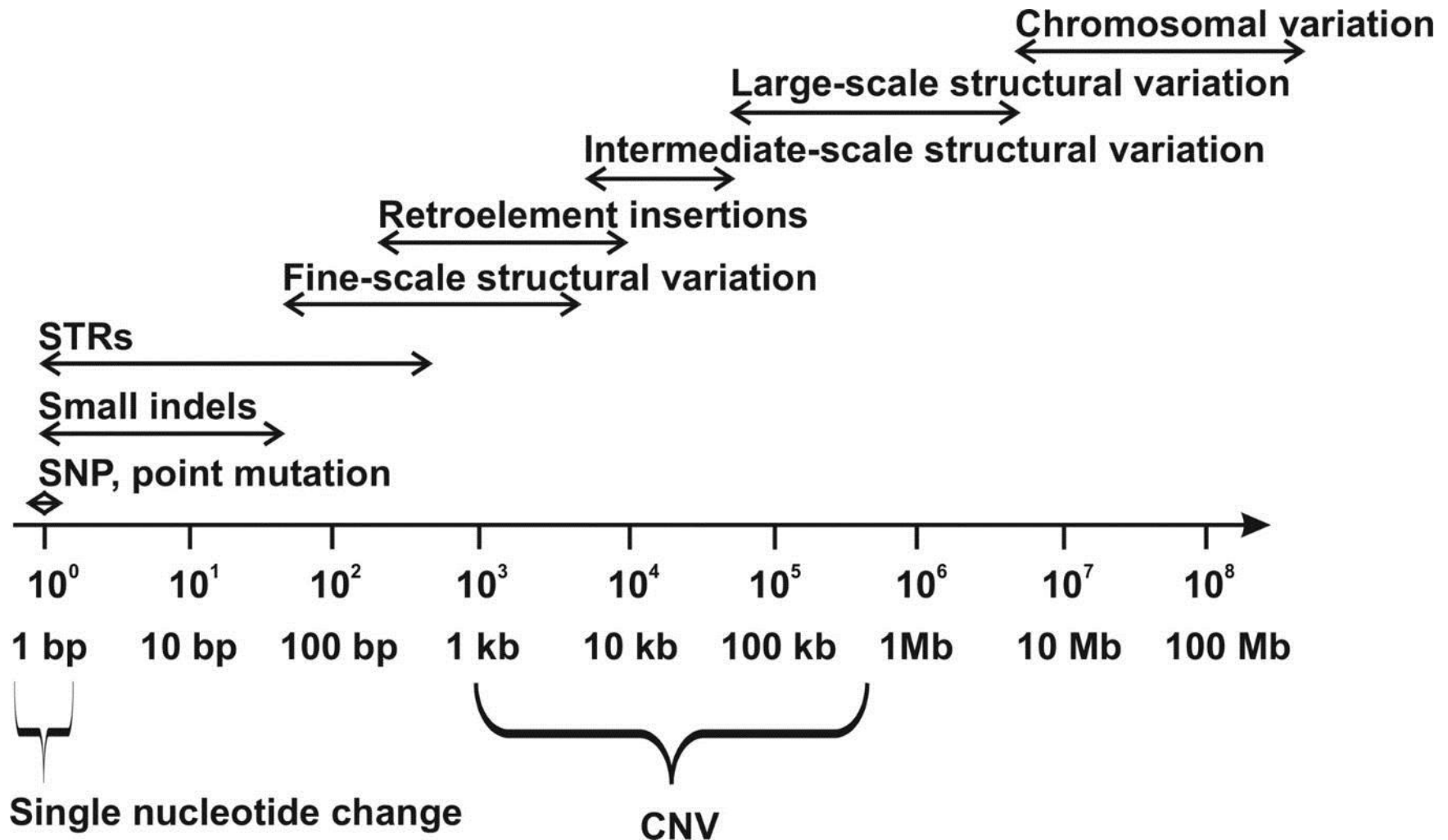


# Genomic technologies in disease studies

2019 Dragon Star Bioinformatics Course (Day 1)

# Human Genetic Variation



# Types of genetic variation

Single Nucleotide Variants (**SNVs**).

Reference Genome

**AGGTCATCGA**

Individual A

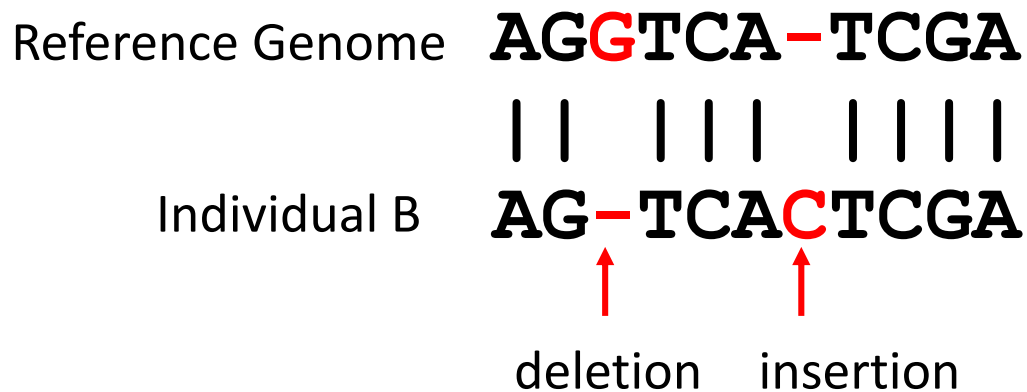
|||||  
**AGGTCCTCGA**



SNV (mismatch in alignment)

# Types of genetic variation

Insertion or deletion (< 50 bp), also known as **Indel**.



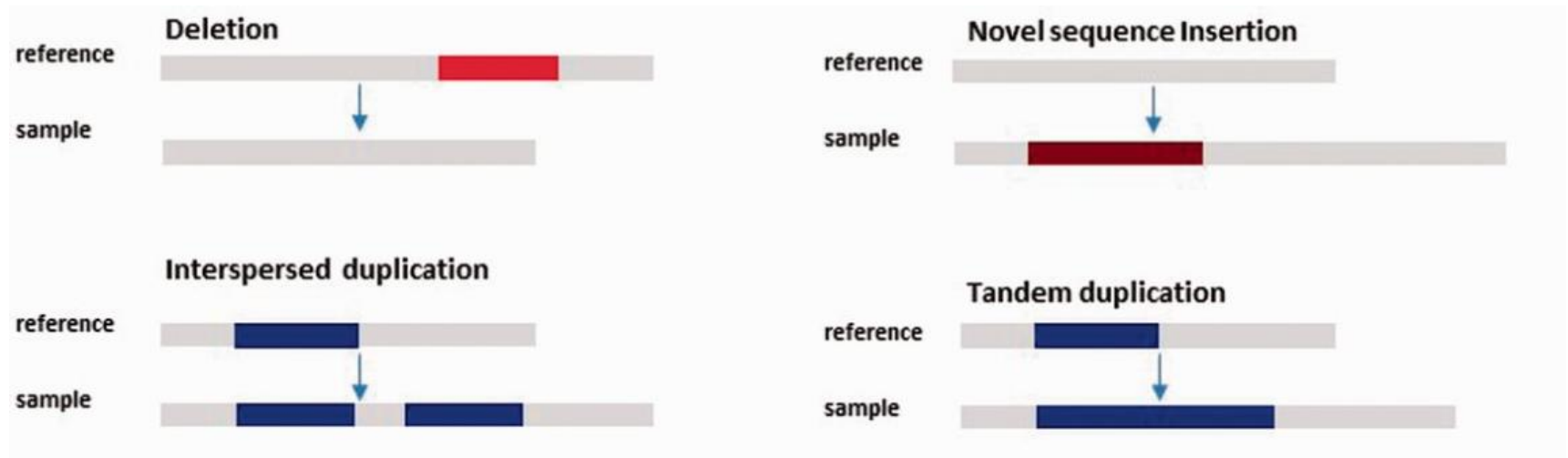
# Types of genetic variation

- **Structural Variants (SV):** generally defined as a region of DNA that shows a change in
  - Copy number (deletions, insertions and duplications)
  - Orientation (inversions) or
  - Chromosomal location (translocations) between individuals.



# Different types of SVs

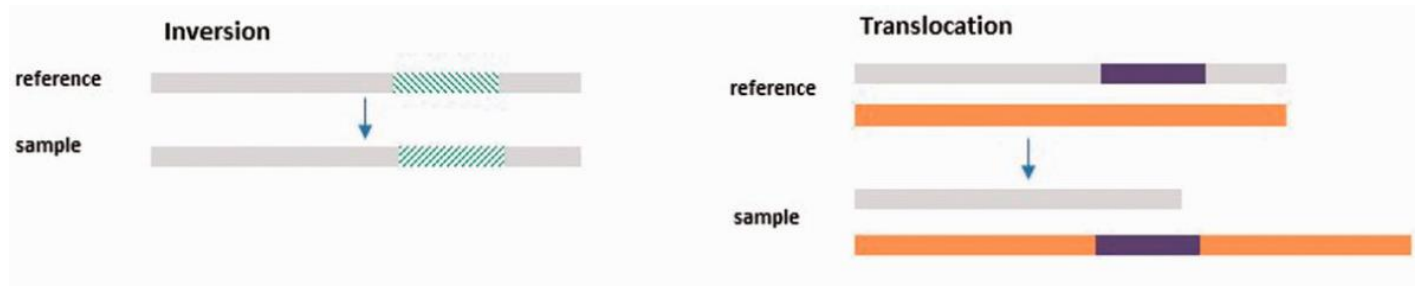
- SV can be *balanced* or *unbalanced*.
  - Unbalanced events: deletions/insertions/duplications
  - Chromosomal aneuploidies (such as trisomy 21) are extreme cases of unbalanced SV.



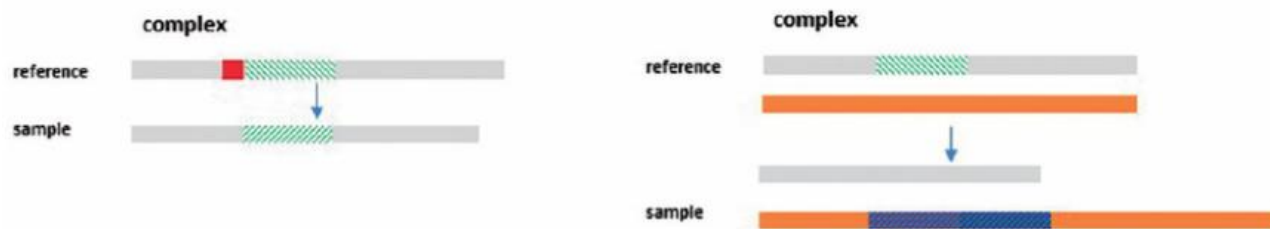
Deletions and duplications are two subtypes of CNVs (**Copy Number Variants**).

# Different types of SVs

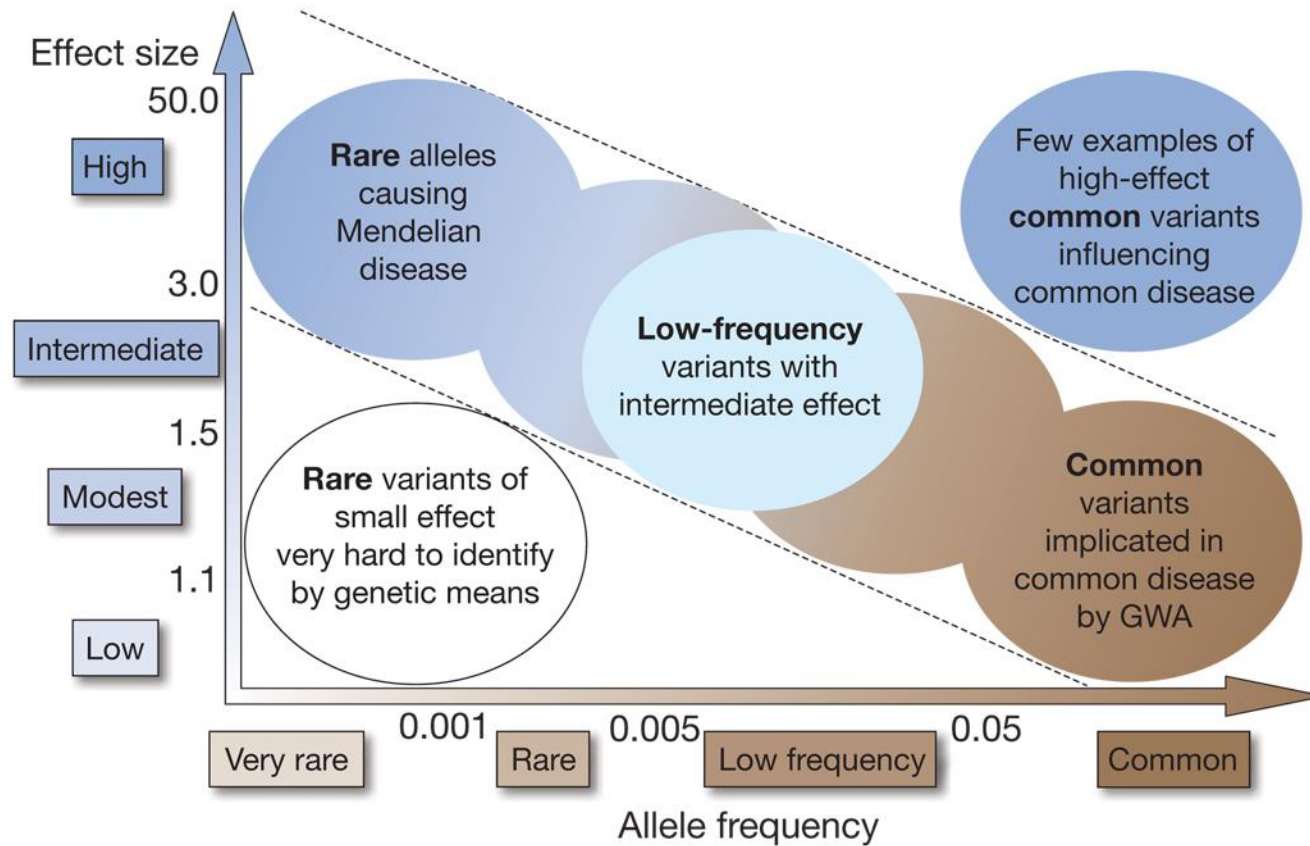
- Balanced events do not involve gain or loss of genetic materials
  - Inversions and translocations



- Complex SVs (several types together)

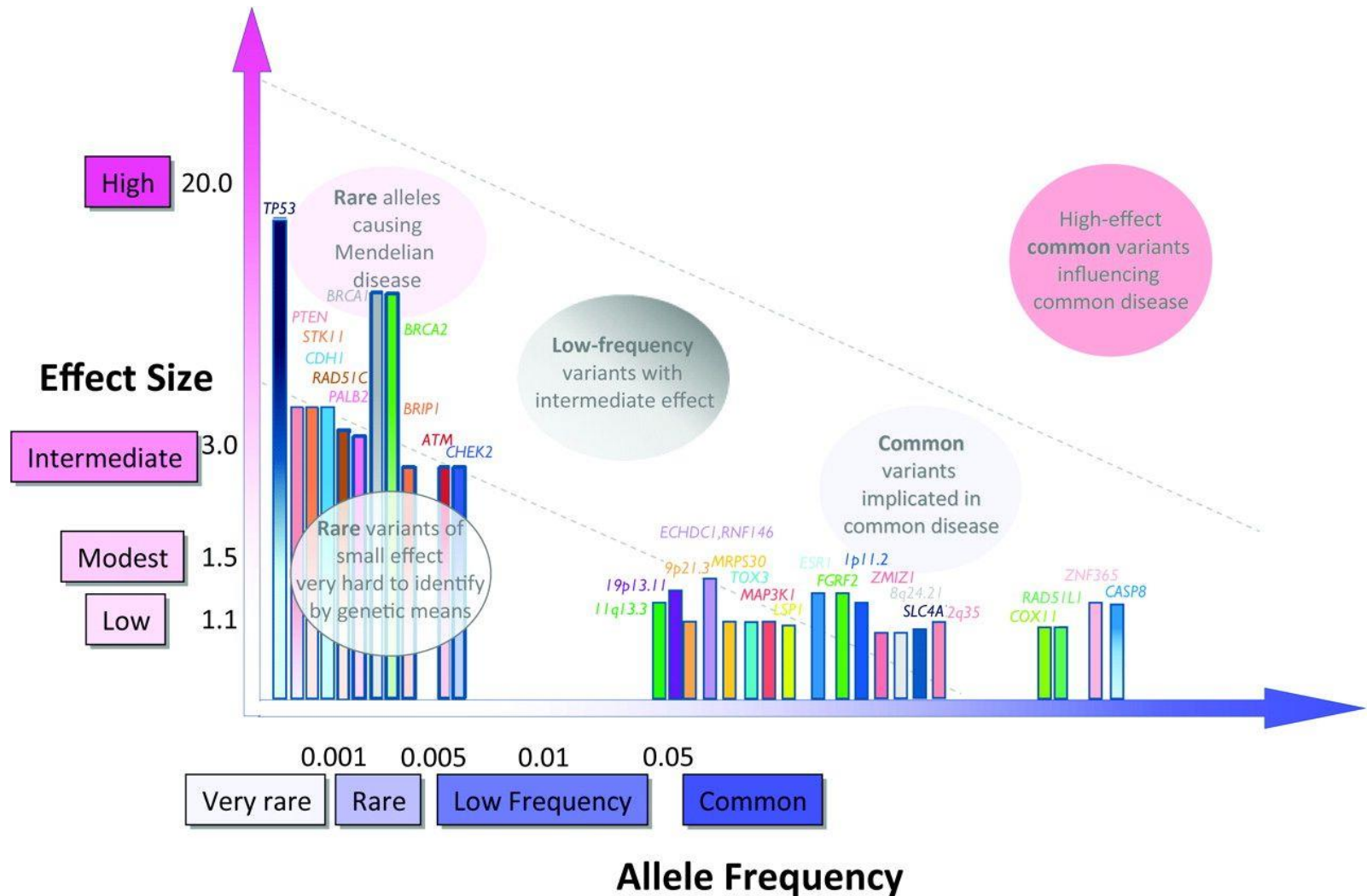


# Allele frequency and effect size

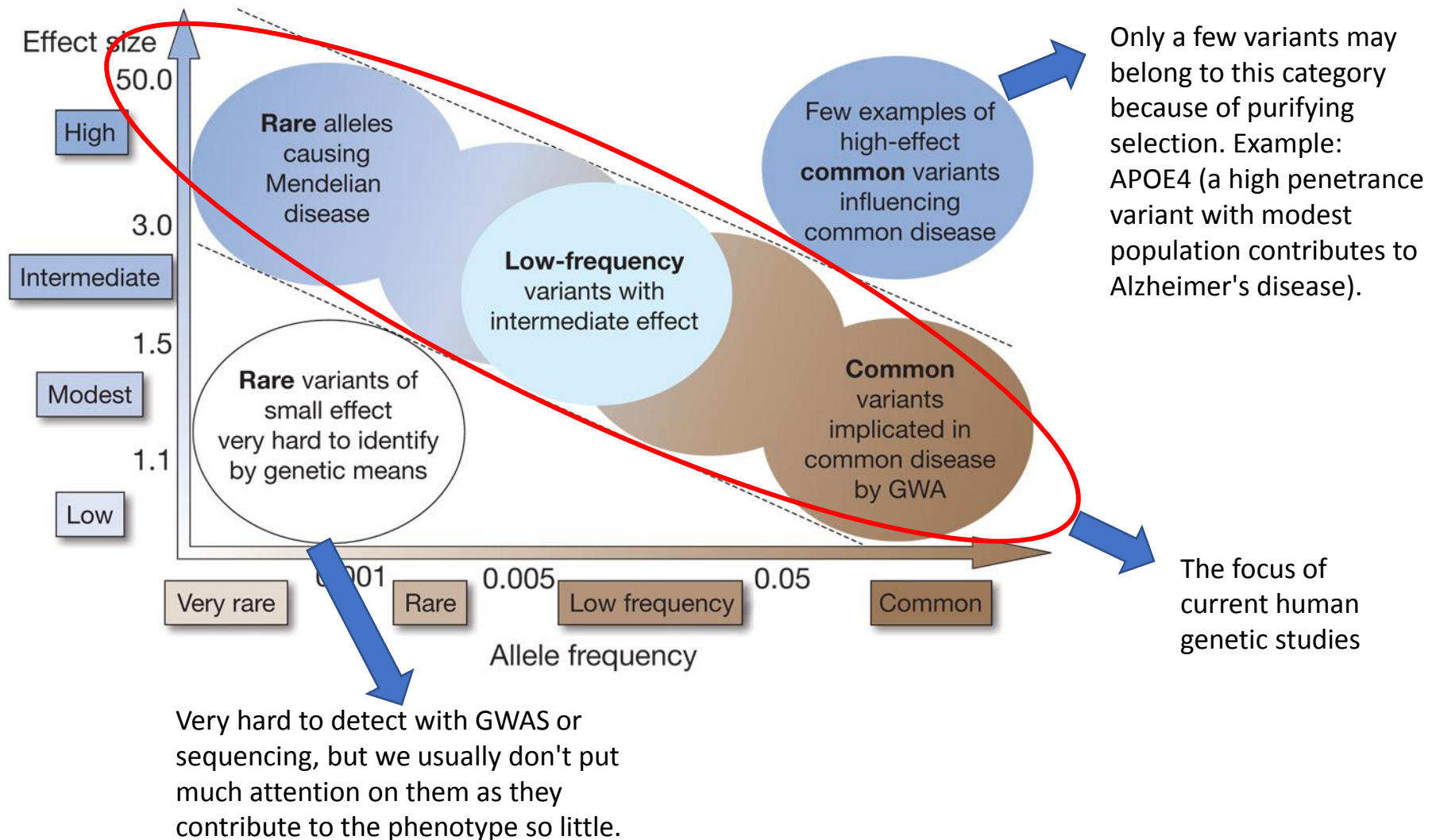




# Breast cancer as an example



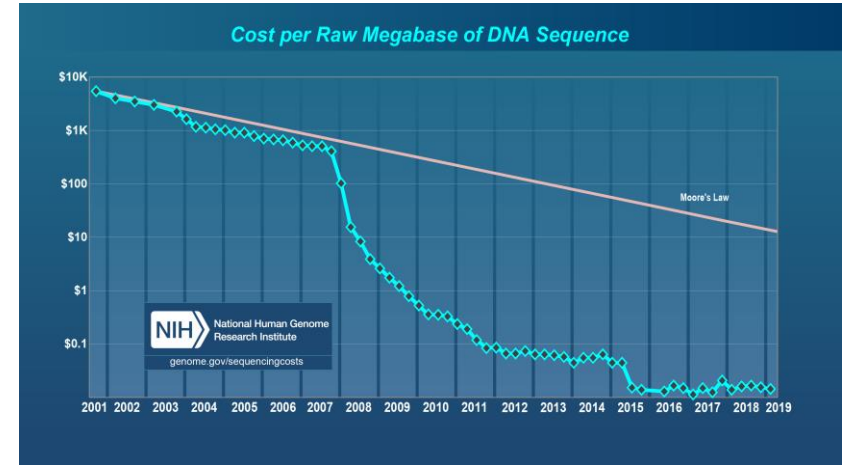
# Variant frequency and effect size



# History of DNA Sequencing

## Technical milestones

- 1953: Sequencing of insulin protein<sup>2</sup>
- 1965: Sequencing of alanine tRNA<sup>4</sup>
- 1968: Sequencing of cohesive ends of phage lambda DNA<sup>6</sup>
- 1977: Maxam–Gilbert sequencing<sup>9</sup>
- 1977: Sanger sequencing<sup>8</sup>
- 1981: Messing's M13 phage vector<sup>12</sup>
- 1986–1987: Fluorescent detection in electrophoretic sequencing<sup>14,15,17</sup>
- 1987: Sequenase<sup>18</sup>
- 1988: Early example of sequencing by stepwise dNTP incorporation<sup>139</sup>
- 1990: Paired-end sequencing<sup>23</sup>
- 1992: Bodipy dyes<sup>140</sup>
- 1993: *In vitro* RNA colonies<sup>37</sup>
- 1996: Pyrosequencing<sup>44</sup>
- 1999: *In vitro* DNA colonies in gels<sup>38</sup>
- 2000: Massively parallel signature sequencing by ligation<sup>47</sup>
- 2003: Emulsion PCR to generate *in vitro* DNA colonies on beads<sup>42</sup>
- 2003: Single-molecule massively parallel sequencing-by-synthesis<sup>33,34</sup>
- 2003: Zero-mode waveguides for single-molecule analysis<sup>57</sup>
- 2003: Sequencing by synthesis of *in vitro* DNA colonies in gels<sup>49</sup>
- 2005: Four-colour reversible terminators<sup>51–53</sup>
- 2005: Sequencing by ligation of *in vitro* DNA colonies on beads<sup>41</sup>
- 2007: Large-scale targeted sequence capture<sup>93–96</sup>
- 2010: Direct detection of DNA methylation during single-molecule sequencing<sup>65</sup>
- 2010: Single-base resolution electron tunnelling through a solid-state detector<sup>141</sup>
- 2011: Semiconductor sequencing by proton detection<sup>142</sup>
- 2012: Reduction to practice of nanopore sequencing<sup>143,144</sup>
- 2012: Single-stranded library preparation method for ancient DNA<sup>145</sup>

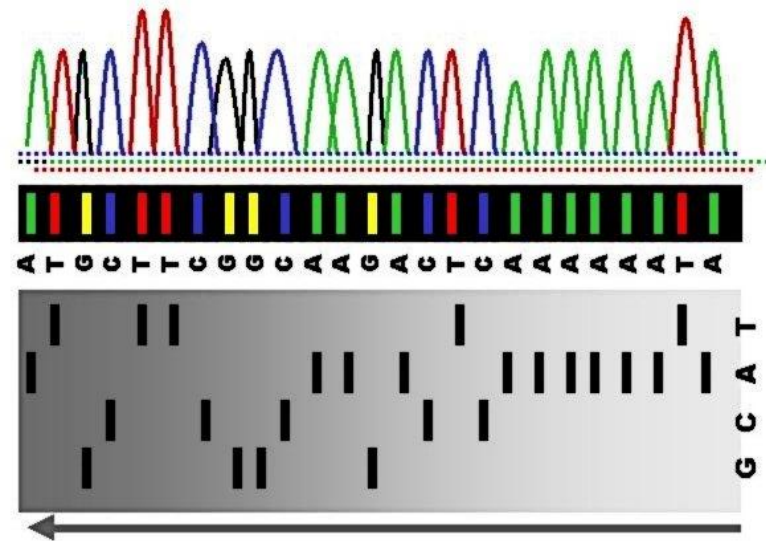
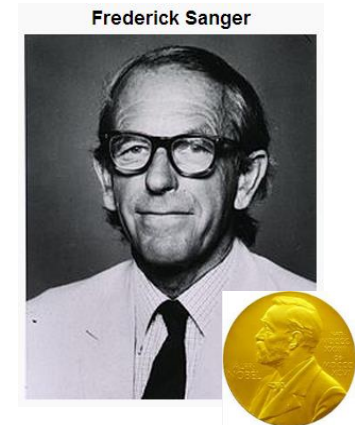


## Illumina NovaSeq

<b>S4</b>	<b>10<sub>B</sub></b>	<b>2000<sub>Gb</sub></b>	<b>3000<sub>Gb</sub></b>
Flow cell type	Single reads*	2 x 100 output	2 x 150 output

# Sanger Sequencing

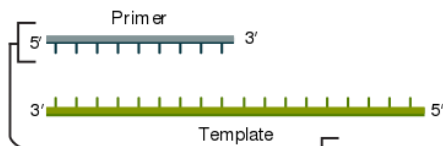
- Developed by Frederick Sanger and colleagues in 1977
- Up to 1,000 bases
- First human genome draft was based on Sanger sequencing
- Remains in wide use today, for smaller-scale projects and for validation of next-generation sequencing results



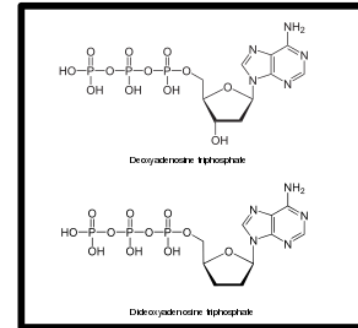
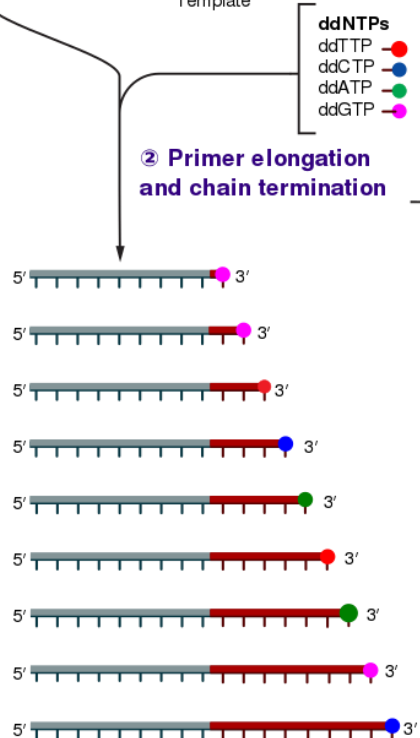
# Sanger Sequencing

## ① Reaction mixture

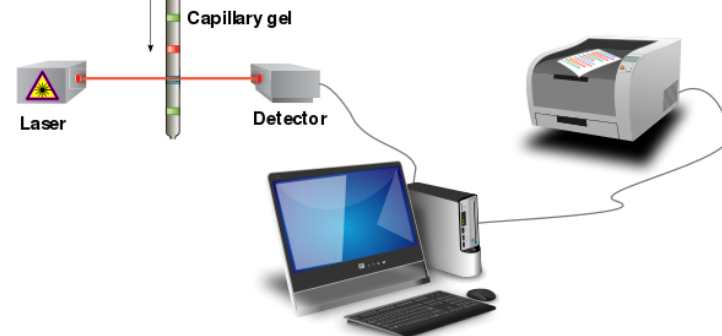
- ▶ Primer and DNA template ▶ DNA polymerase
- ▶ ddNTPs with flouochromes ▶ dNTPs (dATP, dCTP, dGTP, and dTTP)



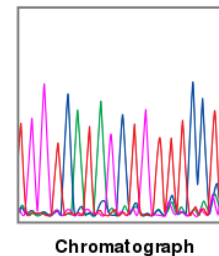
## ② Primer elongation and chain termination



## ③ Capillary gel electrophoresis separation of DNA fragments



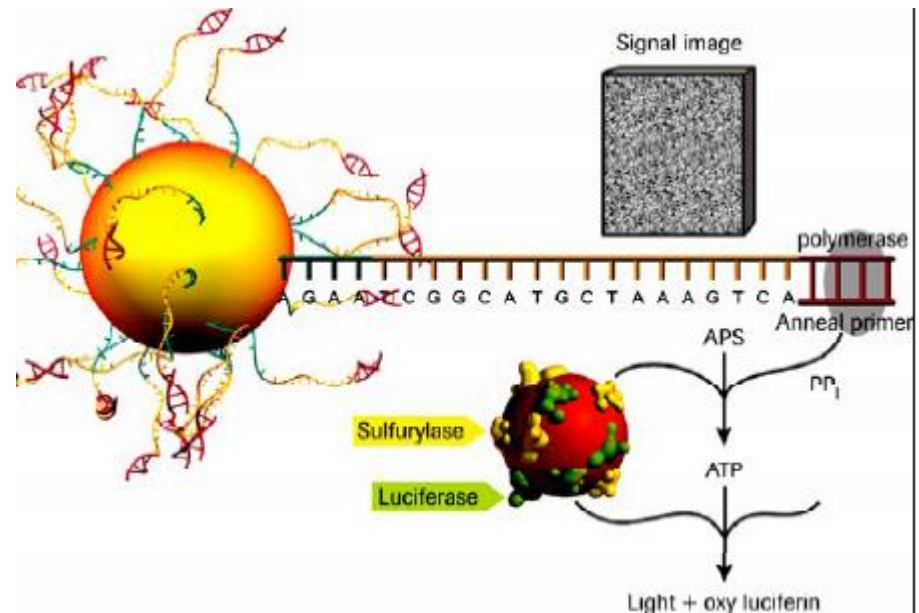
## ④ Laser detection of flouochromes and computational sequence analysis



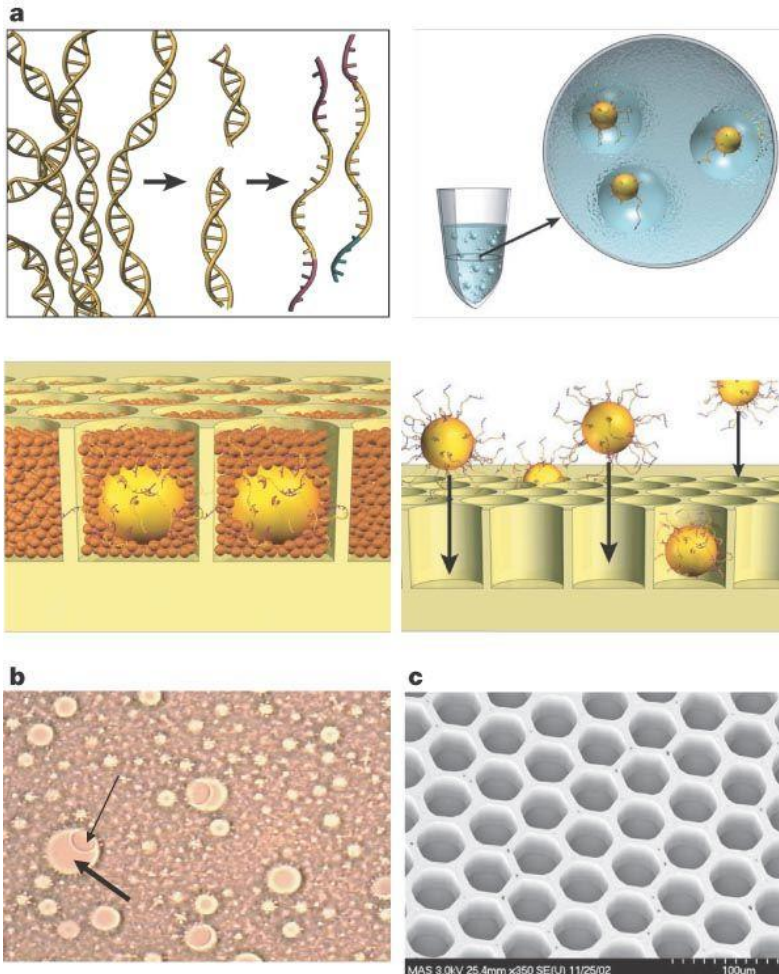


# Next-generation sequencing

- Pyrosequencing: incorporation of nucleotide that results in the release of pyrophosphates which fuels the production of light by firefly enzyme luciferase.
- licensed to 454 Life Sciences, where it evolved into the first major successful commercial 'next-generation sequencing' (NGS) technology.



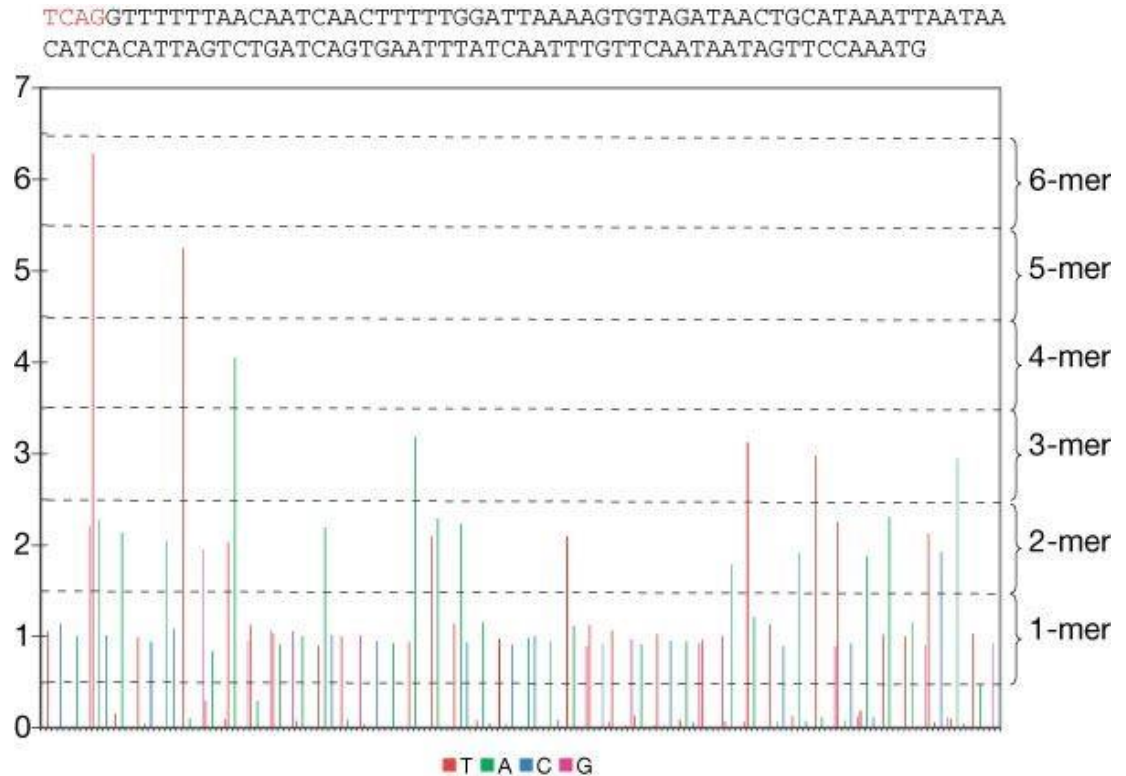
# 454 Sequencing



- Genomic DNA is isolated, fragmented, ligated to adapters and separated into single strands.
- Fragments are bound to beads (one fragment per bead), the beads are captured in the droplets; emulsion PCR occurs within each droplet
- Beads carrying single-stranded DNA clones are deposited into wells of a fibre-optic slide
- After the flow of each nucleotide, a wash containing apyrase is used to ensure that nucleotides do not remain in any well before the next nucleotide being introduced.

# 454 sequencing: base calling

Nucleotide incorporation is detected by the associated release of inorganic pyrophosphate and the generation of photons





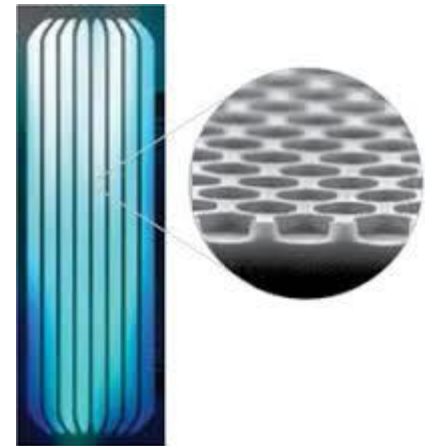
# Illumina short-read sequencing

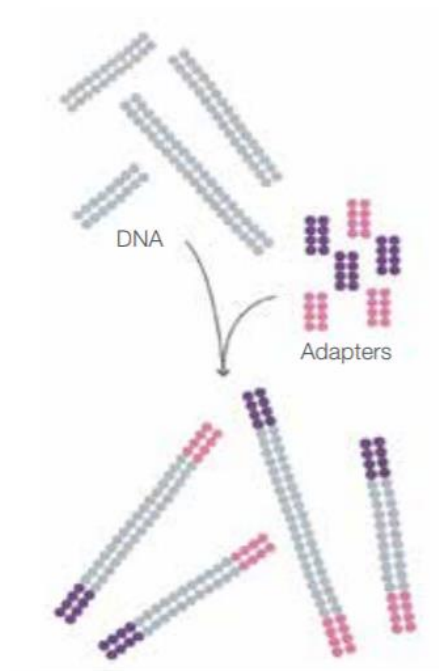
- Illumina sequencing technology, sequencing by synthesis (SBS), is a widely adopted next-generation sequencing (NGS) technology worldwide, responsible for generating more than 90% of the world's sequencing data



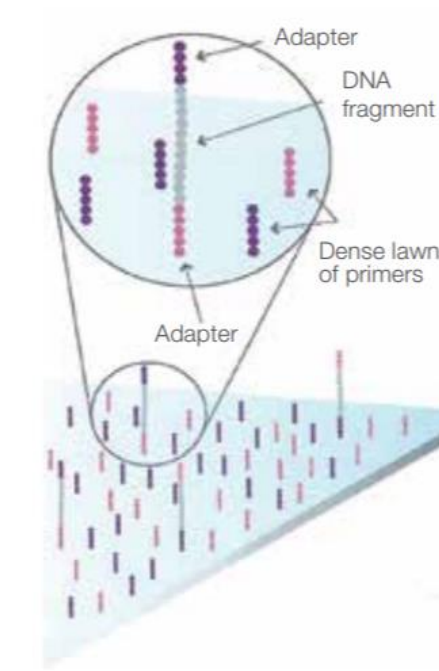
# How Illumina sequencing works

- Cluster generation:
  - Solid-phase amplification creates up to 1,000 identical copies of each single template molecule in close proximity (diameter of 1µm or less), using unlabeled nucleotides.
- Sequencing by synthesis (SBS):
  - Four fluorescently labeled nucleotides to sequence the tens of millions of clusters on the flow cell surface in parallel
  - During each sequencing cycle, a single labeled deoxynucleoside triphosphate (dNTP) is added to the nucleic acid chain. The nucleotide label serves as a terminator for polymerization, so after each dNTP incorporation, the fluorescent dye is imaged to identify the base and then enzymatically cleaved to allow incorporation of the next nucleotide.

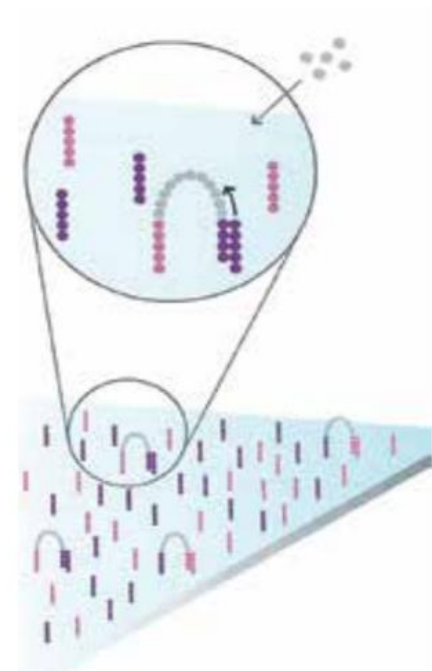




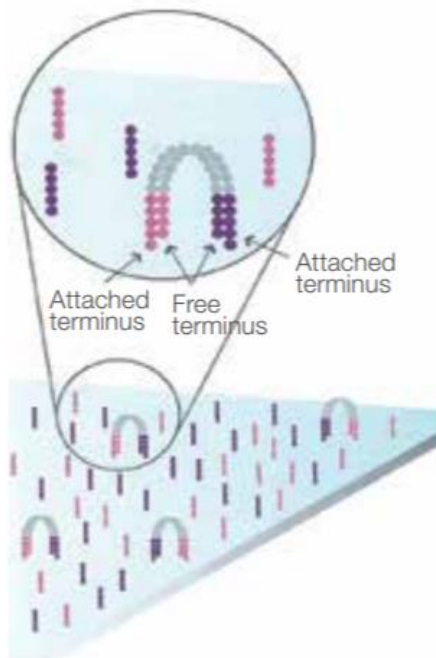
Prepare DNA Sample



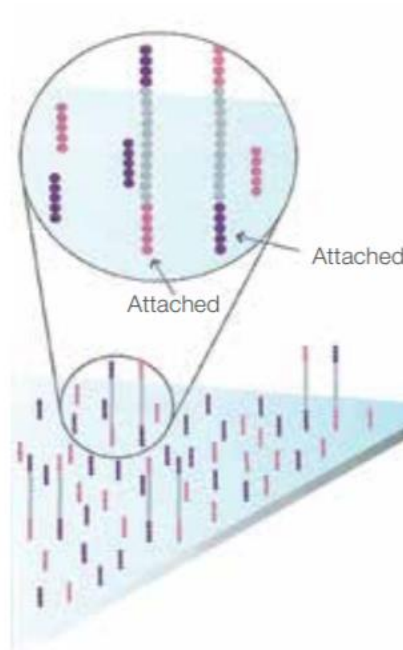
Attach DNA to Surface



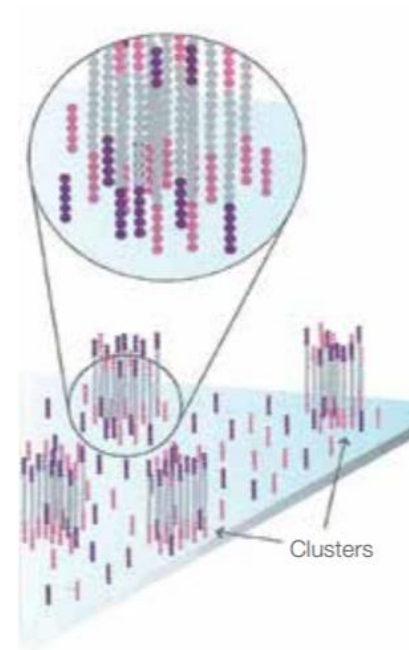
Bridge Amplification



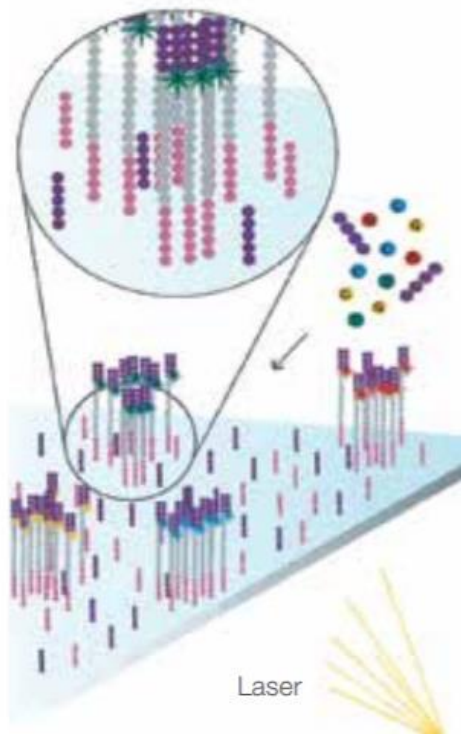
Fragments Become  
Double Stranded



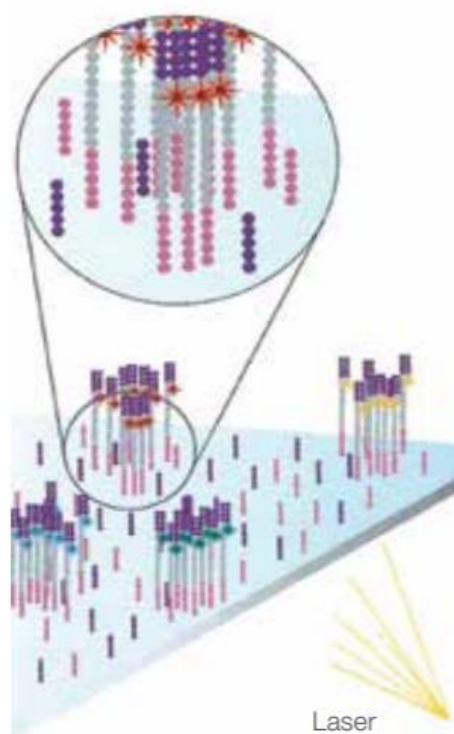
Denature the Double-  
Stranded Molecules



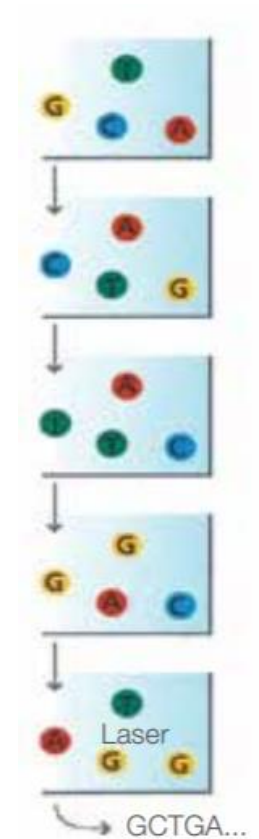
Complete Amplification



Determine First Base



Determine Second Base



Sequencing Over Multiple Cycles

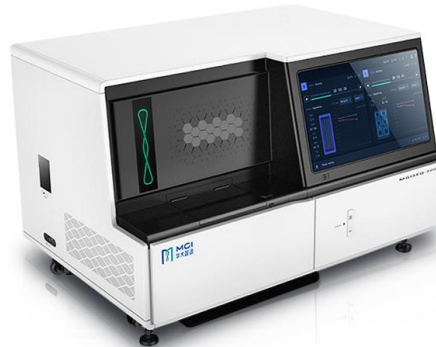
# BGISeq and MGISEq



Complete  
Genomics



BGISEQ-500



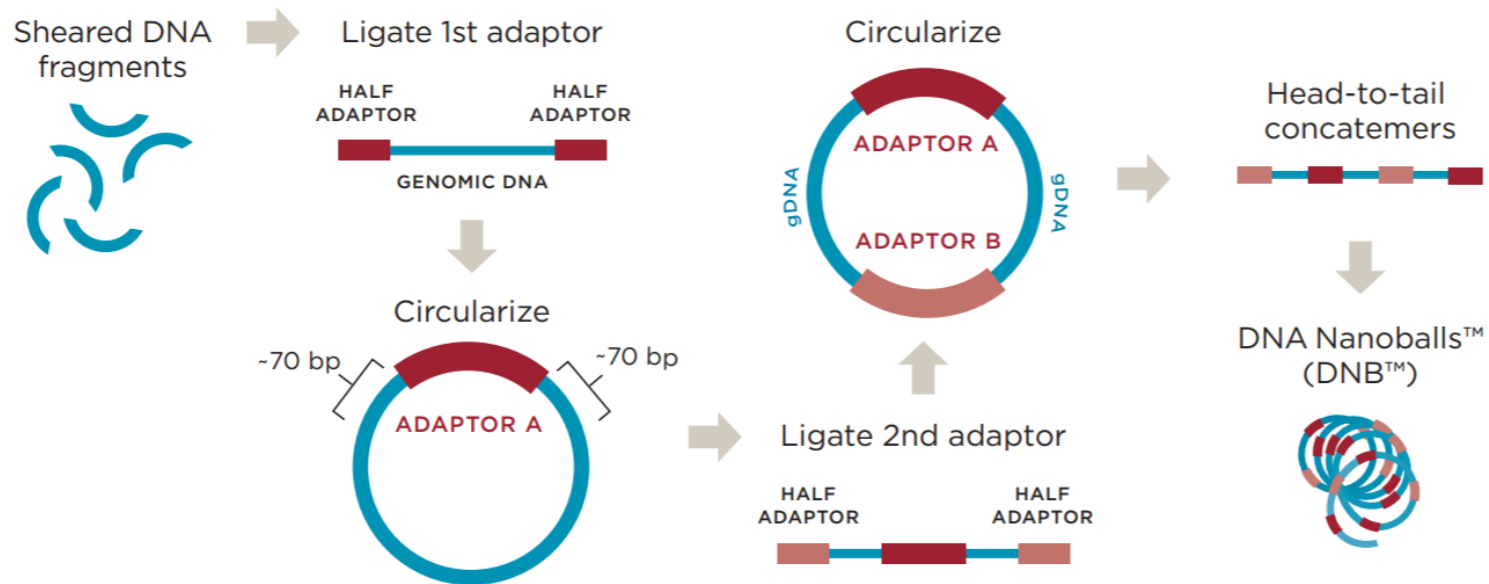
MGISEQ-2000



MGISEQ-T7

# Production of DNA Nanoballs

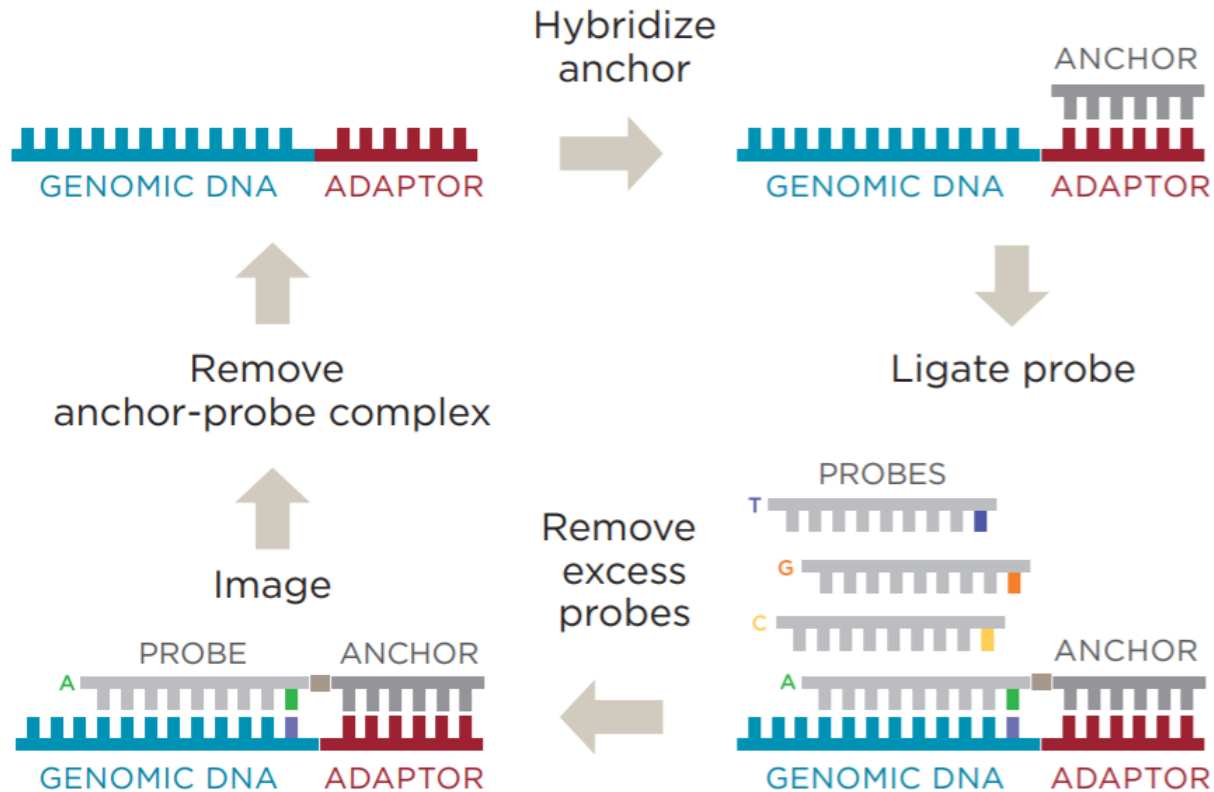
- The circular DNA molecules in the flow cell library are clonally amplified and modified to produce DNA Nanoballs (DNBs), each containing more than 200 copies of the original template





# Ligation-based cPAL (Combinatorial Probe-Anchor Ligation) sequencing chemistry

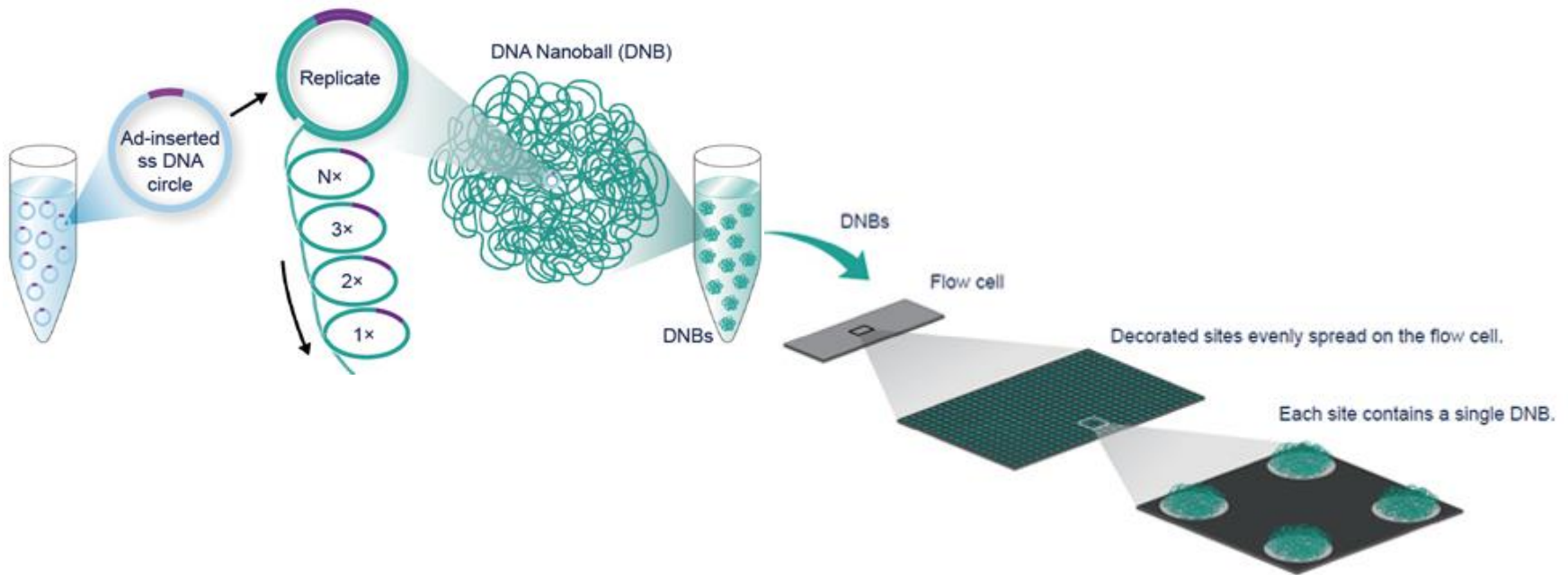
## CPAL SEQUENCING TECHNOLOGY





# From cPAL (hybridization) to cPAS (synthesis)

- Each cycle: addition of fluorescently labelled terminated dNTPs, cleavage of a terminator, and the detection of the produced fluorescent signal



# Revolution: single-molecule long-read sequencing

PacBio

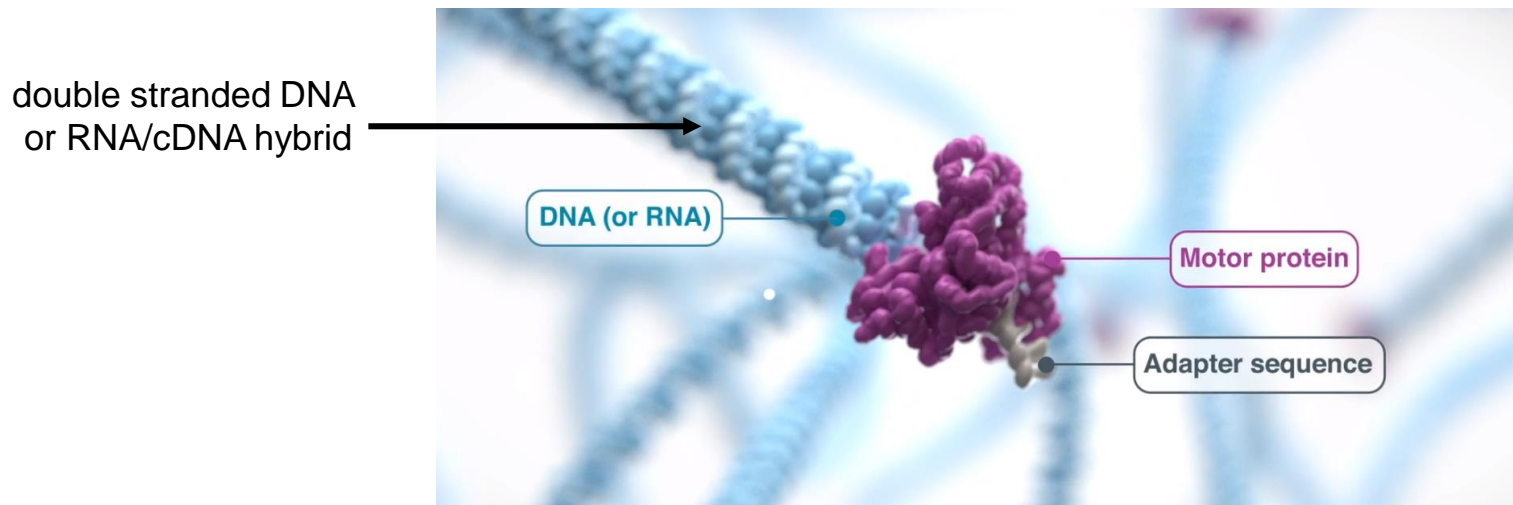


Oxford Nanopore



# Oxford Nanopore Sequencing

- Oxford Nanopore Sequencing is a real-time, direct DNA/RNA sequencing technology.
- The DNA/RNA is sequenced when it is going through a protein pore.

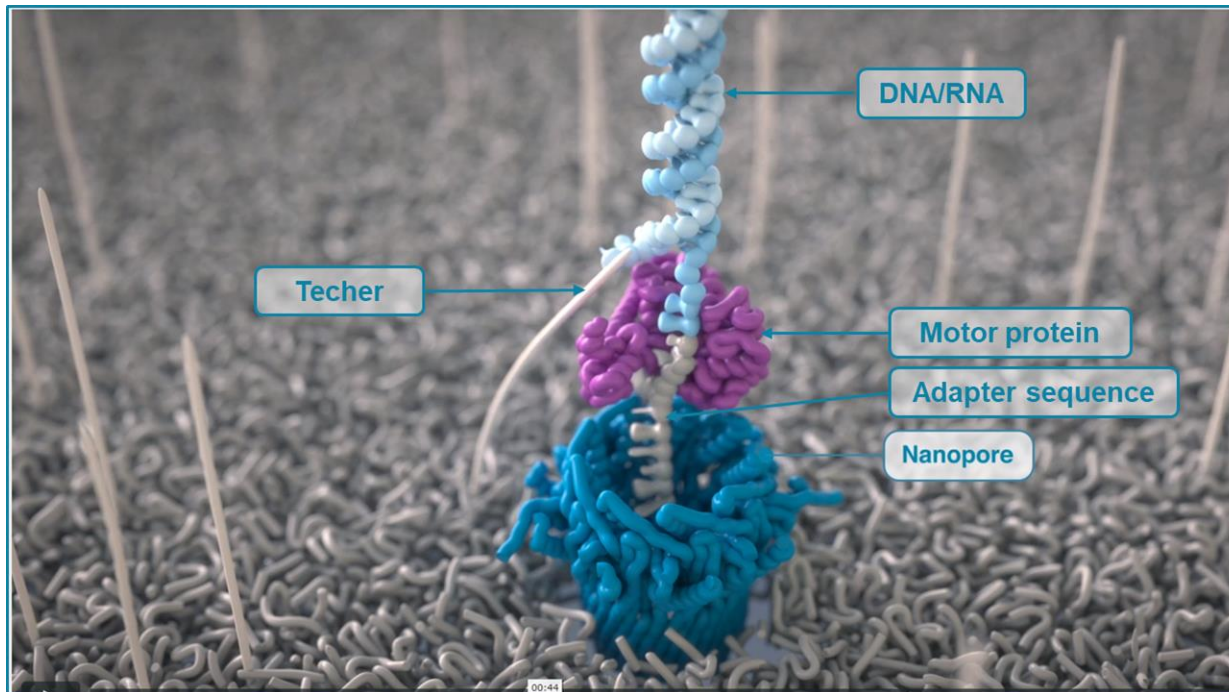


# Oxford Nanopore Sequencing



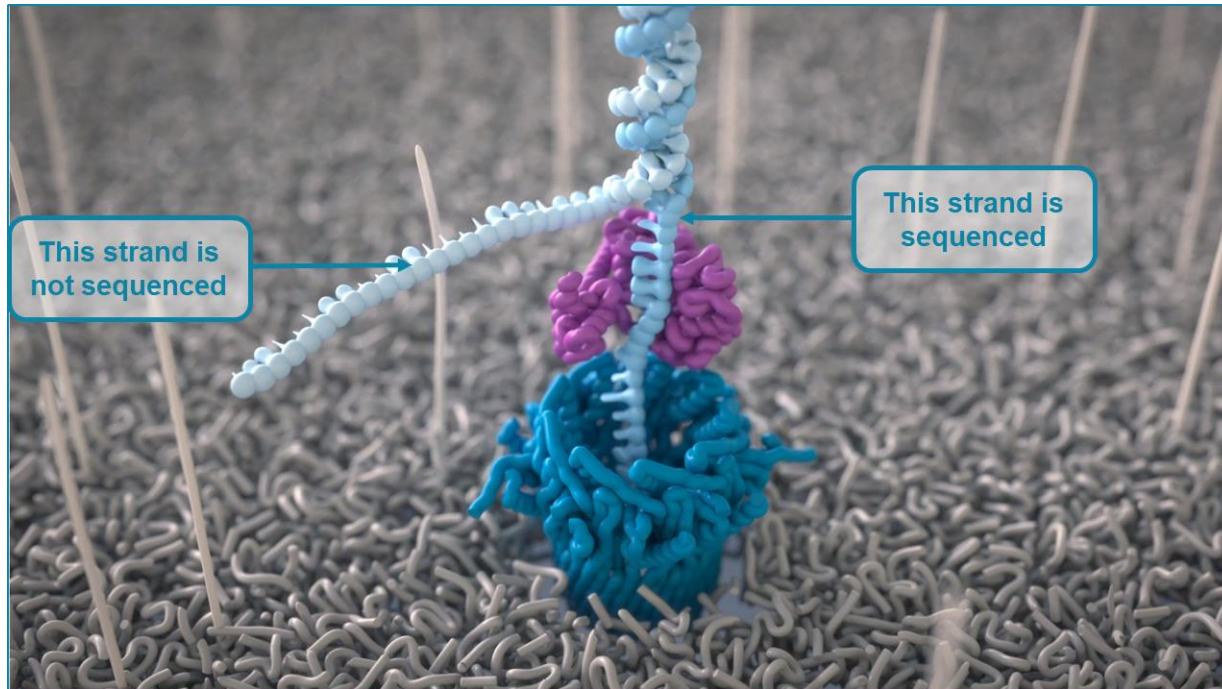
→ Each well contains a nanopore

# Oxford Nanopore Sequencing

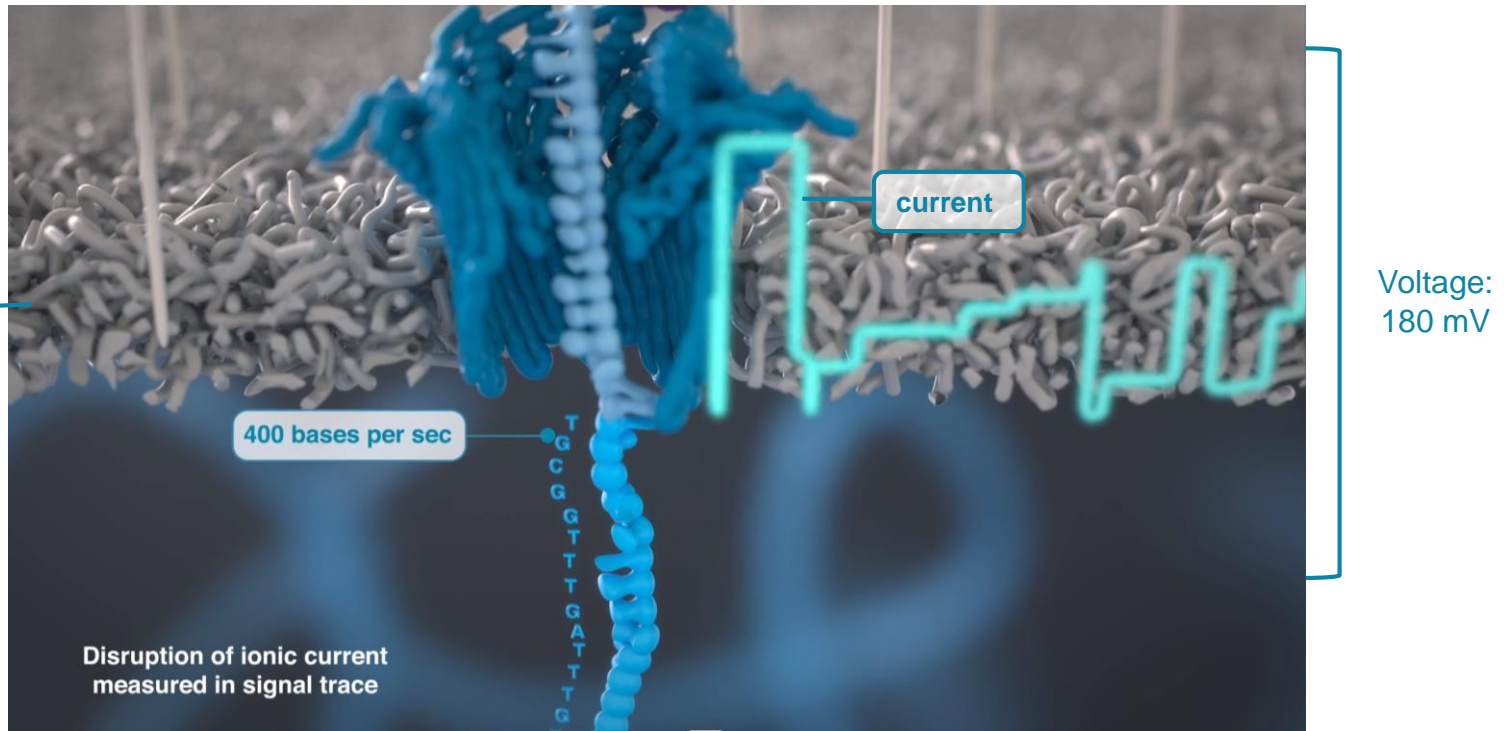




# Oxford Nanopore Sequencing



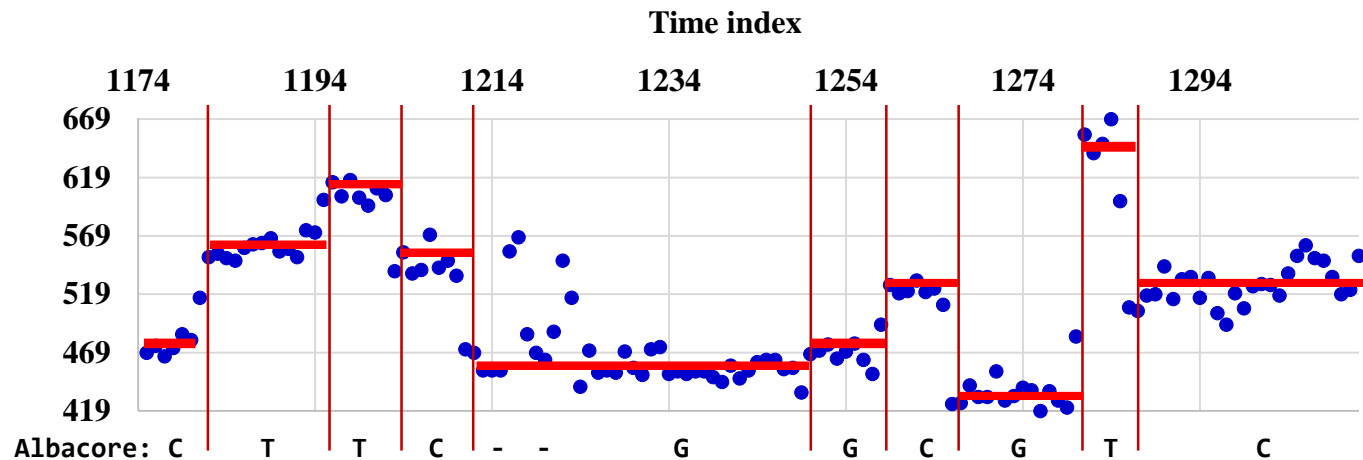
# Oxford Nanopore Sequencing



The nucleotides in the DNA/RNA block the ionic current and induce changes of current, which can be measured.

# How does the data look like?

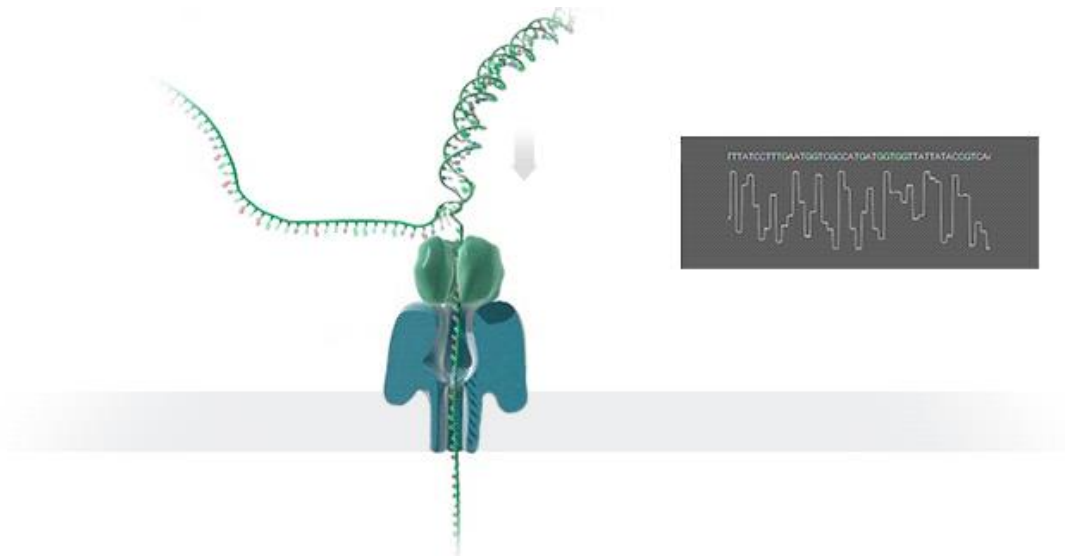
- Nanopore sequencing
  - The raw data is electric current (dot)
  - Event detection (red)
  - Base calling: A, C, G, T



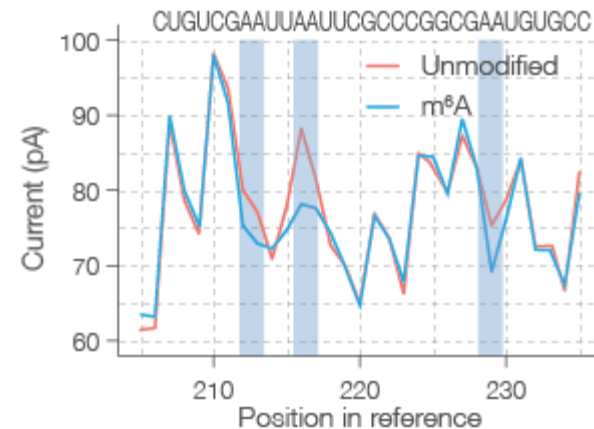
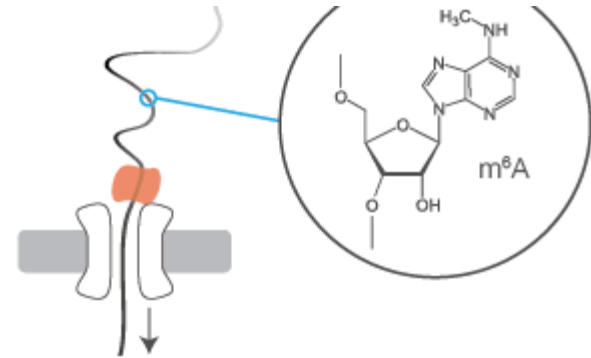
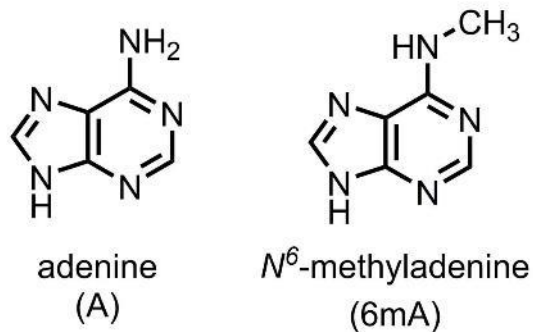
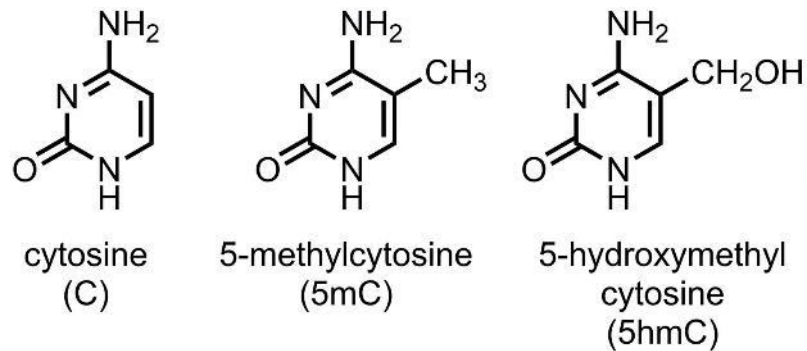


# From A/C/G/T to DNA modifications

- Change of current when a molecule pass through a tiny hole
- Different types of nucleotides and different modifications of nucleotides would generate different signals
- Currently, homopolymer repeats are an issue



# Detect direction of DNA methylations



Shi et al, Front. Genet, 2017

# PacBio Single-molecule real-time (SMRT) sequencing

PacBio RS II



Sequel



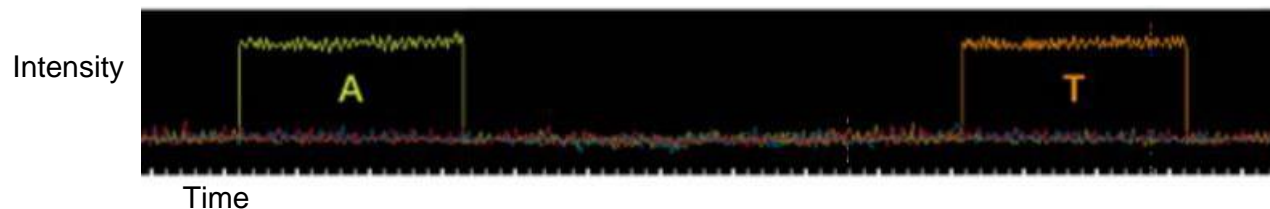
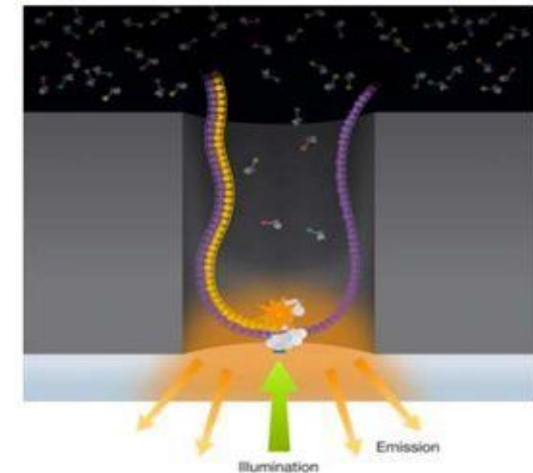
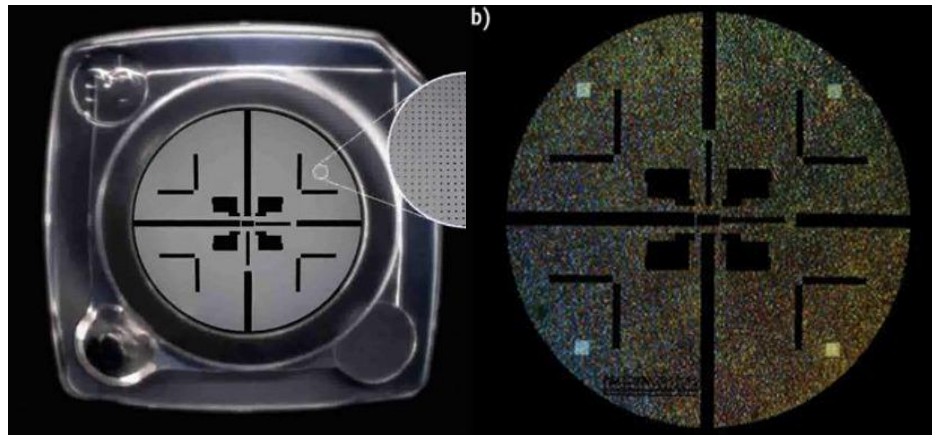
Sequel II



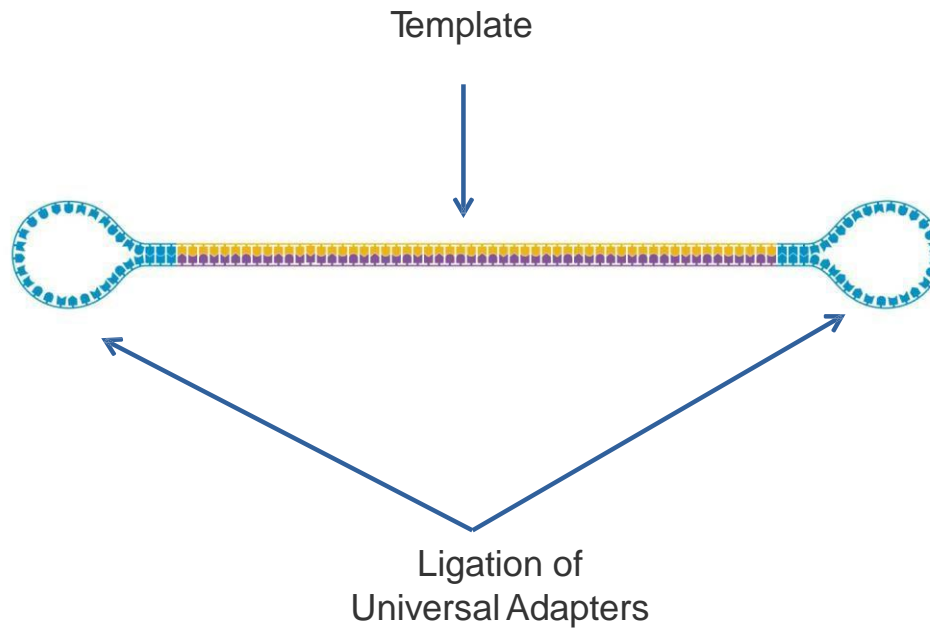
150K/1M/8M zero-mode waveguides (ZMWs)

# SMRT sequencing

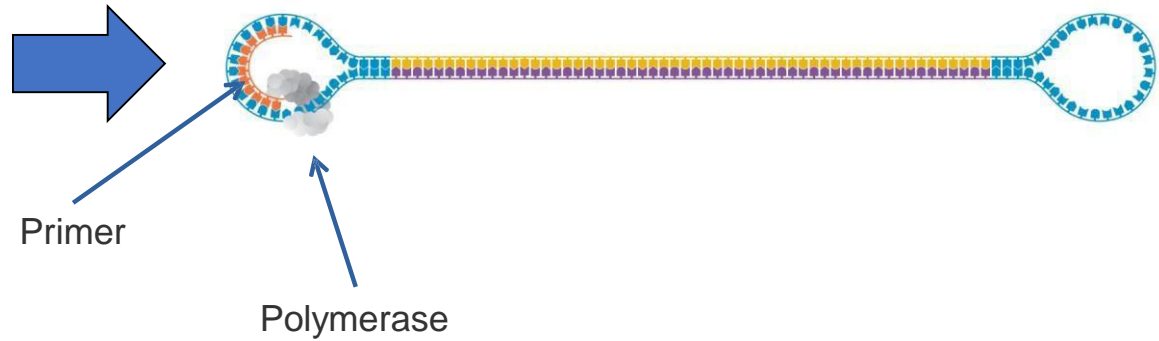
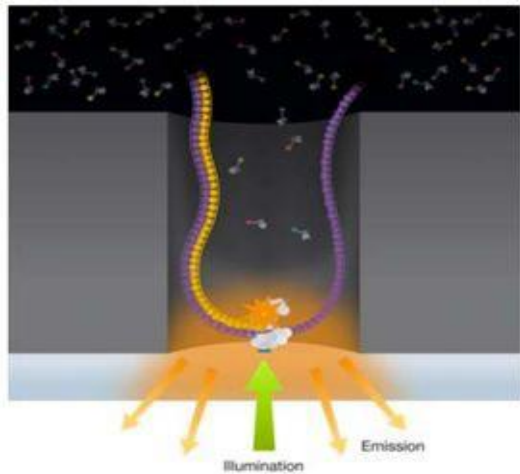
- Imaging of fluorescent phospholinked labeled nucleotides as they are incorporated by a polymerase anchored to a Zero-Mode Waveguide (ZMW).



# SMRTbell library construction



# SMRTbell sequencing



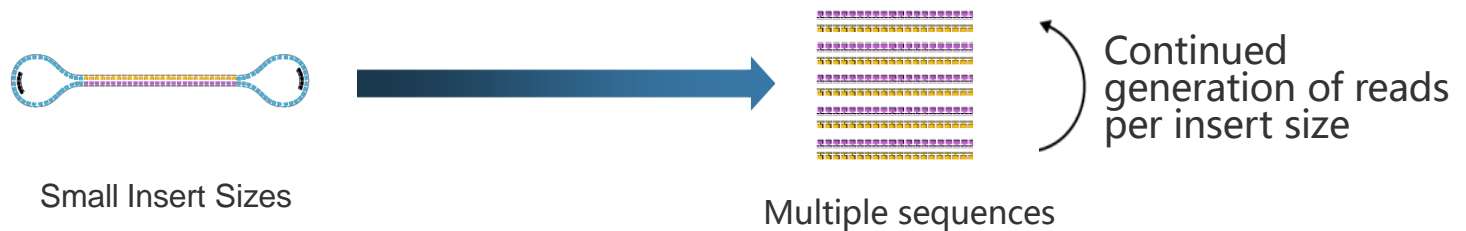
# Types of SMRT sequencing reads

## Continuous Long Reads (CLR)



Long inserts so that the polymerase can synthesize along a single strand

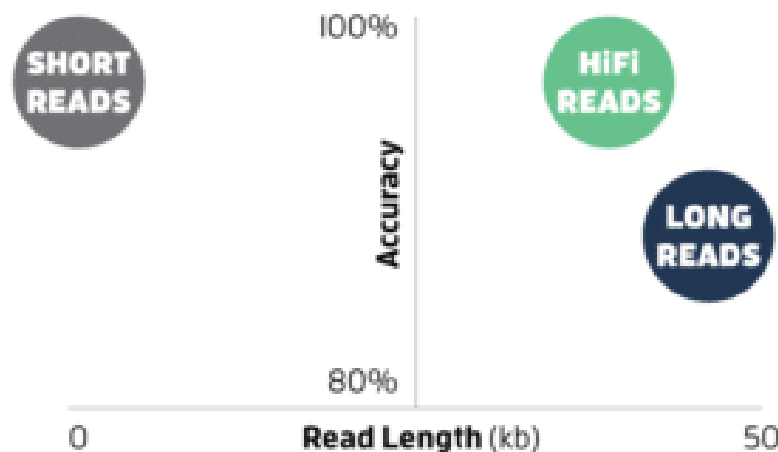
## Circular Consensus Sequencing (CCS)



Short inserts, so polymerase can continue around the entire SMRTbell multiple times and generate multiple sub-reads from the same single molecule

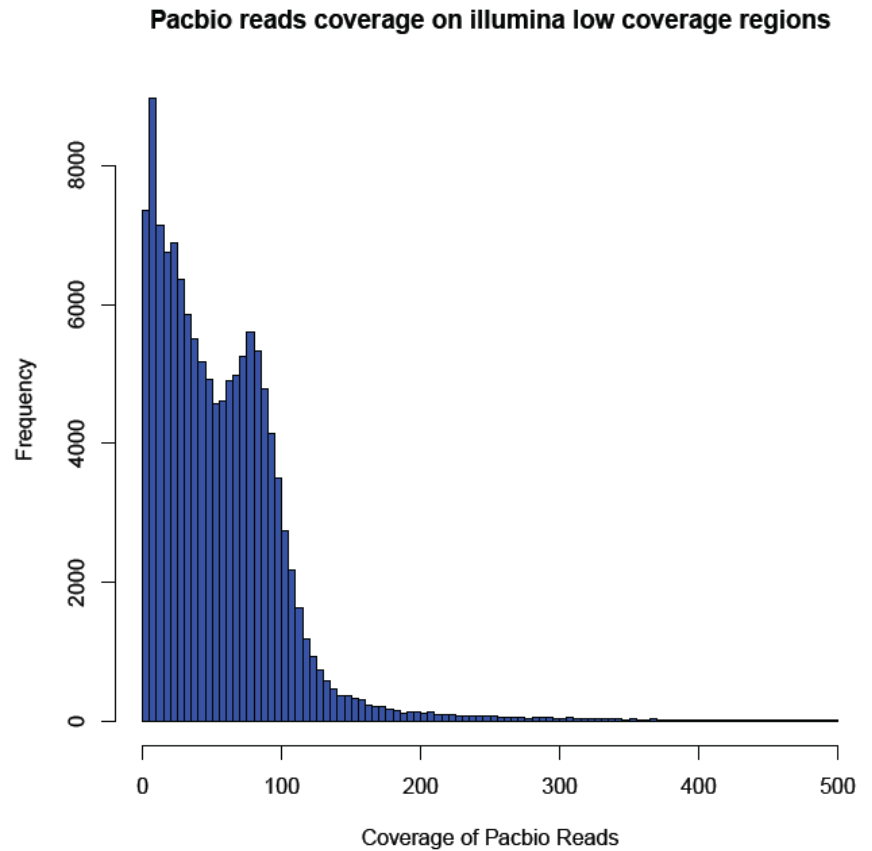
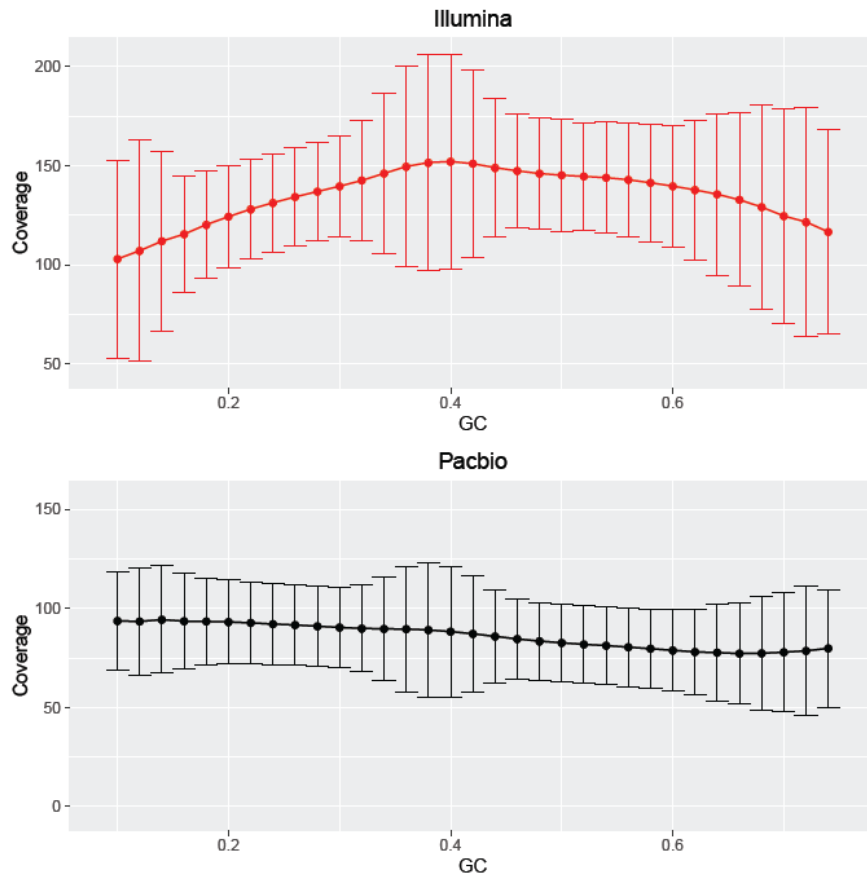
# Difference between CLS and CCS (HiFi)

- On Sequel 2, typically ~100 Gb per 8M SMRT cell using long insert/CLR libraries; or ~12-15 Gb >Q30 HiFi CCS per 8M SMRT cell using HiFi libraries.



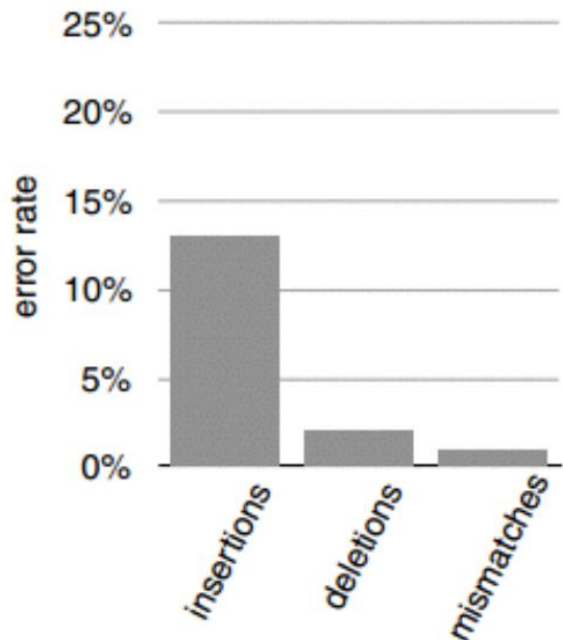


# Impacts of GC on read depth

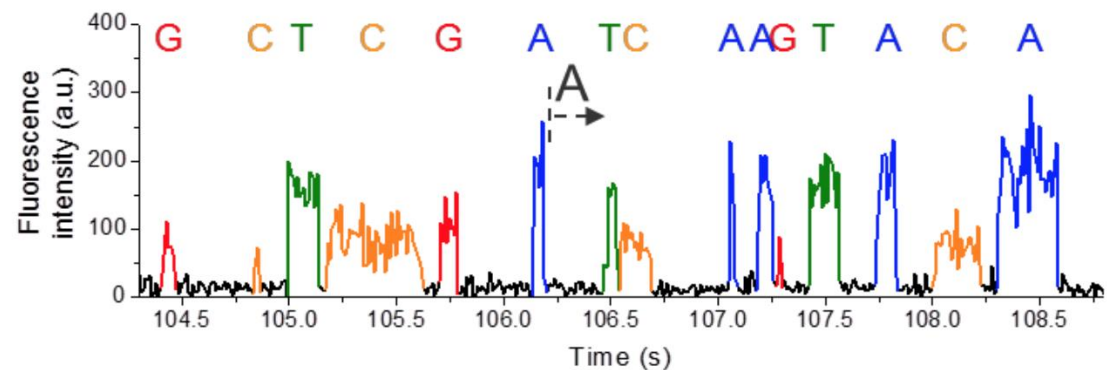
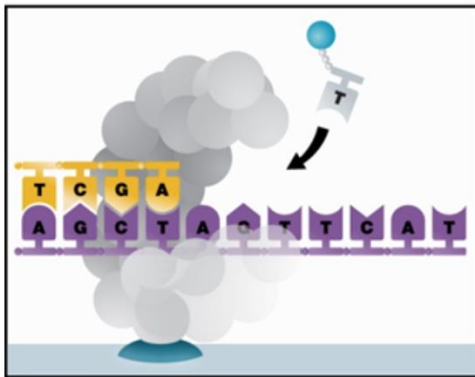
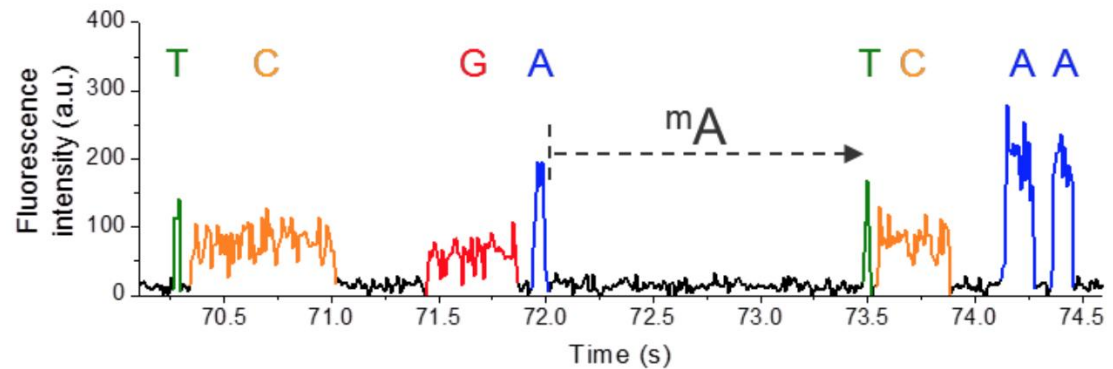
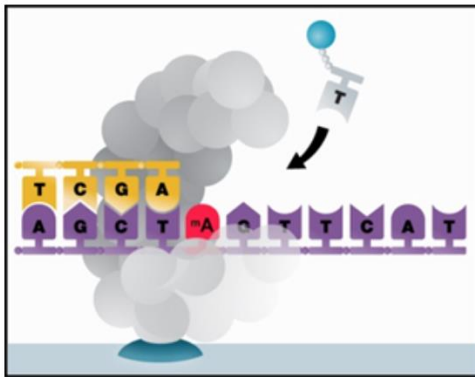


# Error profile

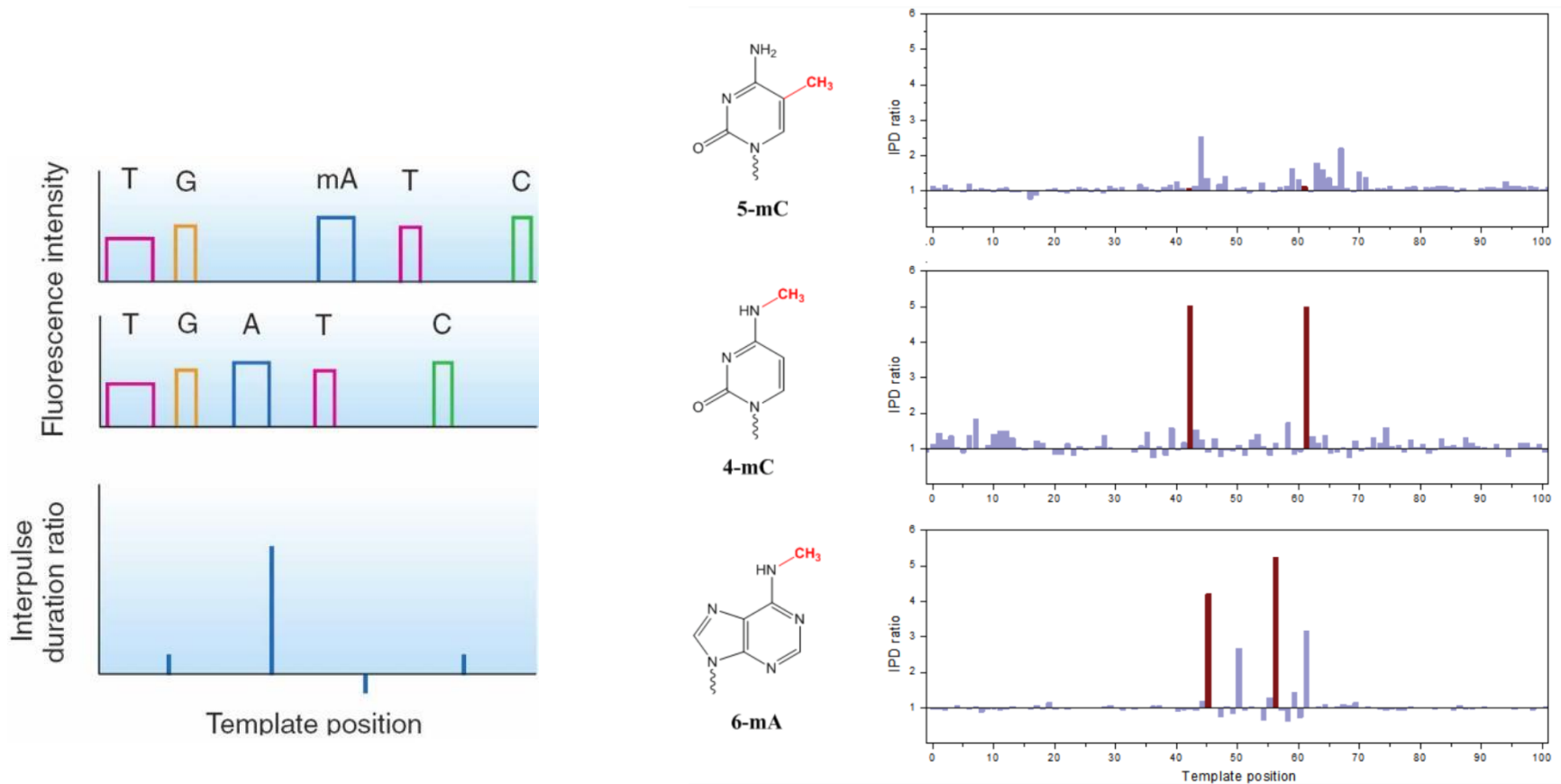
- Insertions tend to be more than deletions and substitutions



# Detection of DNA Base Modifications Using IPD (Interpulse duration ratio)



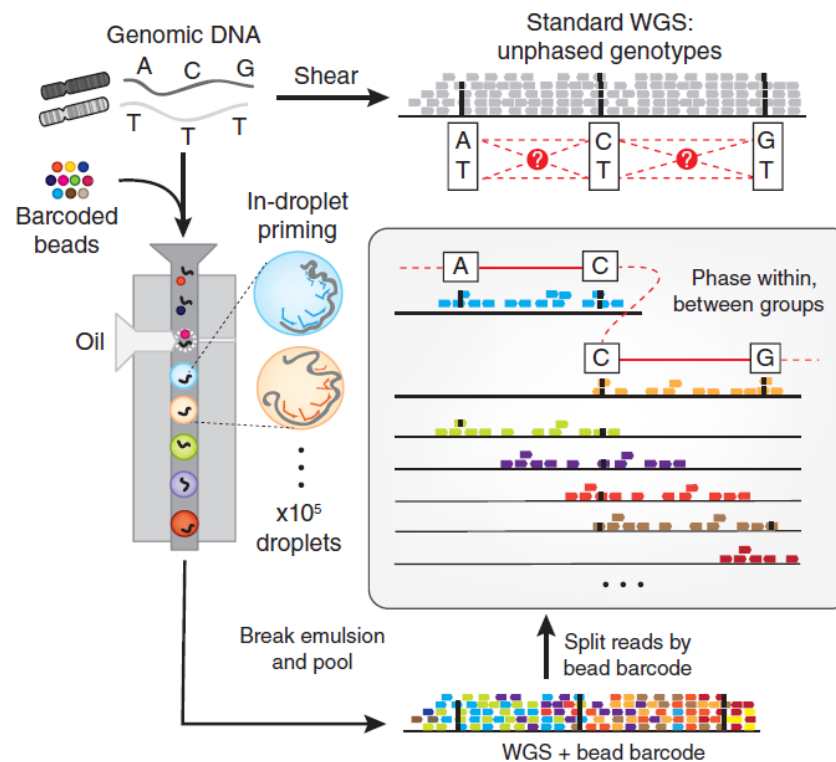
# Different modifications have different IPD patterns



- Coverage needs vary based on the strength of the kinetic signal.
- Kinetic signal strength varies by modification type.

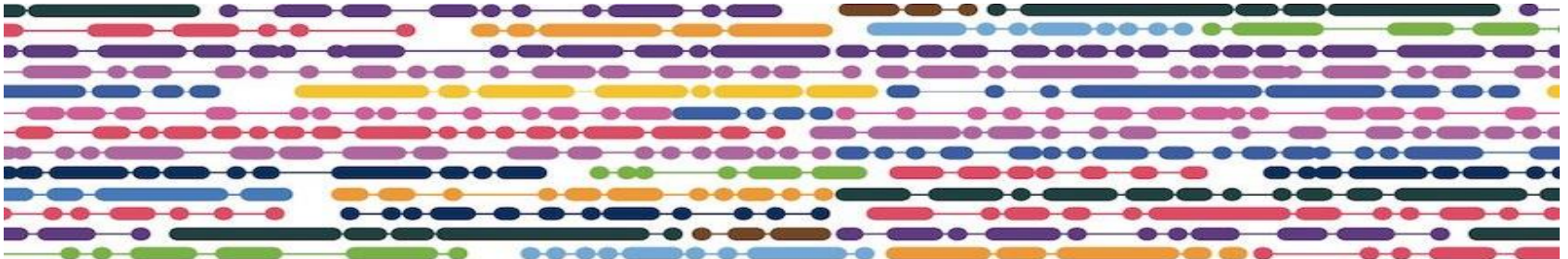
# Linked-read Sequencing

- By adding a unique barcode to every short read generated from an individual molecule, the short reads are linked together.



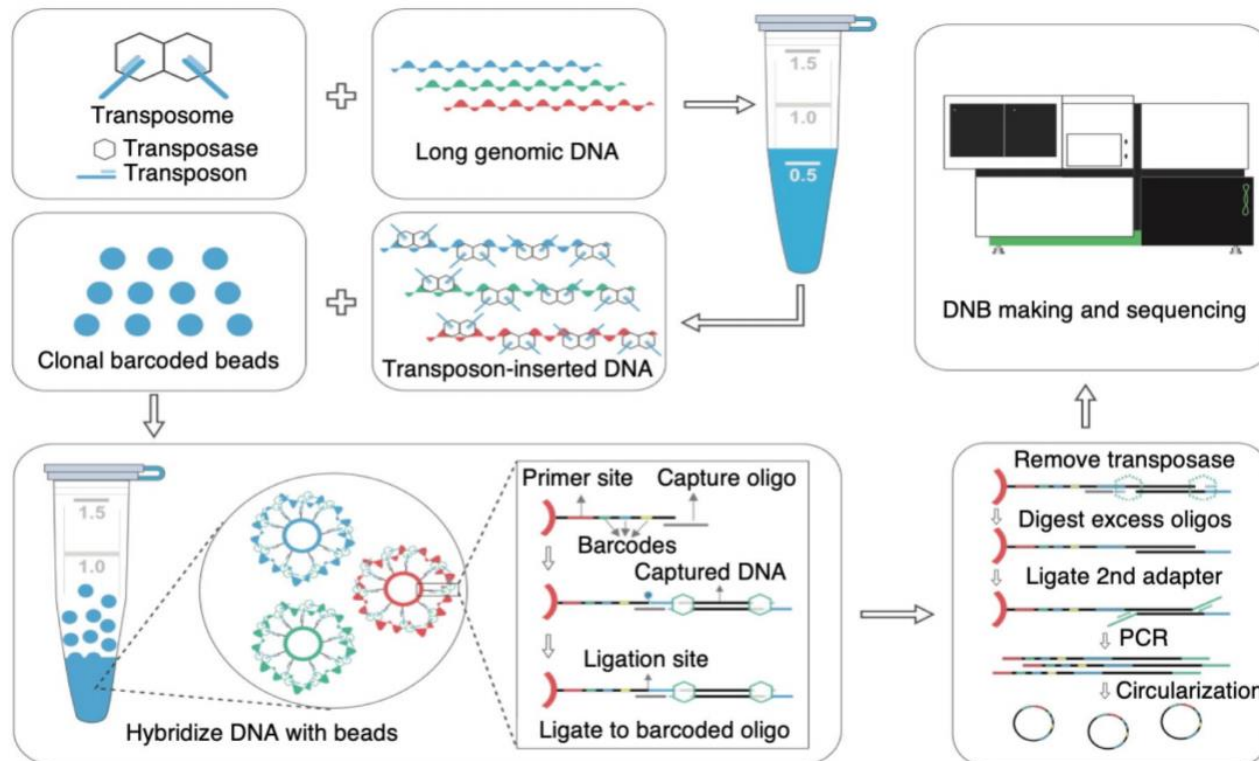
# Linked read use molecular barcoding to preserve long-range information

- 1. Generation of long DNA fragments (weighted mean: ~50 kb)
- 2. The long DNA fragments are randomly dispersed into ~1 million droplet partitions with different barcodes; thus, only a small number (~ 10) of DNA fragments are loaded per partition.
- 3. Short read pairs (2 x150 bp) are generated using barcode-containing primers.
- 4. Short reads that contains the same barcode and within a certain distance can be linked together to “reconstruct” the original long DNA fragment.



# Other linked-reads technologies

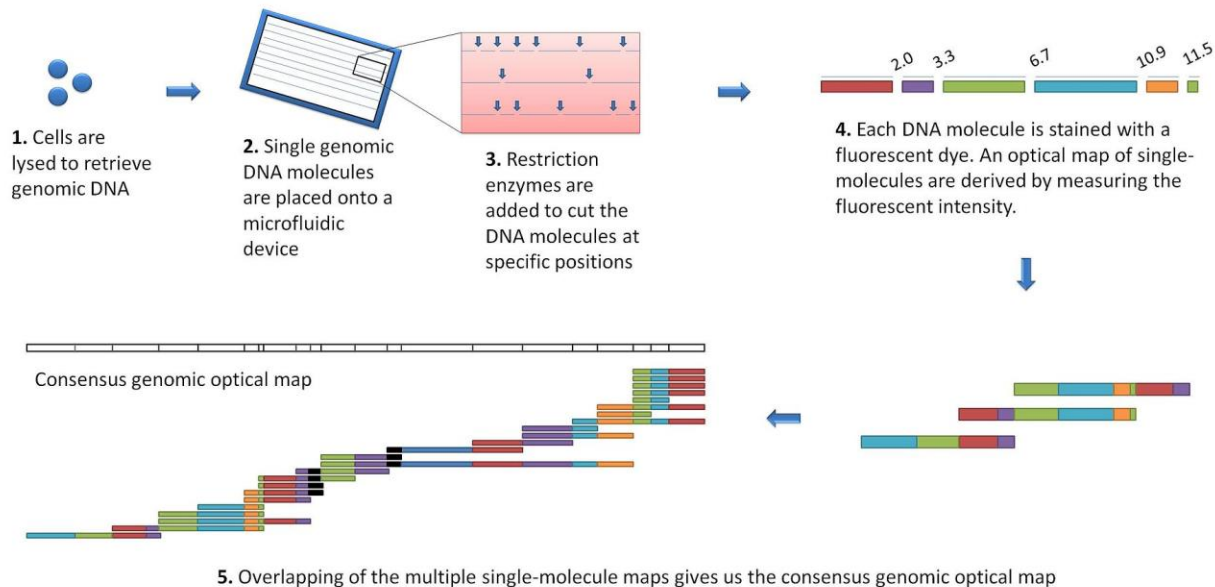
- Single tube long fragment read (stLFR): adding the same barcode sequence to sub-fragments of the original long DNA molecule (DNA cobarcoding).





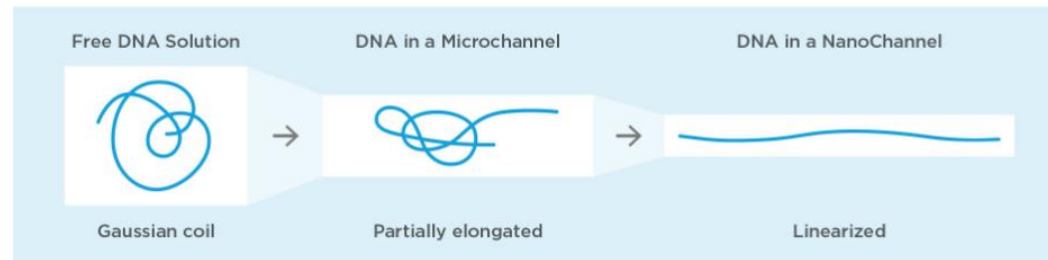
# Optical mapping

- **Optical mapping** is a technique for constructing ordered, genome-wide, high-resolution restriction maps from single, stained molecules of DNA, called "optical maps".



# Single-molecule optical mapping (Bionano Genomics)

- Single DNA molecule linearization in Nanochannel.



- Single DNA molecules are labeled with restriction enzymes. The images are scanned and converted to DNA molecules.

