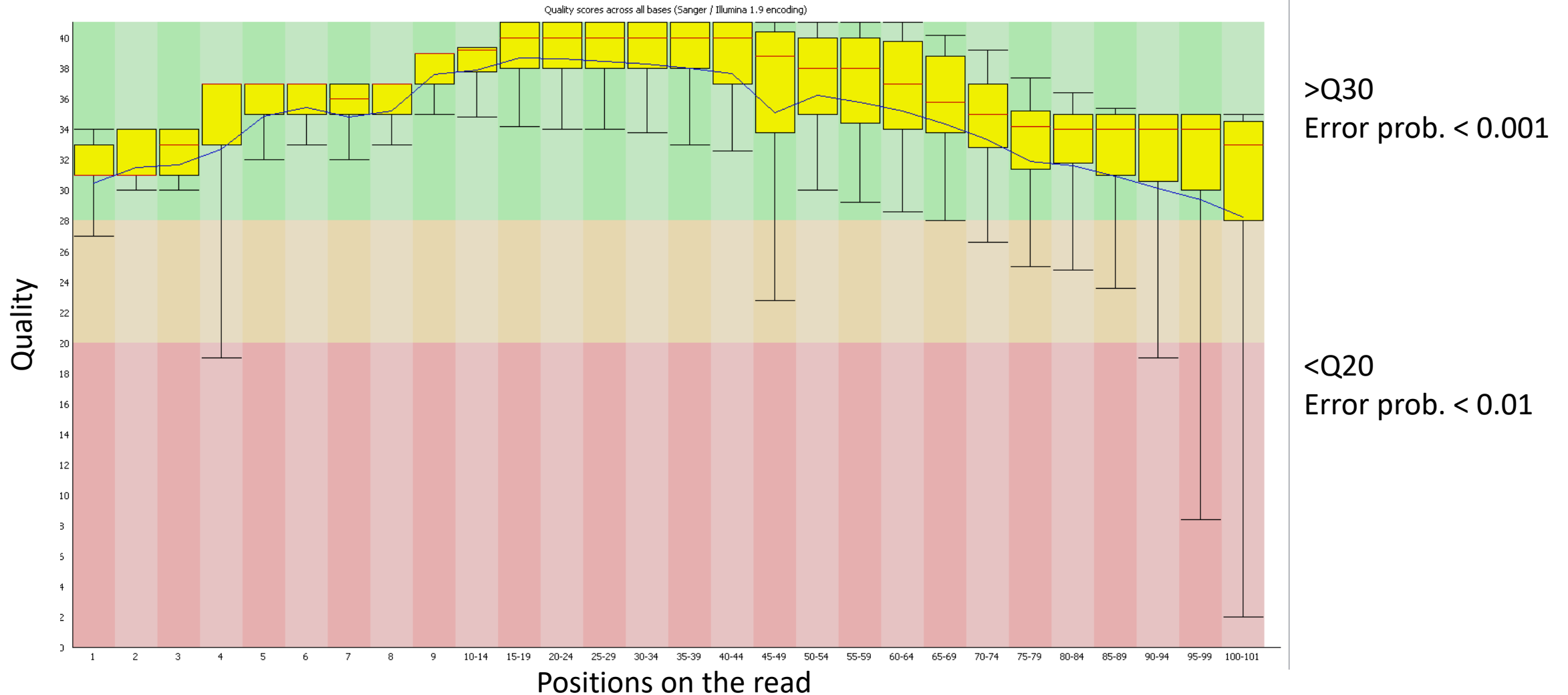
40: @	90: Z	141: a
41: A	91: [	142: b
42: B	92: \	143: c
43: C	93: ]	144: d
44: D	94: ^	145: e
45: E	95: _	146: f
... : ...	... : ...	... : ...

$$\text{Phred} = -10 \log_{10} p$$

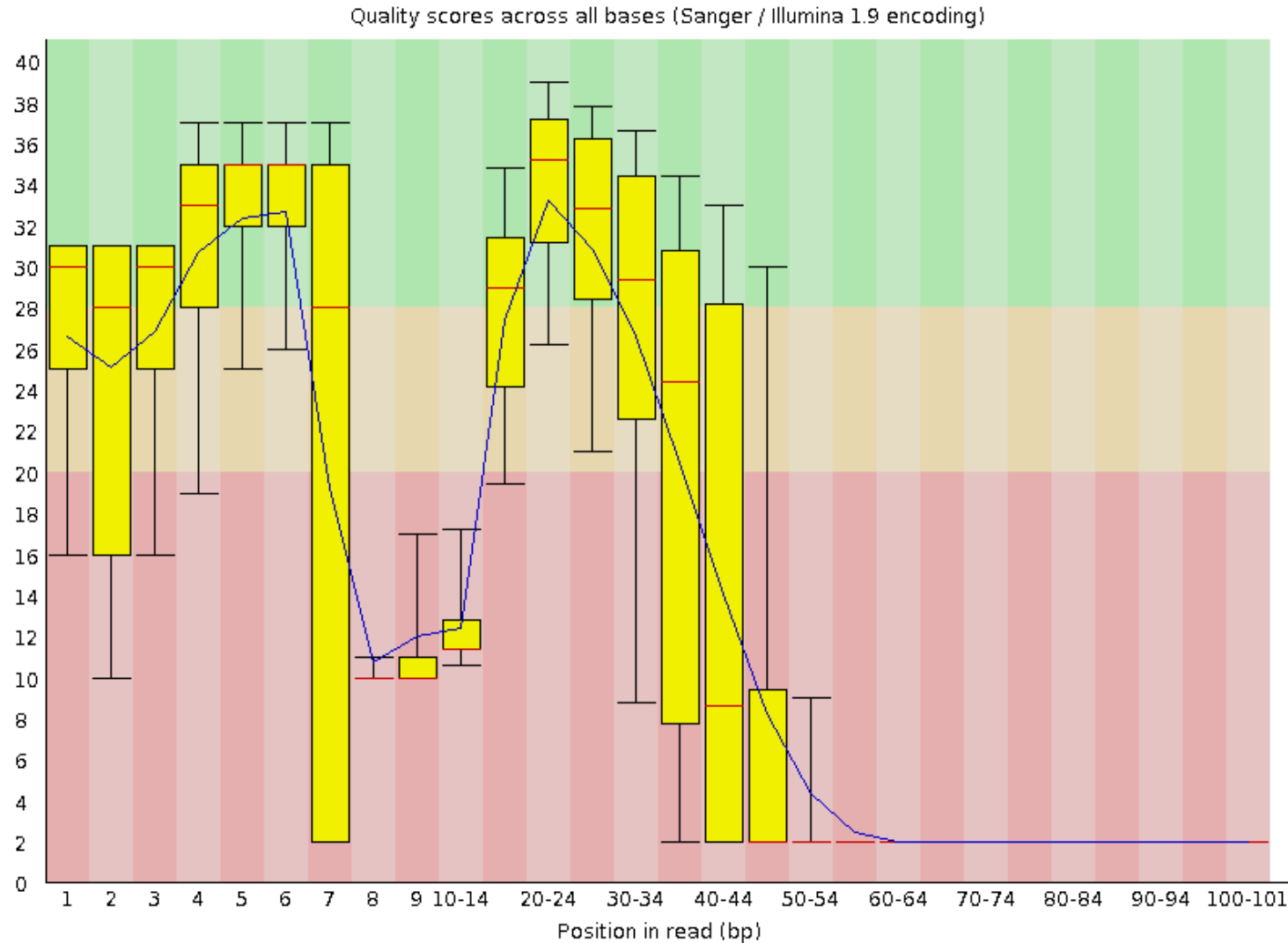
$p$  = Probability call is incorrect

Quality Score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%
40	1 in 10000	99.99%

# FastQC: Quality across bases (good example)



# FastQC: Quality across bases (bad example)



Let's have a look at the first few sequences and check the sequencing quality with fastqc