

# Introduction to Next-Generation Sequencing

Joanna Krupka (updated by Junfan Huang)

CRUK Summer School in Bioinformatics

Cambridge, July 2021



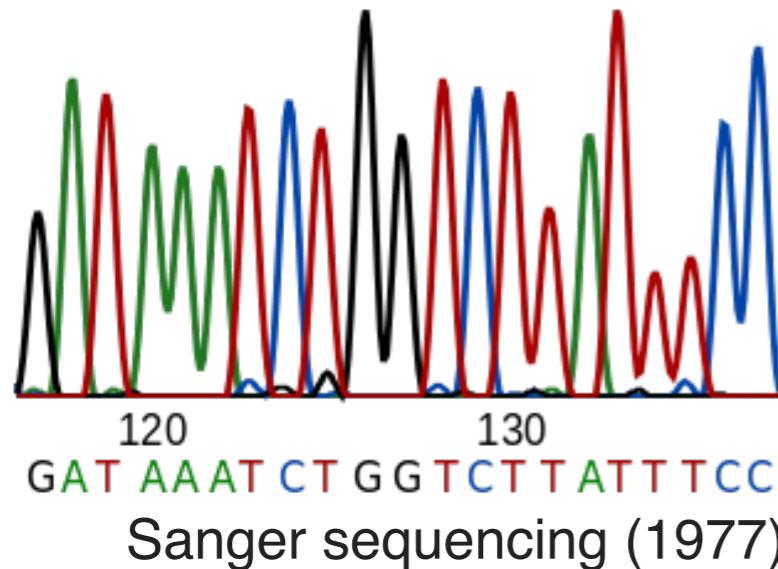
CANCER  
RESEARCH  
UK

MRC | Cancer  
Unit



UNIVERSITY OF  
CAMBRIDGE

# Brave New World of Next Generation Sequencing



## Human Genome Project 1990 - 2006

### DNA Sequencing Technologies Key to the Human Genome Project

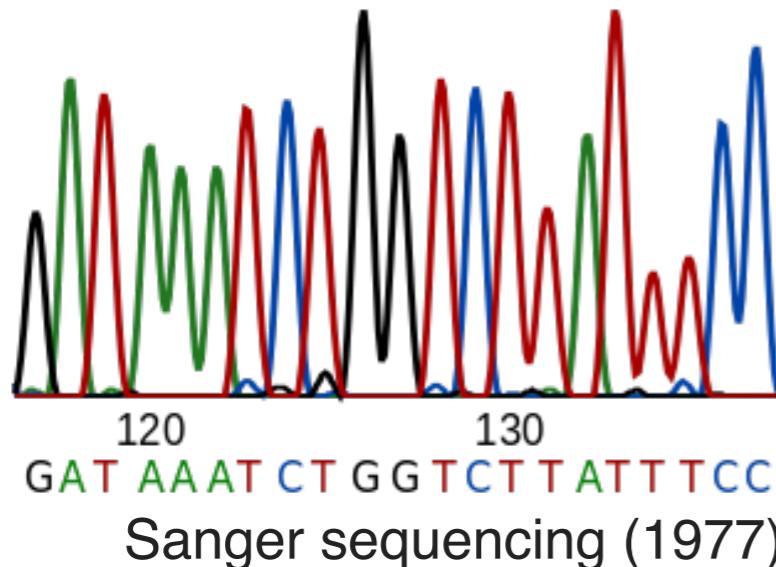
By: Heidi Chial, Ph.D. (*Write Science Right*) © 2008 Nature Education

Citation: Chial, H. (2008) DNA sequencing technologies key to the Human Genome Project. *Nature Education* 1(1):219



Adapted from Joanna Krupka's slides

# Brave New World of Next Generation Sequencing



presenting  
clinically applications  
genetic novo illumina  
human testing rolled targeted  
nanopore data cancer  
using whole today technology  
ecseq new early booth cells  
gs analysis complex  
genomics  
expression genome  
across dna great gene  
cell rna one england reveals  
work visit genomics live health  
van genomic difference  
crisprcas research voor types  
improve genes tests good  
support nhs service  
life bioinformatics  
radygenomics

## Human Genome Project 1990 - 2006

### DNA Sequencing Technologies Key to the Human Genome Project

By: Heidi Chial, Ph.D. (*Write Science Right*) © 2008 Nature Education

Citation: Chial, H. (2008) DNA sequencing technologies key to the Human Genome Project. *Nature Education* 1(1):219



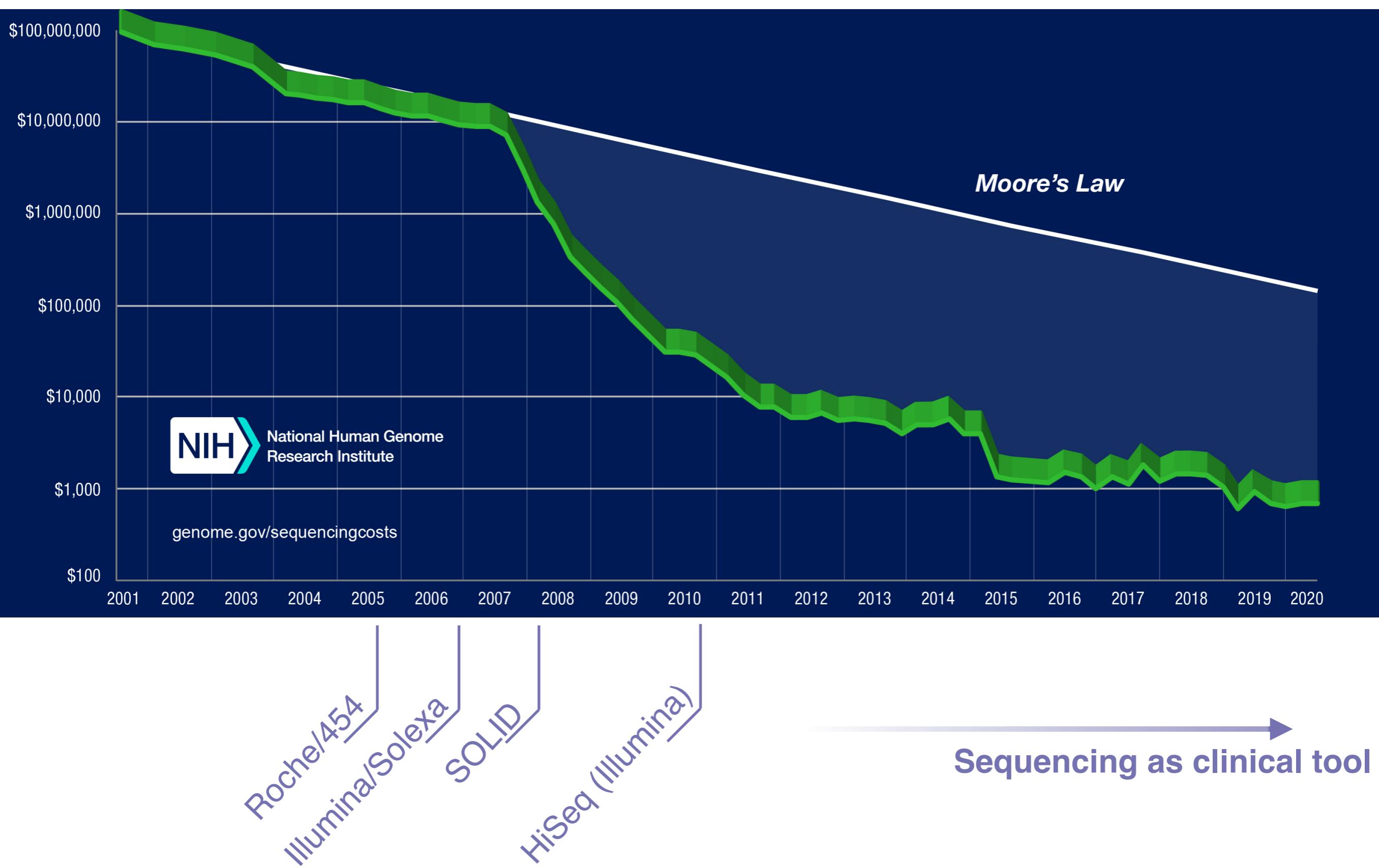
## Next Generation Sequencing mid 2000–present

= high-throughput sequencing

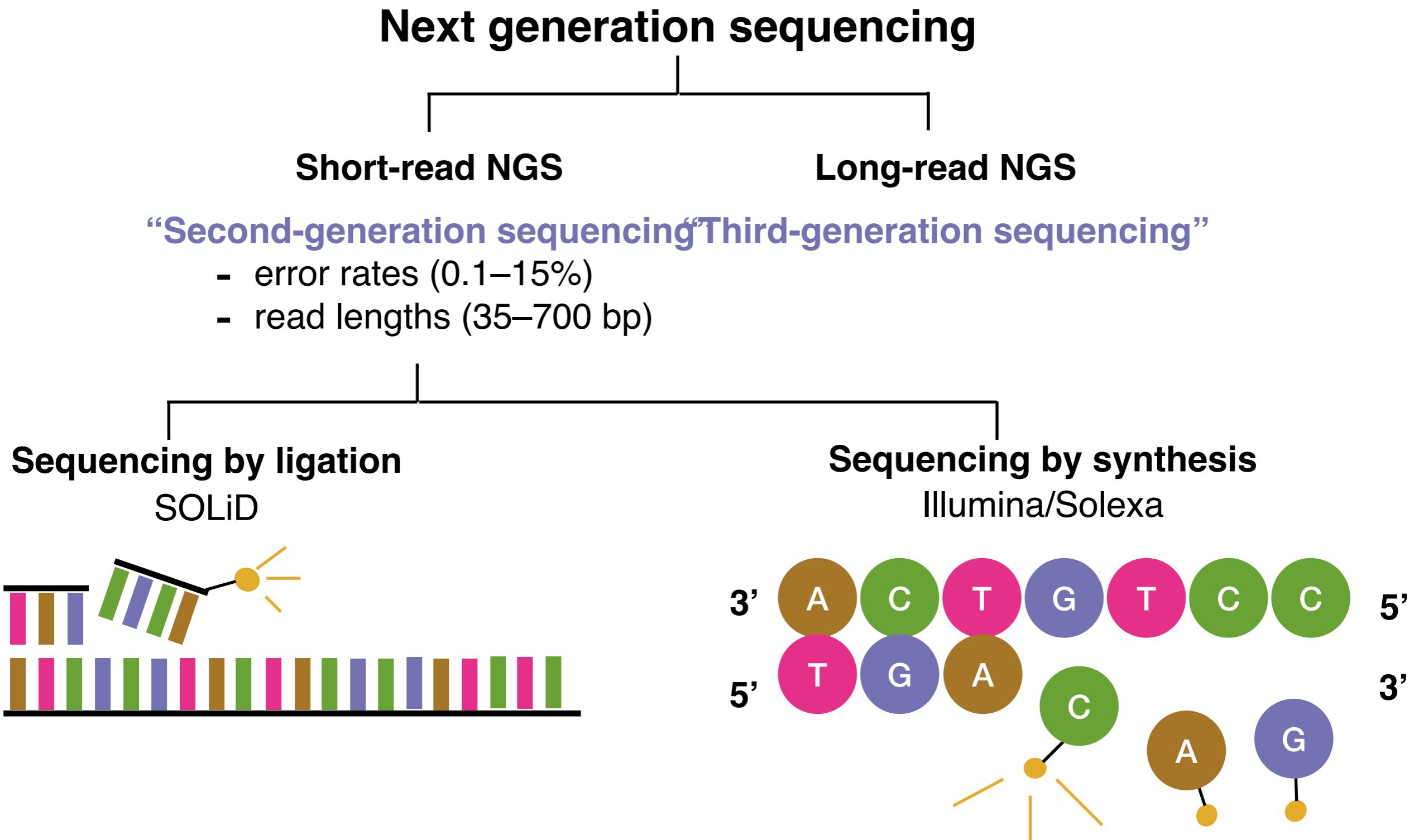
quicker and cheaper parallel sequencing of  
DNA and RNA

Adapted from Joanna Krupka's slides

# Cost of sequencing of human genome

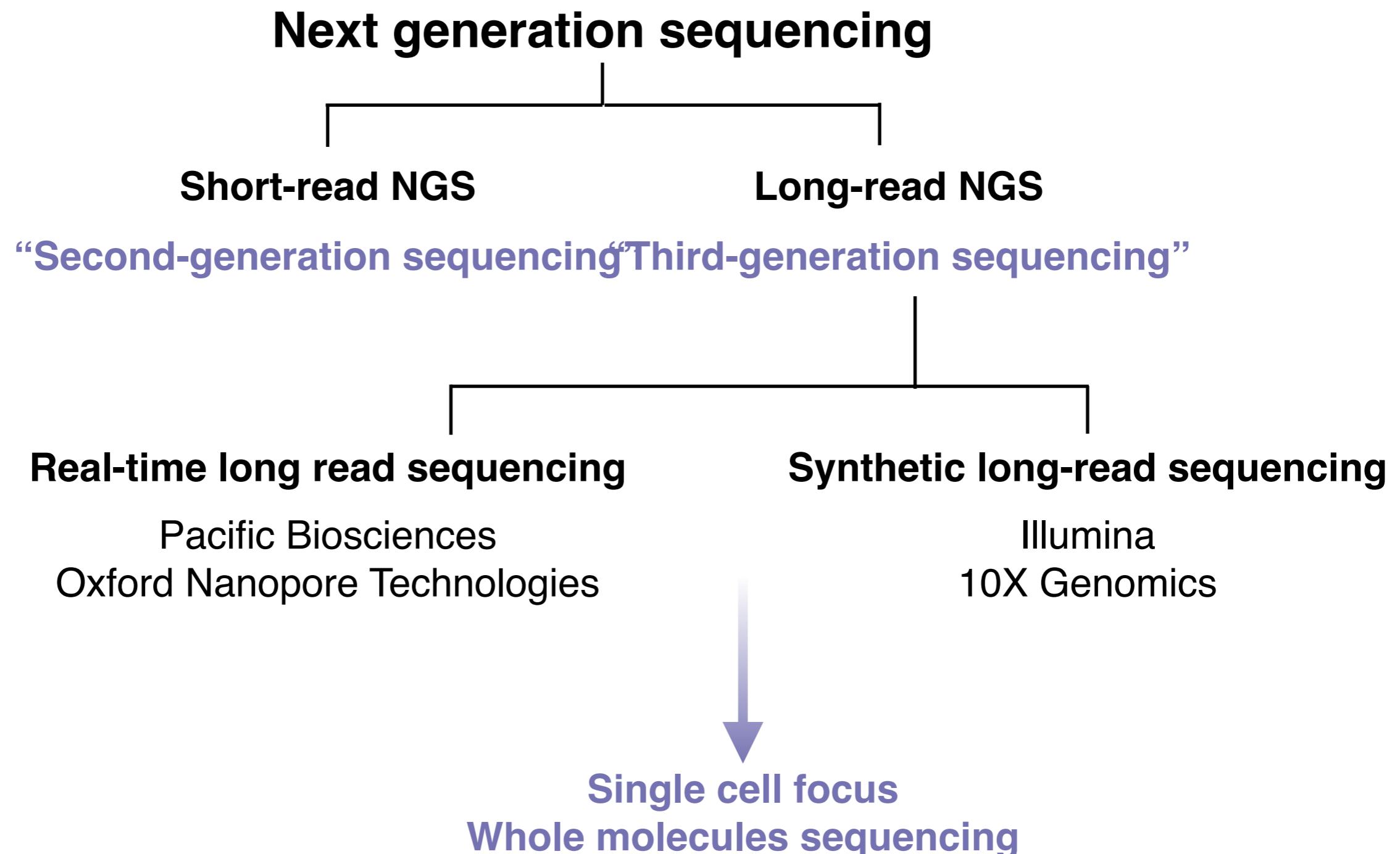


# Next generation sequencing technologies and limitations



Goodwin, S., McPherson, J. D., & McCombie, W. R. (2016). Coming of age: Ten years of next-generation sequencing technologies. *Nature Reviews Genetics*, 17(6), 333–351.

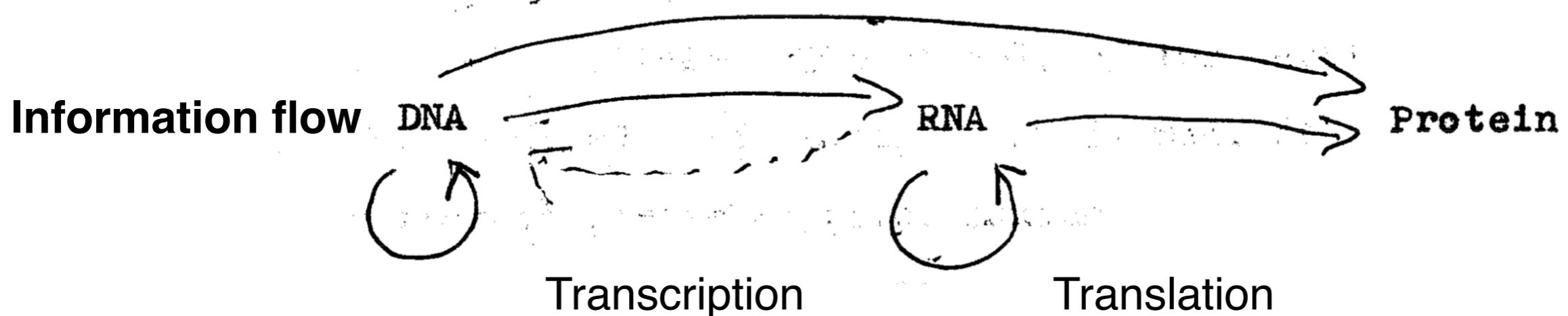
# Next generation sequencing technologies and limitations



Goodwin, S., McPherson, J. D., & McCombie, W. R. (2016). Coming of age: Ten years of next-generation sequencing technologies. *Nature Reviews Genetics*, 17(6), 333–351.

# Sequencing techniques

## Central dogma of molecular biology (Crick F. 1958)



## Whole genome sequencing

Whole exome sequencing

HiC-Seq

ChIP-Seq

ATAC-Seq

...

DNA

RNA

Adapted from Joanna Krupka's slides

scRNA-Seq

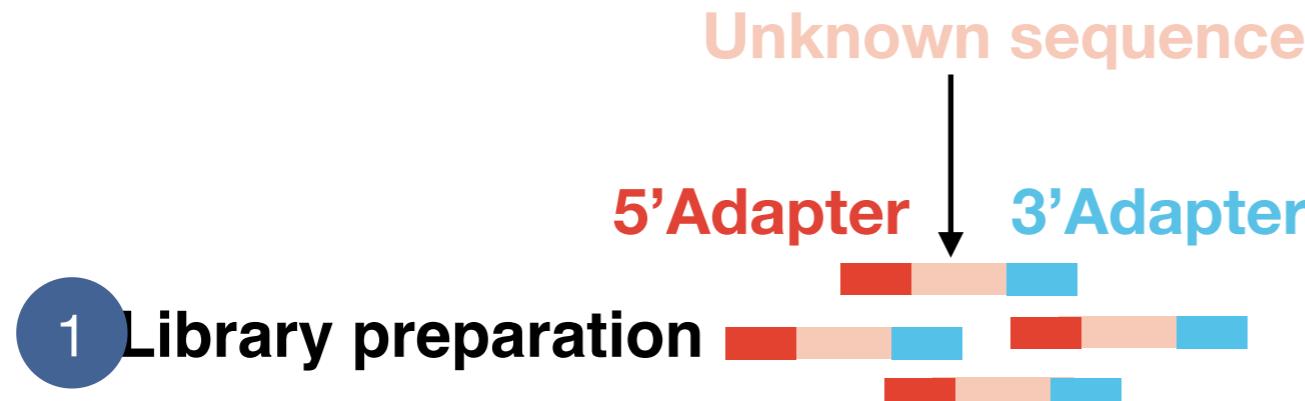
RNA-Seq

Ribo-Seq

SLAM-Seq

...

# Illumina sequencing by synthesis



**NOTE 1:** High quality material needed for high quality experiment!

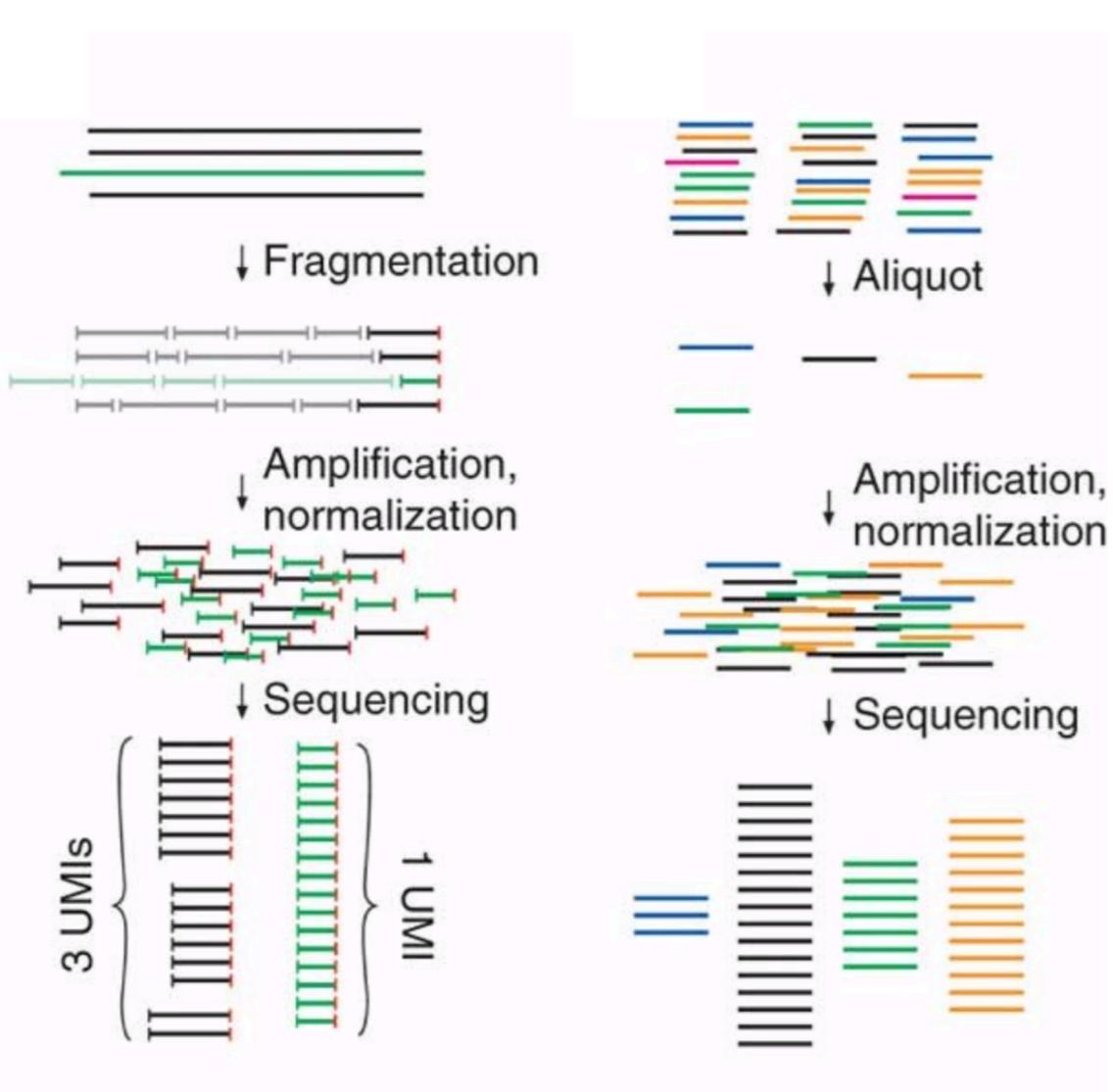
**NOTE 2:** Final step of library preparation is amplification. Some products are preferentially amplified, which introduces **library amplification bias**.

- Fewer cycles - fewer bias
- **Unique molecular identifiers:** oligonucleotides labels to identify duplicated fragments

Adapted from Joanna Krupka's slides

Goodwin, S., McPherson, J. D., & McCombie, W. R. (2016). Coming of age: Ten years of next-generation sequencing technologies. *Nature Reviews Genetics*, 17(6), 333–351.

# Unique molecular identifiers (UMIs)



4 exactly same fragments: unique or duplicates?

4 different UMIs



UNIQUE!

4 same UMIs

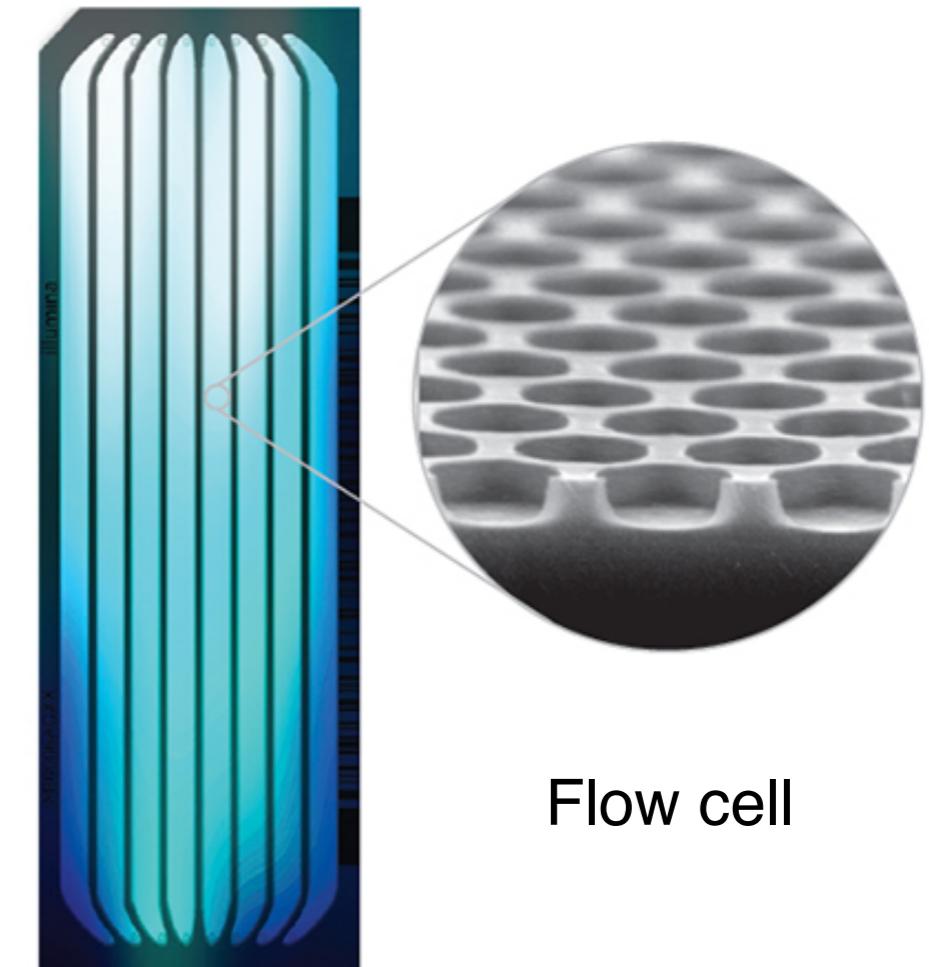


DUPPLICATES!  
!

UMIs help to identify library amplification bias and quantify unique fragments  
(identical fragments with the same UMIs are likely to be duplicates)

# Illumina sequencing by synthesis

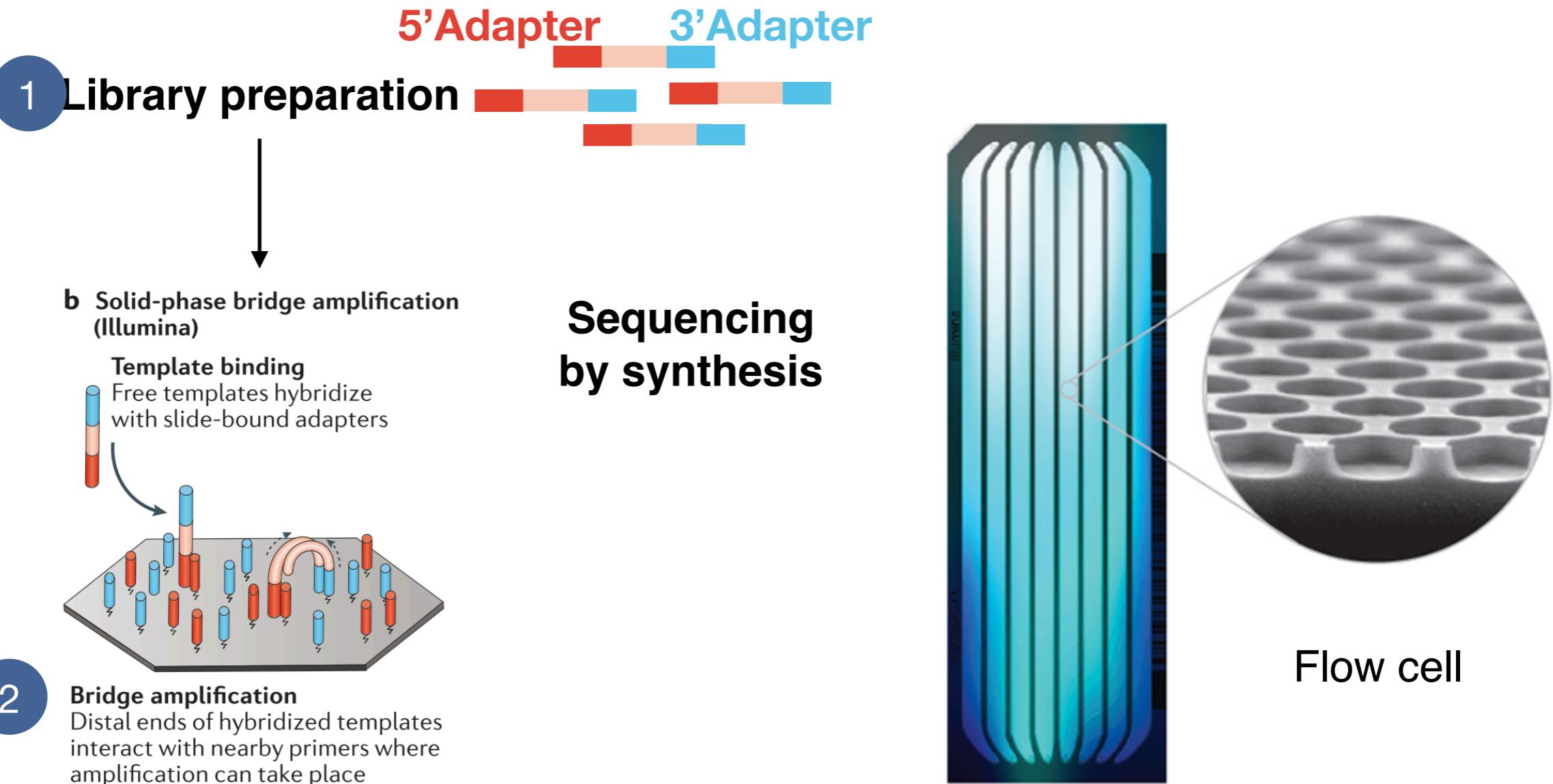
Based on the Solexa technology developed by **Shankar Balasubramanian** and **David Klenerman** at the University of Cambridge (1998)



Flow cell

# Illumina sequencing by synthesis

Based on the Solexa technology developed by **Shankar Balasubramanian** and **David Klenerman** at the University of Cambridge (1998)

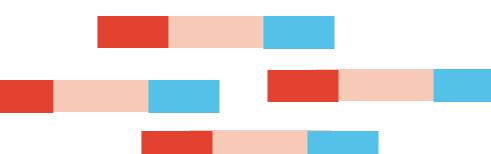


Goodwin, S., McPherson, J. D., & McCombie, W. R. (2016). Coming of age: Ten years of next-generation sequencing technologies. *Nature Reviews Genetics*, 17(6), 333–351.

# Illumina sequencing by synthesis

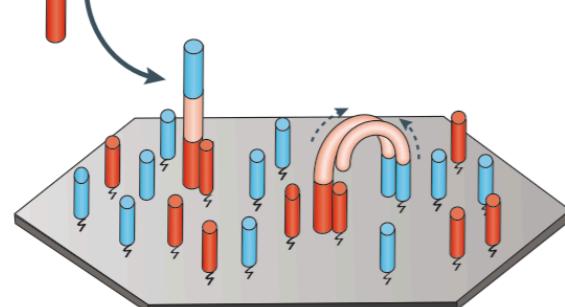
Based on the Solexa technology developed by **Shankar Balasubramanian** and **David Klenerman** at the University of Cambridge (1998)

## 1 Library preparation

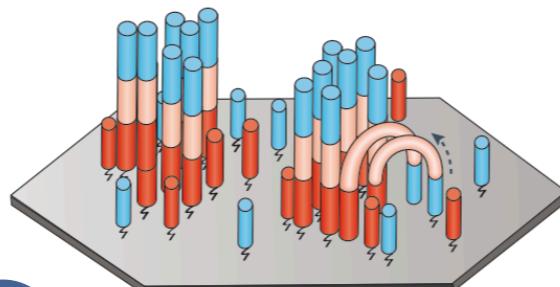


## b Solid-phase bridge amplification (Illumina)

Template binding  
Free templates hybridize with slide-bound adapters



## Sequencing by synthesis



## 2

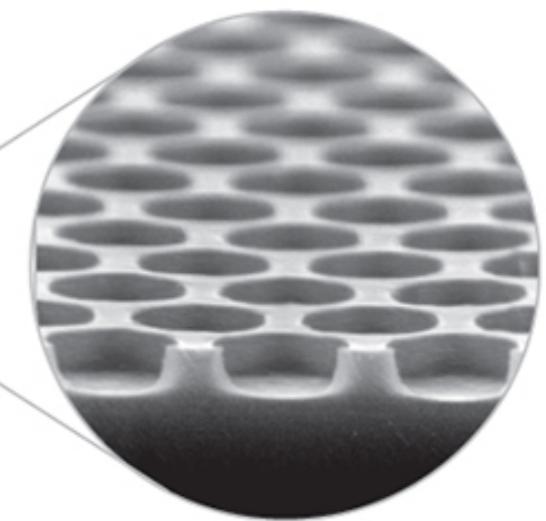
### Bridge amplification

Distal ends of hybridized templates interact with nearby primers where amplification can take place

## 3

### Cluster generation

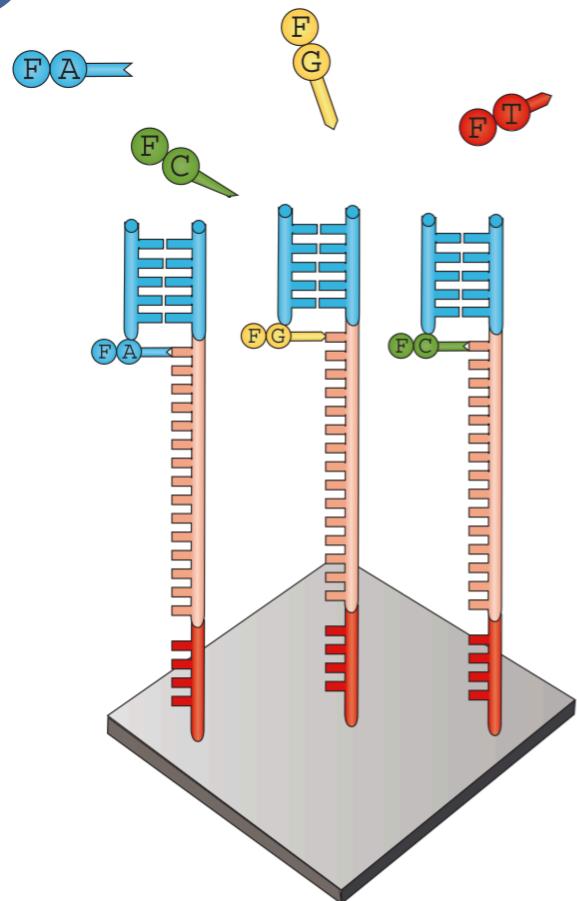
After several rounds of amplification, 100–200 million clonal clusters are formed



Flow cell

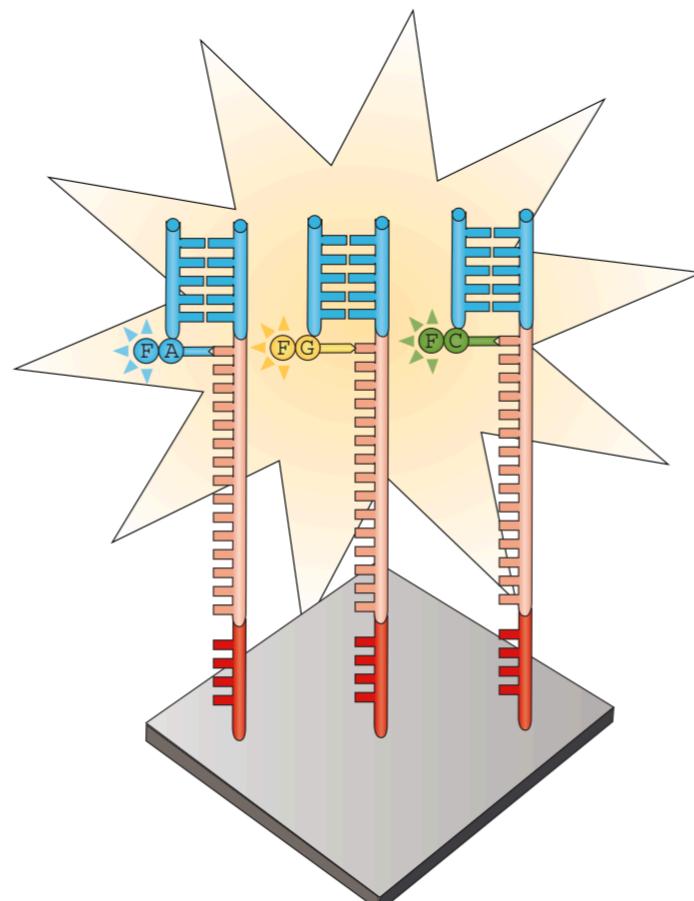
# Illumina sequencing by synthesis

## 4 Sequencing using reversible terminators



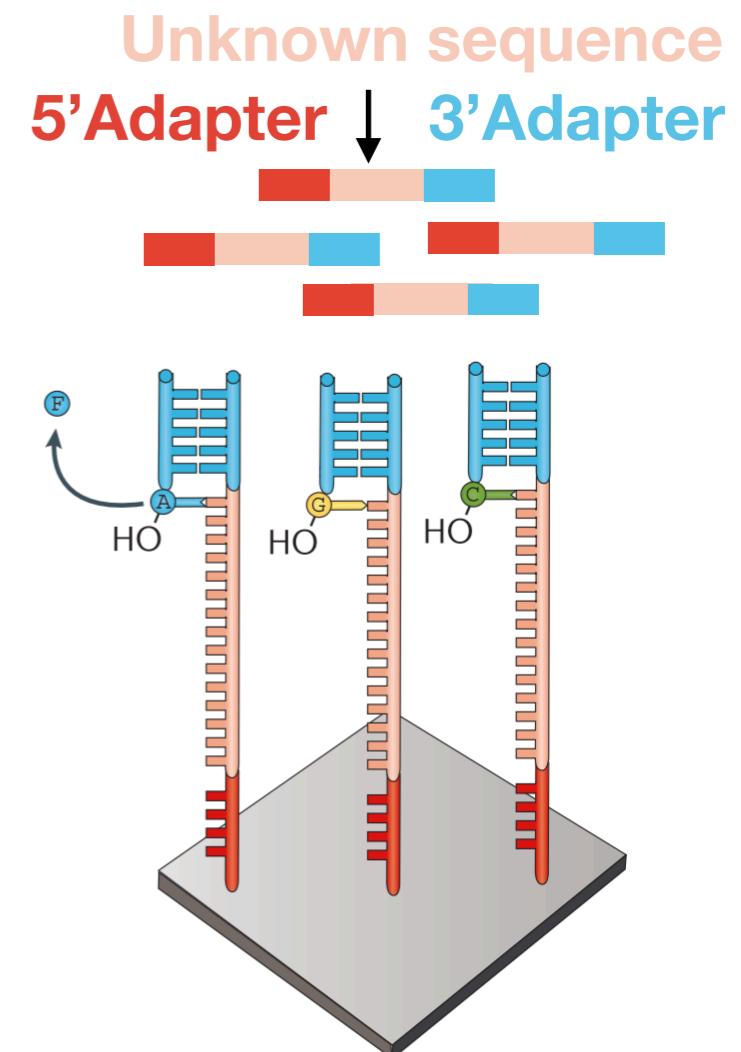
### Nucleotide addition

Fluorophore-labelled, terminally blocked nucleotides hybridize to complementary base. Each cluster on a slide can incorporate a different base.



### Imaging

Slides are imaged with either two or four laser channels. Each cluster emits a colour corresponding to the base incorporated during this cycle.

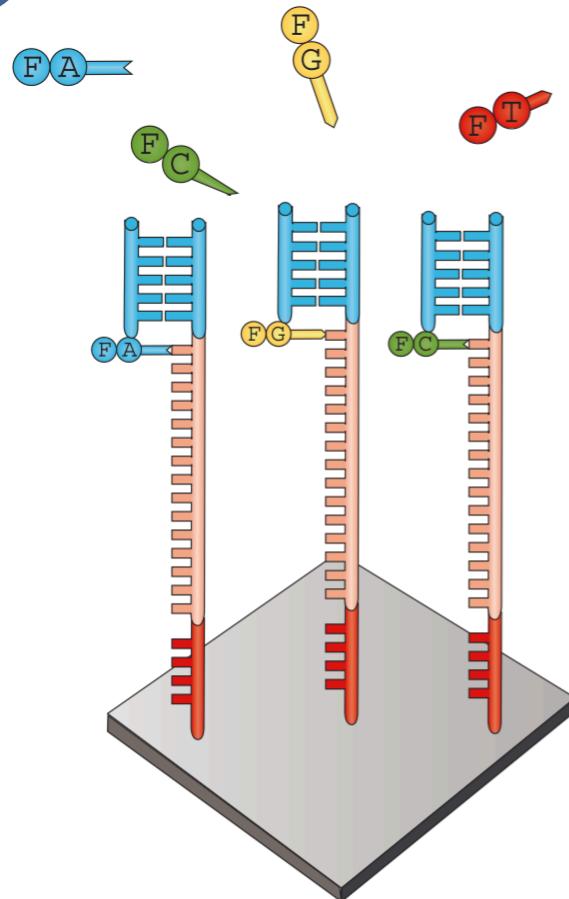


### Cleavage

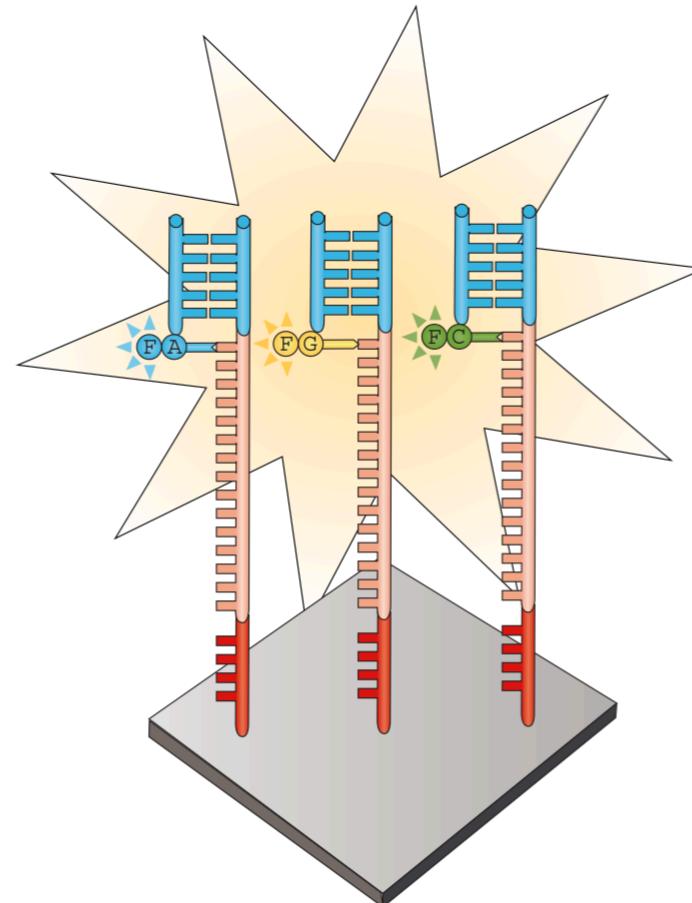
Fluorophores are cleaved and washed from flow cells and the 3'-OH group is regenerated. A new cycle begins with the addition of new nucleotides.

# Illumina sequencing by synthesis

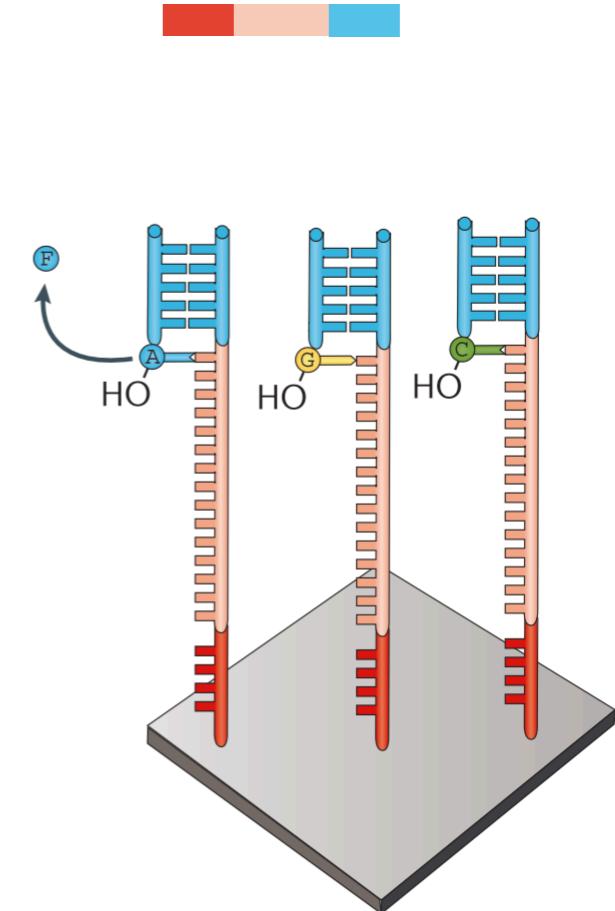
## 4 Sequencing using reversible terminators



**Nucleotide addition**  
Fluorophore-labelled, terminally blocked nucleotides hybridize to complementary base. Each cluster on a slide can incorporate a different base.



**Imaging**  
Slides are imaged with either two or four laser channels. Each cluster emits a colour corresponding to the base incorporated during this cycle.



**Cleavage**  
Fluorophores are cleaved and washed from flow cells and the 3'-OH group is regenerated. A new cycle begins with the addition of new nucleotides.

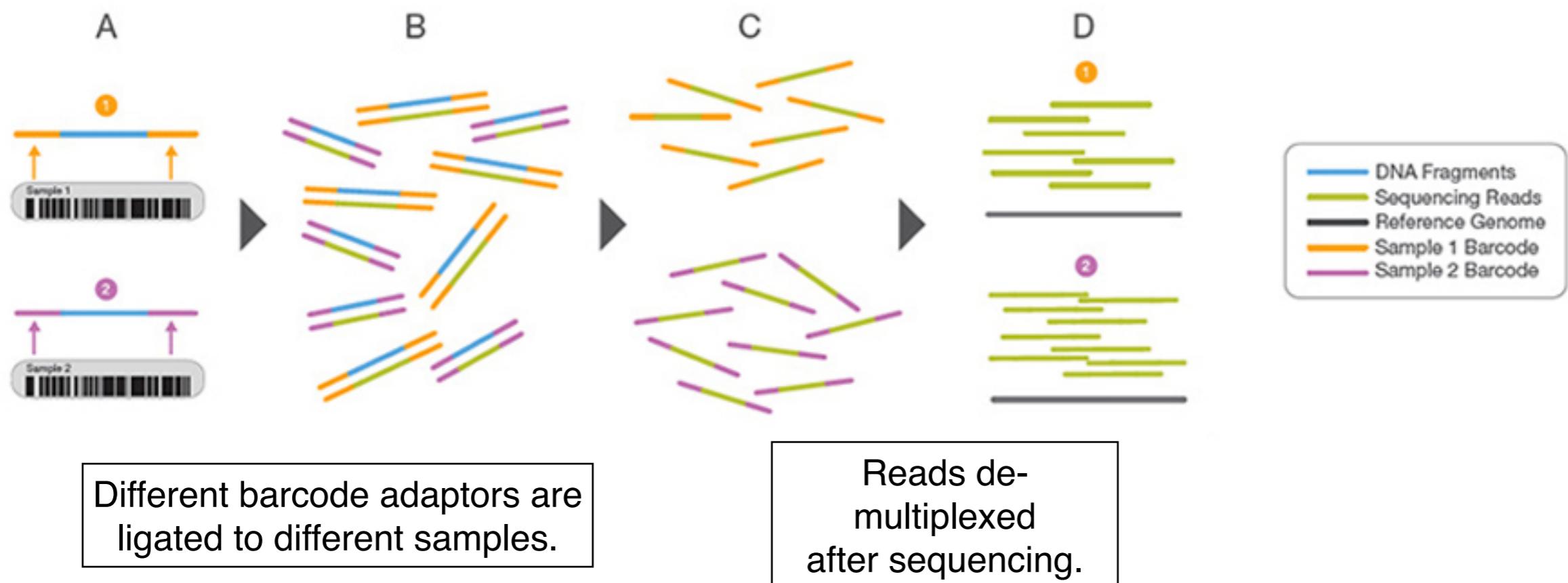
## 5 Output: sequence saved in FASTQ format

## 6 Bioinformatic analysis: quality check, alignment and data analysis

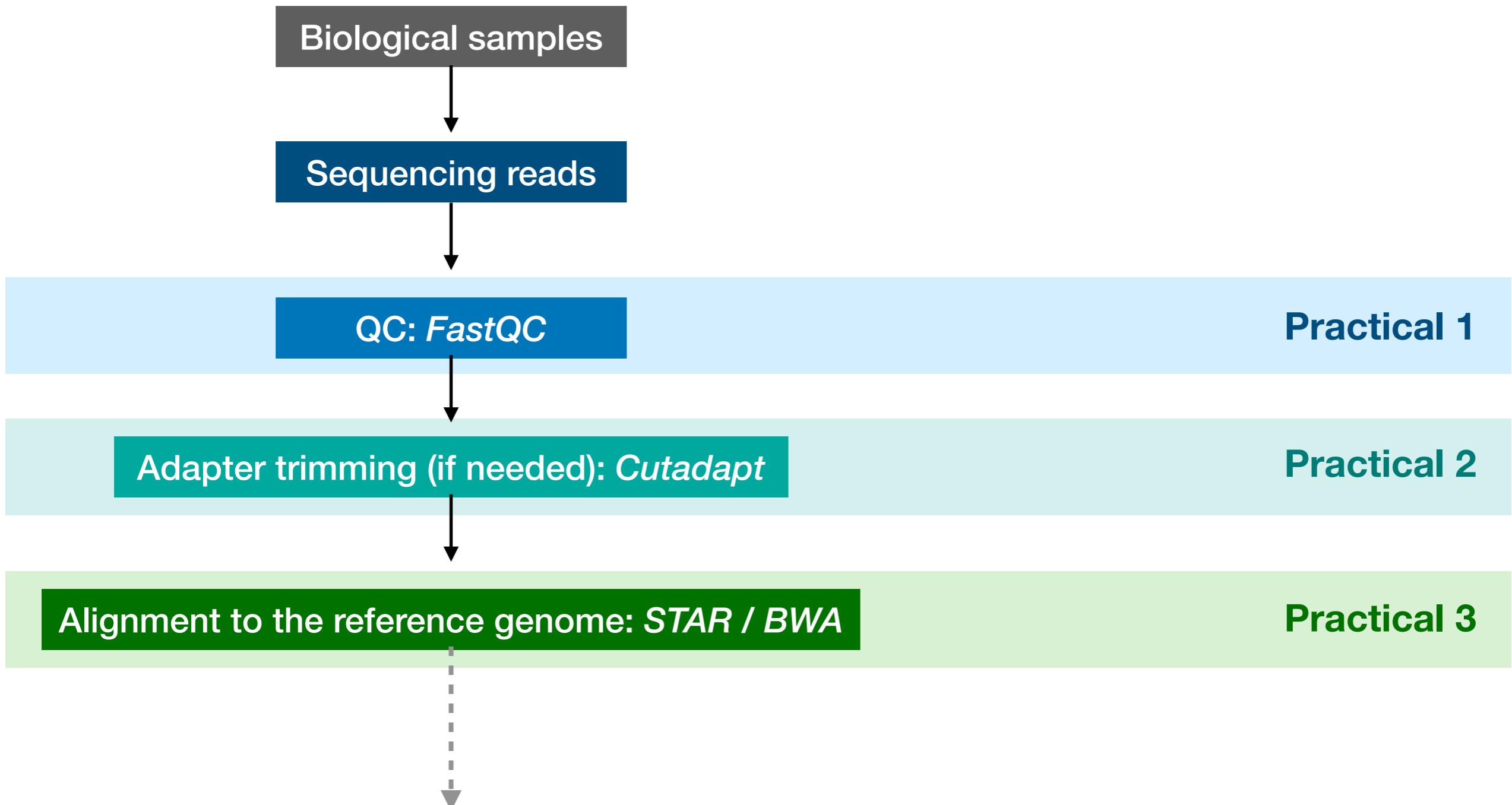
Goodwin, S., McPherson, J. D., & McCombie, W. R. (2016). Coming of age: Ten years of next-generation sequencing technologies. *Nature Reviews Genetics*, 17(6), 333–351.

# Multiplexing

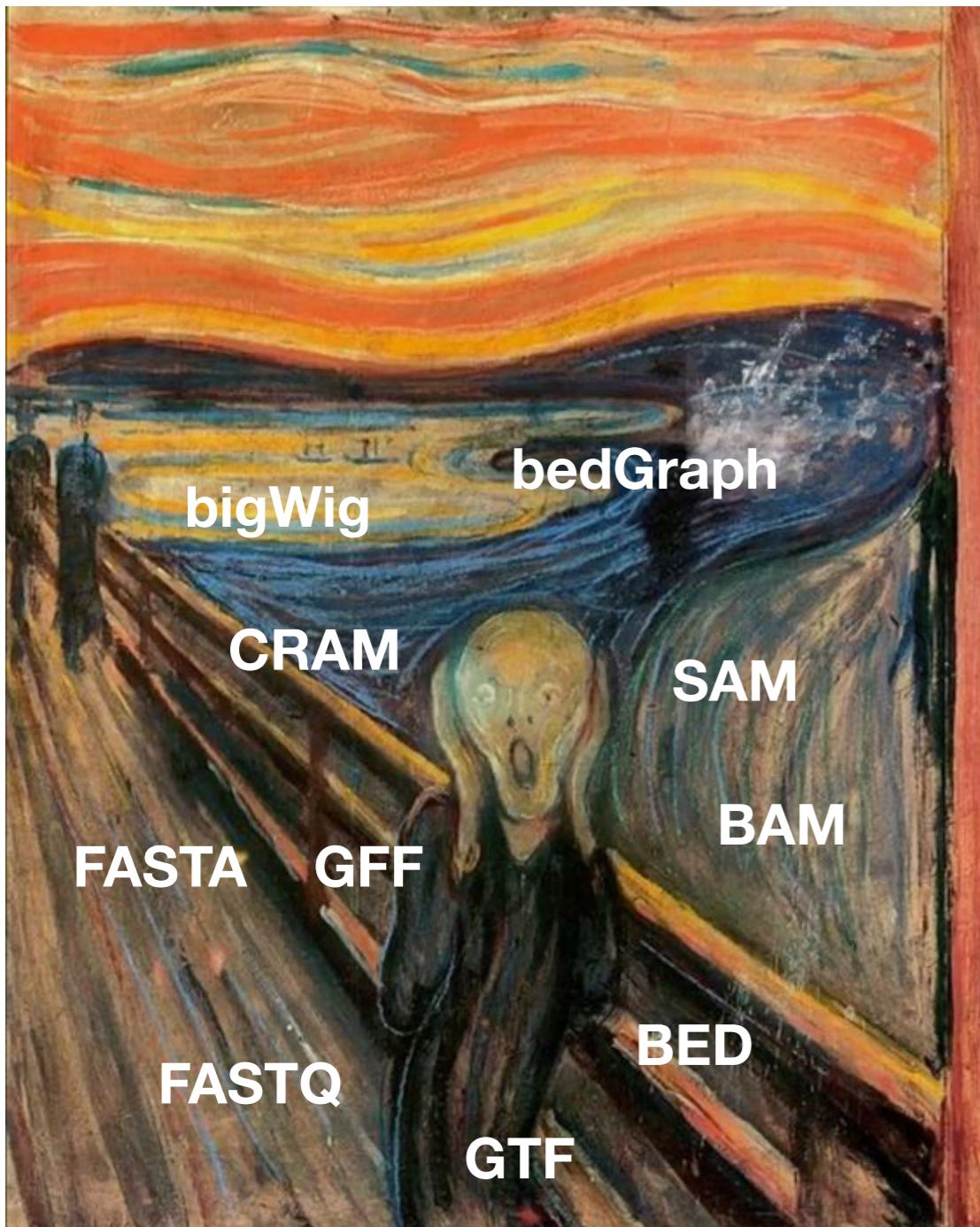
- Multiplexing gives the ability to sequence multiple samples at the same time
- Blocks against possible technical bias caused by differences between flow cell lanes
- Useful when sequencing small genomes or specific genomic regions.



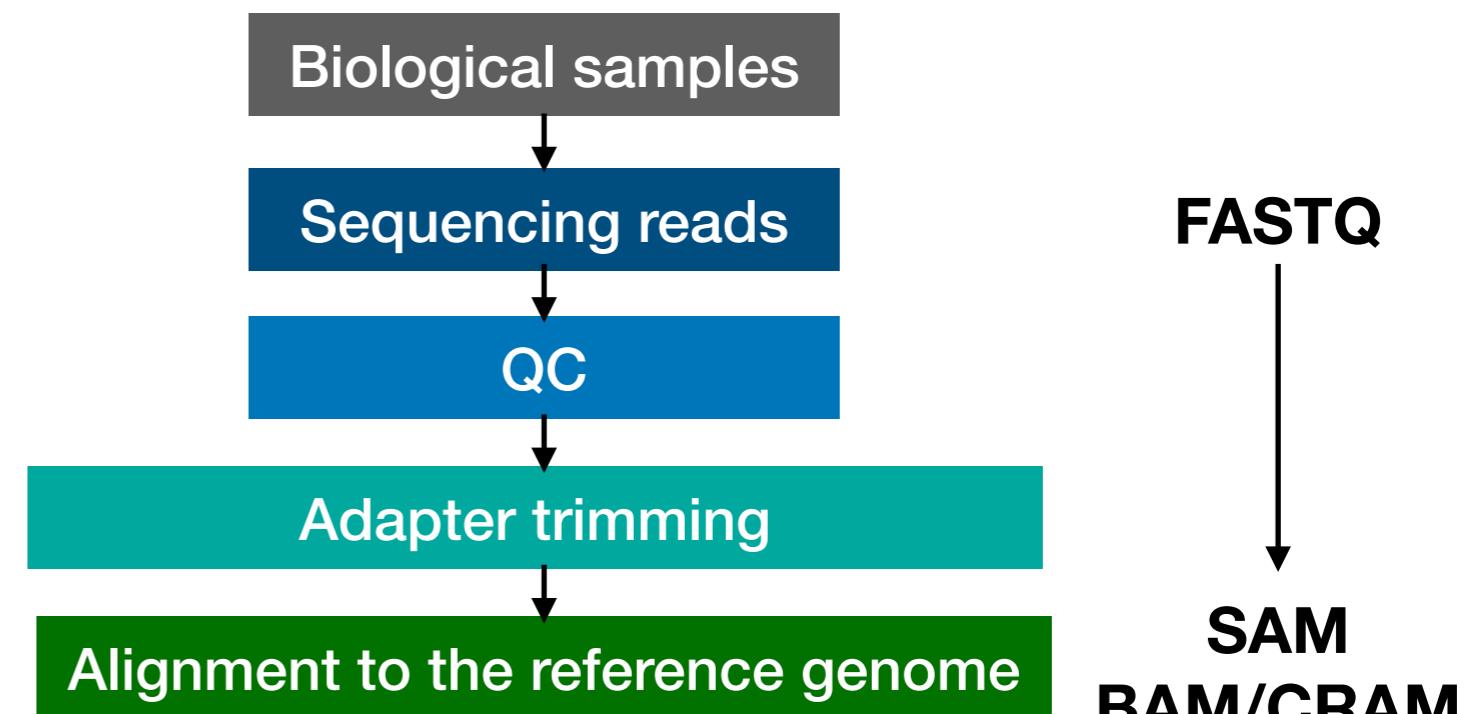
# Workflow for today



# Common file formats: why so many?



**Different formats - different informations**



# Nucleotide/peptide sequences: FASTA

A sequence in FASTA format consists of:

**1st line** starting with “>” followed by the sequence name

# **2nd line** with the sequence itself

A single FASTA file may contain > 1 sequence

# Unaligned sequence: FASTQ

Unaligned sequence (reads) files generated from NGS machines

A sequence in FASTQ format consists of:

**1st line** starting with "@" followed by the read identifier.

**2nd line** with the sequence itself.

**3rd line** “+”

**4th line** Quality scores encoded as ASCII characters

```
@K00359:71:HJJL7BBXX:3:1101:1996:1508 1:N:0:ATCACG
AAAATTCCAAGCTGGTTCAACAGTACTTGTTCAGAACAAAGAAATG
+
AAAFFJJJJJJFJJ<J<FJJJJJJJJJJJJFJJFJJJJFJJFJJJJ<
@K00359:71:HJJL7BBXX:3:1101:2240:1508 1:N:0:ATCACG
GTAAAGGATGCGTAGGGATGGGAGGGCGATGAGGACTAGGATGATGGCGG
+
AAFFFJJJJJJF<J7JJFJJJJJJFFFJFJJJJJJJJJJJJJJJJJJJJ
@K00359:71:HJJL7BBXX:3:1101:2402:1508 1:N:0:ATCACG
GTCGACCATGTGGGCAGAACCTTGATGTTGGATTCCAGCAGGACCTGTCC
+
AAFFFJJJJJJJJ<JJJJJJJJJJ<JFJJJJJJJJJJJJJJJJFJJJJJJ
@K00359:71:HJJL7BBXX:3:1101:2463:1508 1:N:0:ATCACG
ATGTGGTGTATGCATGGGGTAGTCCGAGTAACGTCGGGCATTCCGGAT
+
AAAFFFFJJJJJJJJJJJJFJJJJJJJJJJJJJJJJFJ7JJJJJJJJJJ
```

## Unaligned sequence: FASTQ

FASTQ header decoded (Illumina example):

Machine ID	Run	Flow cell ID	Lane	Tile	Tile coordinates	Read	Barcode
					X	Y	Idx Filter
@K00359:71:HJJL7BBXX:3:1101:1996:1508							
AAAATTCCAAGCTGGTTAACAGTACTTGTTCCAGAACAAAGAAATG							
+							
AAAFFJJJJJJFJJ<J<FJJJJJJJJJJJJFJJFJJFFJJFJJJJJ<							

## Unaligned sequence: FASTQ

Quality scores come after the "+" line

Quality  $Q$  is proportional to -log<sub>10</sub> probability of sequence base being wrong  $e$

$$Q = -10 \cdot \log_{10}(e)$$

```
@K00359:71:HJJL7BBXX:3:1101:1996:1508 1:N:0:ATCACG  
AAAATTCCAAGCTGGTTAACACAGTACTTGTTCCAGAACAAAGAAATG  
+
```

```
AAAFFJJJJJJFJJ<J<FJJJJJJJJJJJJJJFJJFJJJJFFJFJJJJJ<
```

Encoded in ASCII to save space:

Quality encoding:	!"#\$%&'()*+,.-./0123456789:;=>?@ABCDEFGHI
Quality score:	0.....10.....20.....30.....40

Used in quality assessment and downstream analysis

# SAM - Sequence Alignment Map

Unaligned sequence files generated from NGS machines are mapped to a reference genome to produce aligned sequence:

**FASTQ(unaligned sequences) → SAM (aligned sequences)**

# FASTA + quality

# FASTQ + location

SAM:

- Standard format for aligned sequence data
  - Recognised by majority of software and browsers
  - Starts with a header section followed by alignment information as tab separated lines for each read.

*Header section*

```
@HD     VN:1.3      SO:coordinate  
@SQ     SN:contigA   LN:443  
@SQ     SN:contigB   LN:1493  
@SO     SN:contigC   LN:328
```

### *Tab-delimited read alignment information lines*

```
readID43GYAX15:7:1:1202:19894/1      256      contig43      613960      1      65M      *      0      0  
CCAGCGCGAACGAAATCCGATGCGTCTGGTCGTTGCACGGAACGGCGGCGGTGTGATGCACGGC      EDDEEDEE=EE?DE??  
DDDBADEBEFFFDBEFFEBCBC=?BEEEE@=:?:?:?7?:8-6?7?@??#      AS:i:0      XS:i:0      XN:i:0      XM:i:0  
XO:i:0      XG:i:0      NM:i:0      MD:Z:65      YT:Z:UU
```

# SAM - Sequence Alignment Map

## SAM header

- Header lines start with '@'

@HD	VN:1.4	SO:coordinate
@SQ	SN:chr1	LN:248956422
@SQ	SN:chr2	LN:242193529
@SQ	SN:chr3	LN:198295559
@SQ	SN:chr4	LN:190214555
@SQ	SN:chr5	LN:181538259
@SQ	SN:chr6	LN:170805979
@SQ	SN:chr7	LN:159345973
@SQ	SN:chr8	LN:145138636
@SQ	SN:chr9	LN:138394717
@SQ	SN:chr10	LN:133797422
@SQ	SN:chr11	LN:135086622
@SQ	SN:chr12	LN:133275309
@SQ	SN:chr13	LN:114364328
@SQ	SN:chr14	LN:107043718
@SQ	SN:chr15	LN:101991189
@SQ	SN:chr16	LN:90338345
@SQ	SN:chr17	LN:83257441
@SQ	SN:chr18	LN:80373285
@SQ	SN:chr19	LN:58617616
@SQ	SN:chr20	LN:64444167
@SQ	SN:chr21	LN:46709983
@SQ	SN:chr22	LN:50818468
@SQ	SN:chrX	LN:156040895
@SQ	SN:chrY	LN:57227415
@SQ	SN:chrM	LN:16569

← **File-level metadata**  
VN: format version, SO: sorting order

← **Reference sequence dictionary**  
SN : name (eg. chr1), LN : length

Full format specification:

<https://samtools.github.io/hts-specs/SAMv1.pdf>

# SAM - Sequence Alignment Map

## Aligned reads

- Organised as tab-delimited text
  - Each alignment line has 11 mandatory fields for essential alignment information such as mapping position, and variable number of optional fields for flexible or aligner specific information.

# Read informations (as in FASTQ):

## QNAME: read ID

## SEQ: read sequence

## QUAL: read quality

# SAM - Sequence Alignment Map

# Aligned reads

- Organised as tab-delimited text
  - Each alignment line has 11 mandatory fields for essential alignment information such as mapping position, and variable number of optional fields for flexible or aligner specific information.

**RNAME:** reference seq name (eg. chromosome, transcript)

# CIGAR: summary of alignment (eg. insertion/deletion)

**POS:** position of 5' end of a read

## CIGAR string encoding:

**50M** - continuous match of 50 bases

**28M1D72M** - 28 bases continuously match, 1 deletion from reference, 72 base match

## Full format specification:

<https://samtools.github.io/hts-specs/SAMv1.pdf>

# SAM - Sequence Alignment Map

# Aligned reads

- Organised as tab-delimited text
  - Each alignment line has 11 mandatory fields for essential alignment information such as mapping position, and variable number of optional fields for flexible or aligner specific information.

**Bit flag** - TRUE/FALSE for pre-defined read criteria, like: is it paired? duplicate?

## Paired read position and insert size

K00359:71:HJJL7BBXX:3:1209:18436:10229

11S29M10S \* 0 0

TGCCTGGGAGGCCGGACCTTGGAGACTGTGTGGGGCCTGGCAC

AS:i:28 NM:i:0 MD:Z:29

chr1

chr1

16079

16079

255

255

# Mapping quality

## Flags explained:

<https://broadinstitute.github.io/picard/explain-flags.html>

# Compressed aligned sequences - BAM and CRAM format

SAM files can be large, so to save space people usually store some compressed versions of them instead:

## BAM

- Binary SAM file
- You also need to store an index file

## CRAM

- Another way to compress alignment files
- The compression is driven by the reference the sequence data is aligned to, so it is very important that the exact same reference sequence is used for compression and decompression
- Typically 40-50% space saving compared to BAM files
- Full compatibility with BAM files
- For further information: <http://samtools.github.io/hts-specs/>

# **10 min break!**