

# The Chemistry of Advanced Molecular Detection

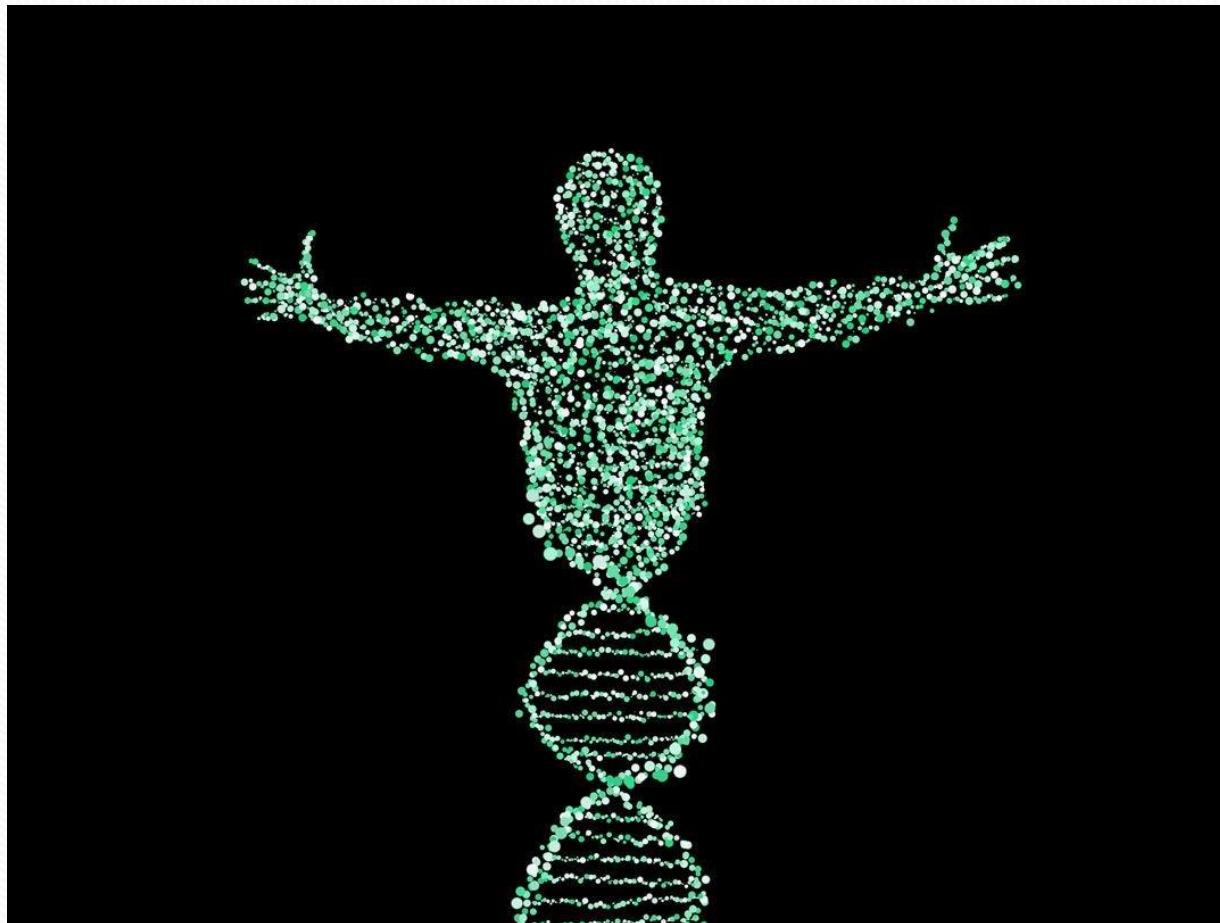
Rebecca Kramer & Kevin Rodeman

Michigan Department of Health and Human Services  
Bureau of Laboratories

Prevent Disease – Promote Wellness – Improve Quality of Life



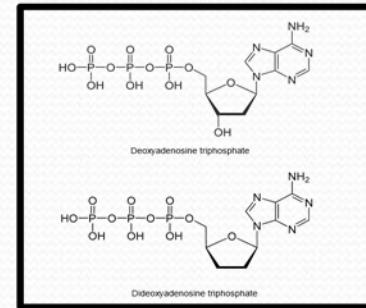
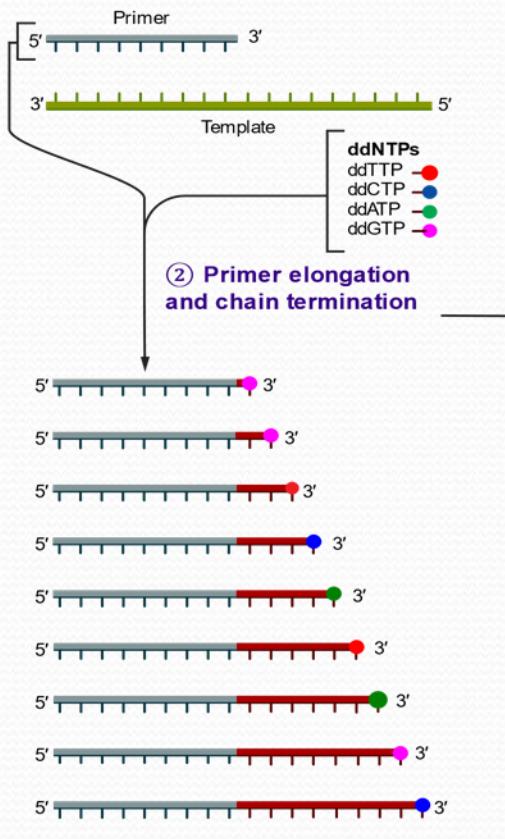
# What is Sequencing?



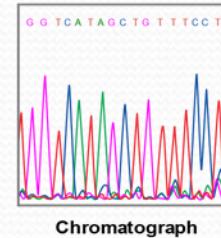
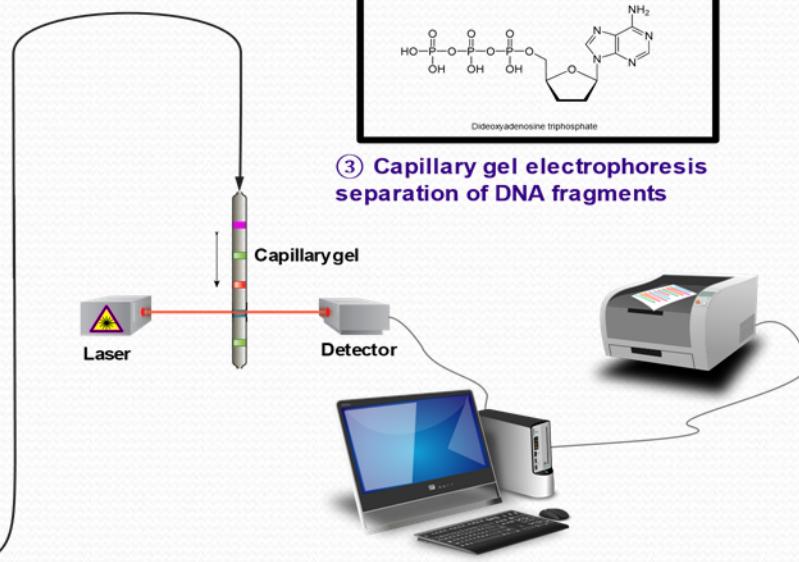
# Sanger Sequencing

## ① Reaction mixture

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



## ③ Capillary gel electrophoresis separation of DNA fragments



## ④ Laser detection of flurochromes and computational sequence analysis

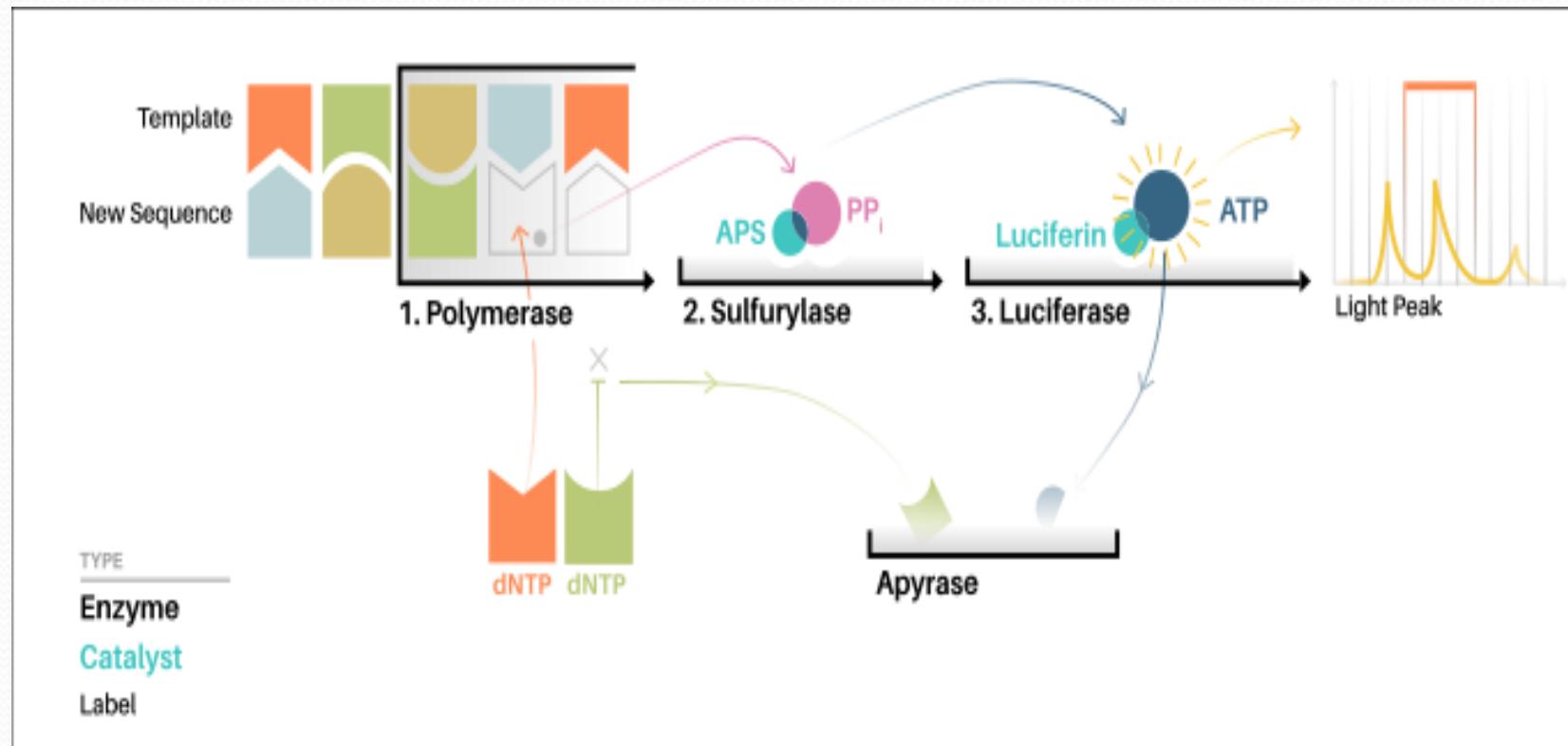
# Sequence Analysis



# Pyrosequencing - the Start of NGS

- Sequencing method based on the “sequencing by synthesis” principle.

# Pyrosequencing Schematic



# Roche 454 (Life Sciences): First Next Generation Sequencer

- Alternative approach by attaching DNA for sequencing on a microfabricated array.
- This allowed for high-throughput DNA sequencing.
- Automated instrument allowed for a new era of genomics research.
- Limitations



# Second Wave of NGS Devices

## Illumina : MiSeq and NextSeq

- Attach indices to each DNA specimen.
- Specimens are pooled, diluted and denatured to create a library.
- The Library is added to the Flow Cell. Bridge amplification creates clusters of clones for each index set (specimen).
- Sequencing by Synthesis is done on the forward strand then on the reverse strand. Similar to pyrosequencing but chemistry is proprietary.
- Run time 4-48 hours; error rate 0.1%; 300bp read length

# Ion Torrent



Ion Torrent™ next-generation sequencing  
[Discover the technology >](#)

# Second Wave of NGS Devices

## Ion Torrent Ion Proton

- Double stranded DNA is fragmented and double stranded adapter is added.
- The double stranded DNA is denatured and a sequencing primer is added.
- Emulsion PCR creates many clones and the amplicons attached to an acrylamide ball at 5' end. Only the forward strands are sequenced.
- The ball is added to a microchip with thousands of wells with each well containing an ion-sensitive field effect transistor.
- As dNTPs are added in waves, when a complimentary base is added a proton is detected as an electrical signal.
- Run time 2-3 hours; 200bp read length; 2% error rate.

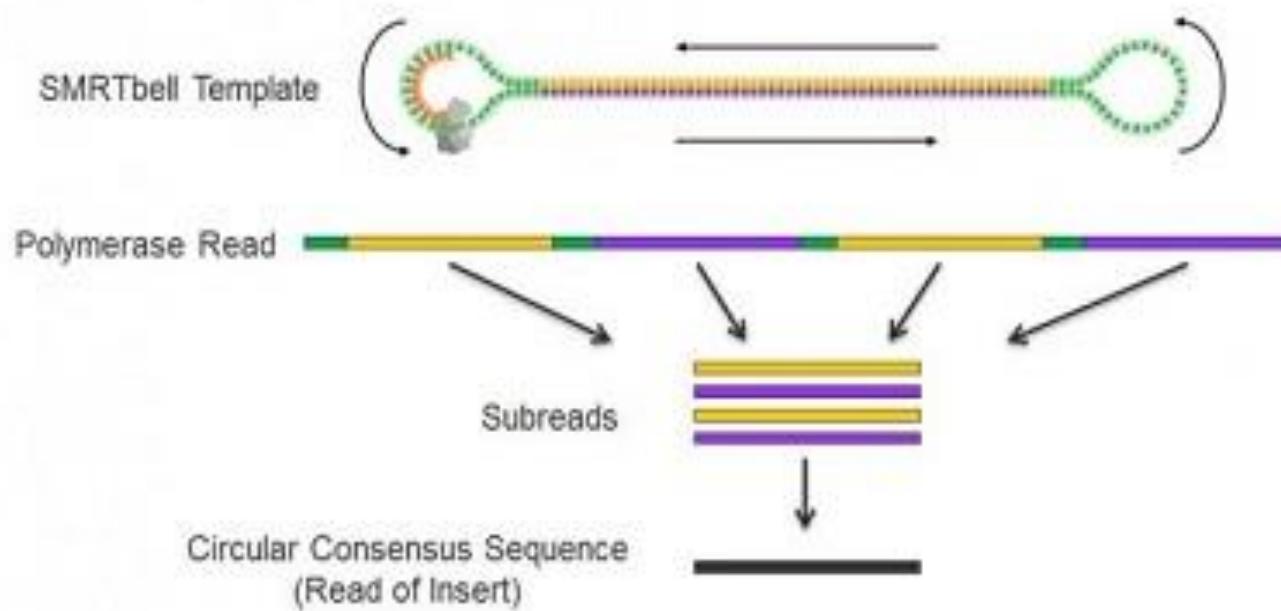
# 3rd Generation Sequencing aka Long Read Sequencing



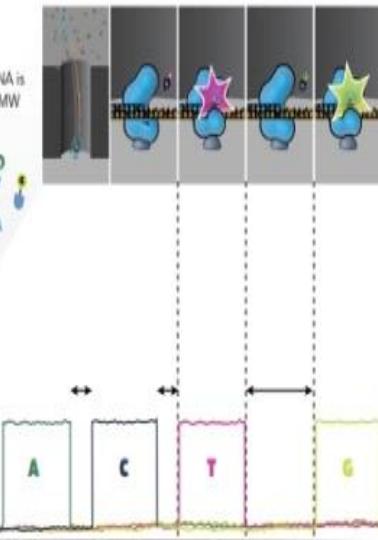
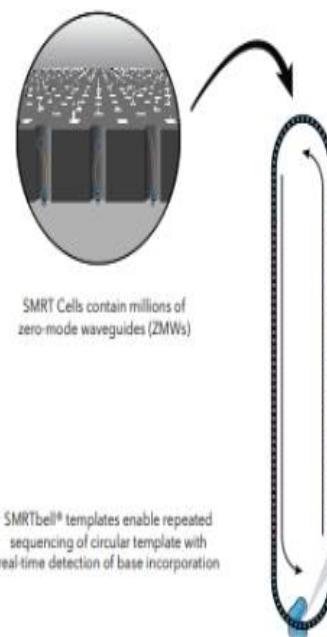
## Third Generation Sequencing: PAC BIO

- SMRTbell™, stream lined protocol to create libraries of varying insert length 250bp to >20,000bp.
- DNA is sheared into fragments of desired size.
- Resulting fragments are repaired by treating the sample with DNA damage repair mix to repair nicks, abasic sites, and oxidation damage.
- Blunt ends are created on each end and the hairpin adapters are ligated to each end.
- In the final step, sequencing polymerase is bound to the template to sequence.

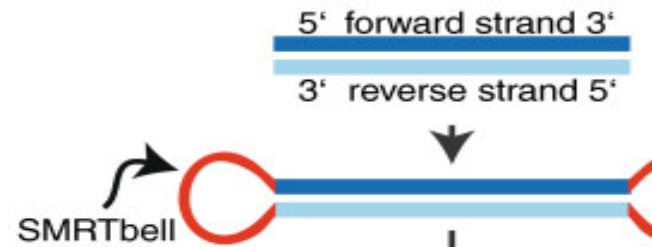
# PAC BIO Sequencing



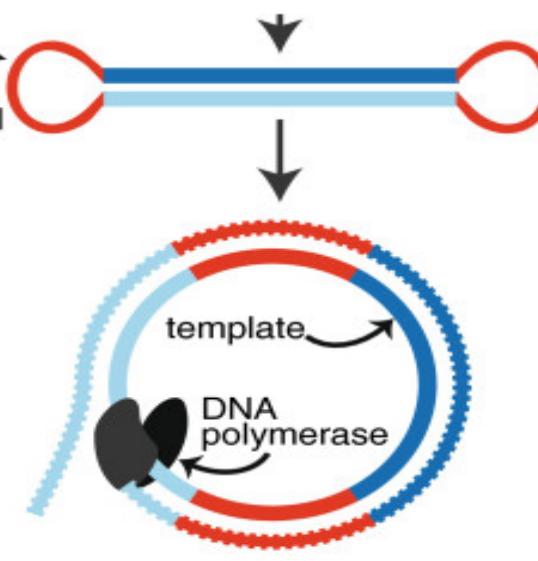
## How SMRT Sequencing Works



## 1. generate amplicon



## 2. ligate adaptors



## 3. sequence

## 4. data analysis

*raw long read*

*processed long read*

*single-molecule fragments*

*circular consensus sequence (ccs)*

*1<sup>o</sup> analysis*

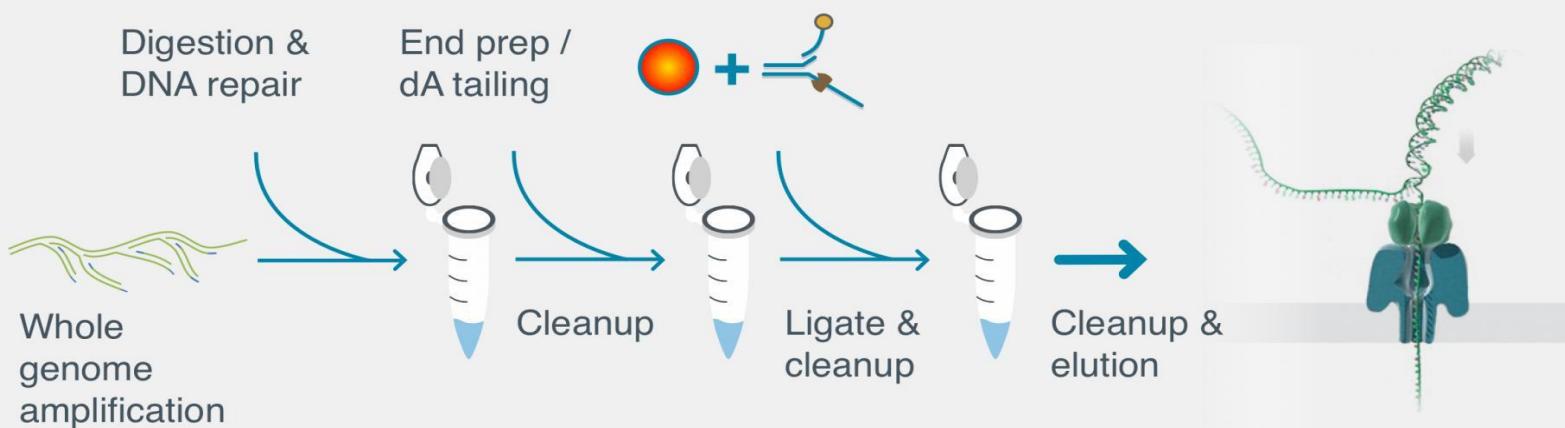
# Oxford Nanopore MinION



# Oxford Nanopore MinION



New: Sequence from as little as 10 pg genomic DNA



# Extractions

Manual versus Robotics



# Instrument Extraction

- Roche: MagNA Pure Compact, MagNA Pure 24 and 96 System
- BIOMERIEUX NUCLISENS® easyMAG® and eMAG™
- QIAGEN: EZ1 Advanced XL, MagAttract 48/96 kits, QIAcube, and QIAsymphony

# Manual Extractions

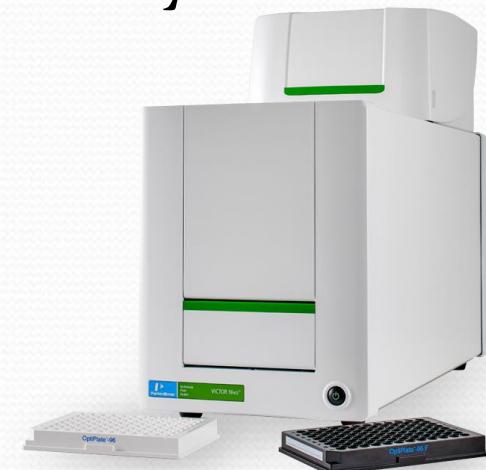
- QIAGEN QIAamp kits: DNA Mini and DNA Blood Mini Kits, and Viral RNA Mini Kits
- Roche: High Pure Kits
- Phenol/Chloroform Extraction
- Nuclease Treatment before extraction and DNase Treatment On-column for RNA Viruses

# QC Checks

- Quality



- Quantity



METHOD	EXAMPLE	BRIEF DESCRIPTION	BENEFITS/LIMITATIONS
Spectrophotometry (260/280)	NanoDrop™	This method detects the absorption of UV light by the macromolecules in the sample.	<ul style="list-style-type: none"> <li>✓ Low cost, as most laboratories already have access to UV/vis spectrophotometers</li> <li>✗ Not specific for DNA</li> <li>✗ Results can be skewed by RNA or protein contamination</li> <li>✗ Cannot determine fragment size</li> </ul>
Fluorimetry	Qubit®	This method measures the enhanced fluorescence of a dye upon binding to DNA/ macromolecules.	<ul style="list-style-type: none"> <li>✓ Low cost, as most laboratories already have access to fluorimeters</li> <li>✓ Can quantitate specifically dsDNA, ssDNA, RNA or protein</li> <li>✗ Quantitates all nucleic acid present in sample, not just molecules to be sequenced</li> <li>✗ Cannot determine fragment sizes</li> </ul>
Electrophoretic	Bioanalyzer®, TapeStation®, Fragment Analyzer™	This method relies on capillary electrophoresis of DNA fragments for size estimation, as well as intercalating dyes for quantity determination.	<ul style="list-style-type: none"> <li>✓ Accurate determination of fragment size distribution</li> <li>✗ Less reliable quantitation</li> <li>✗ Requires expensive equipment</li> </ul>

# Manual Pipetting

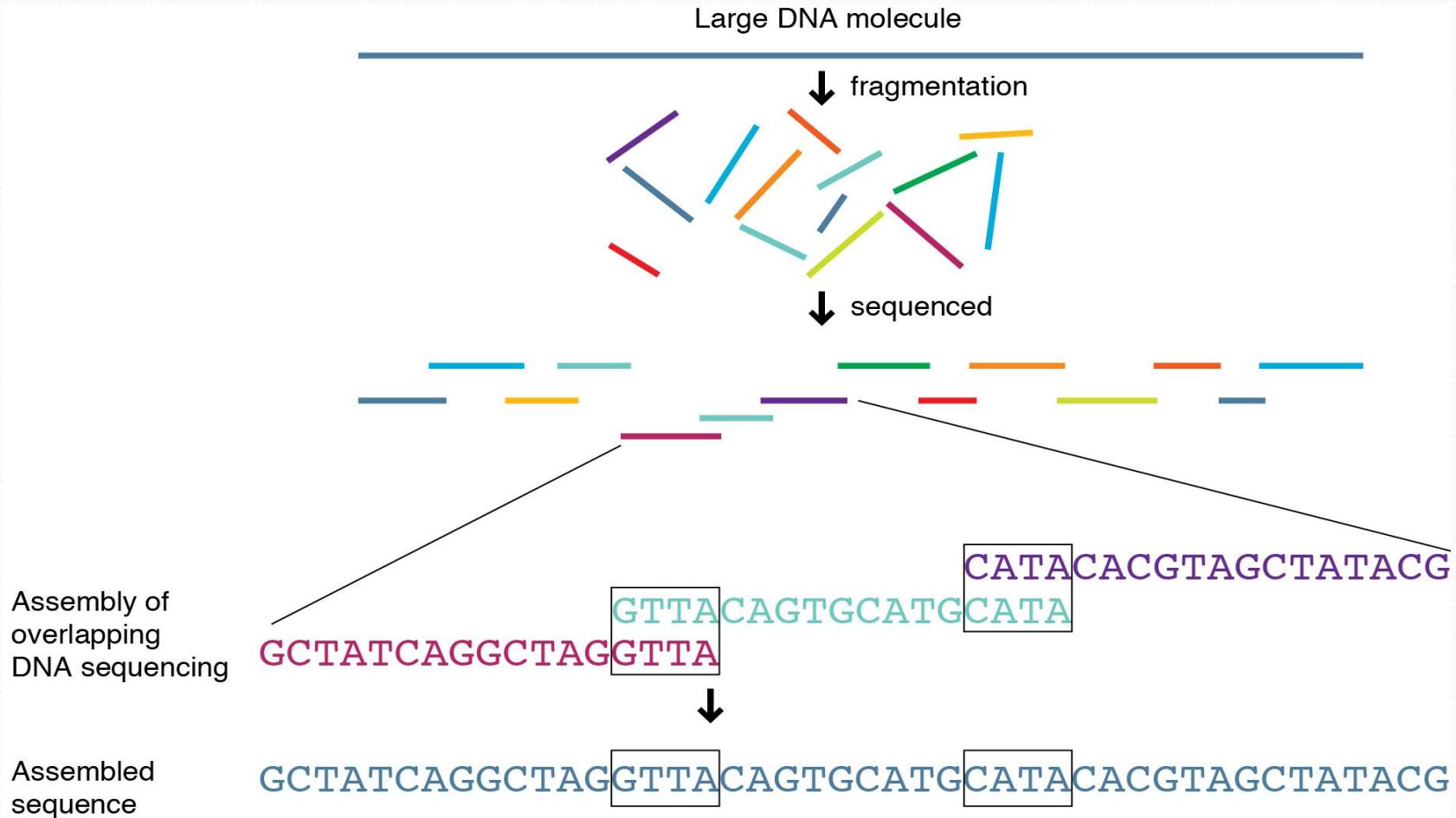
vs.



# Liquid Handler

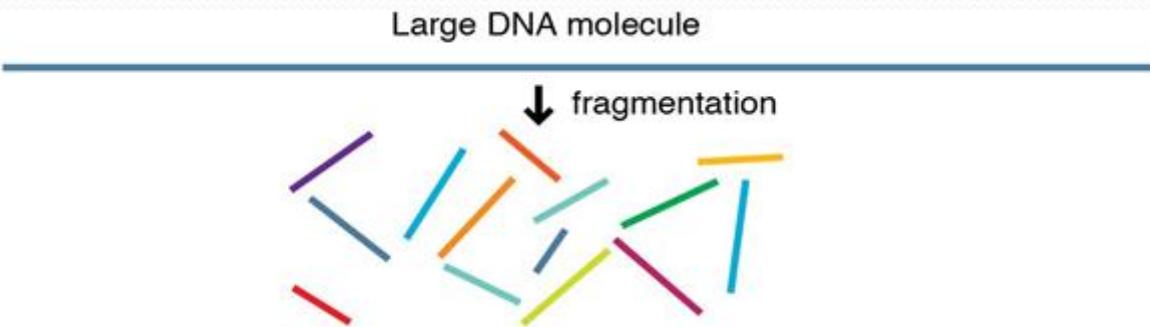
# Basics of NGS Chemistry

1. Library Prep
2. Cluster Generation
3. Sequencing
4. Data Analysis



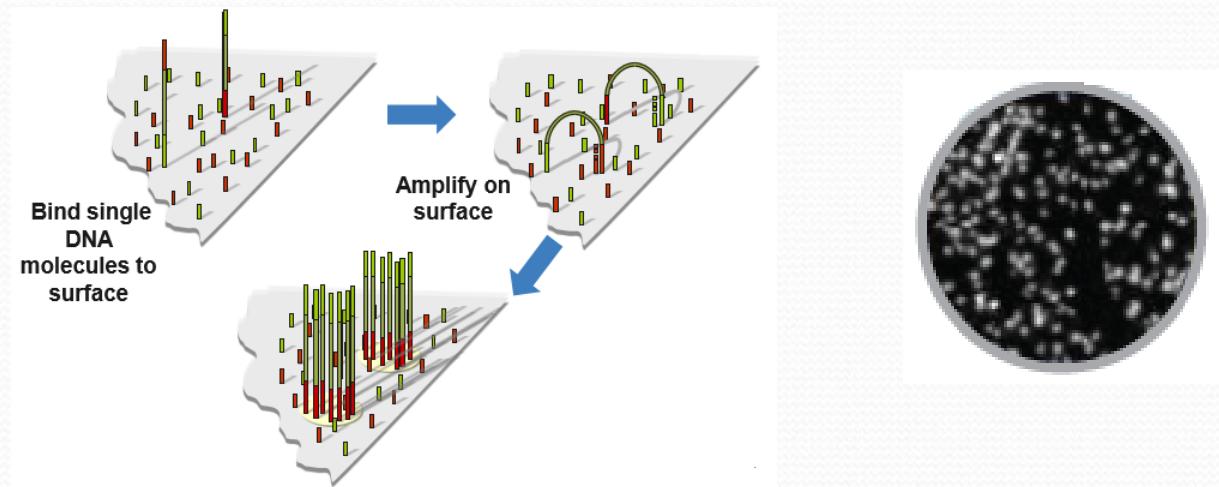
# 1. Library Prep

- Random Fragmentation of DNA
- 5' and 3' Adapter Ligation
- PCR Amplification
- Purification
- Pooling



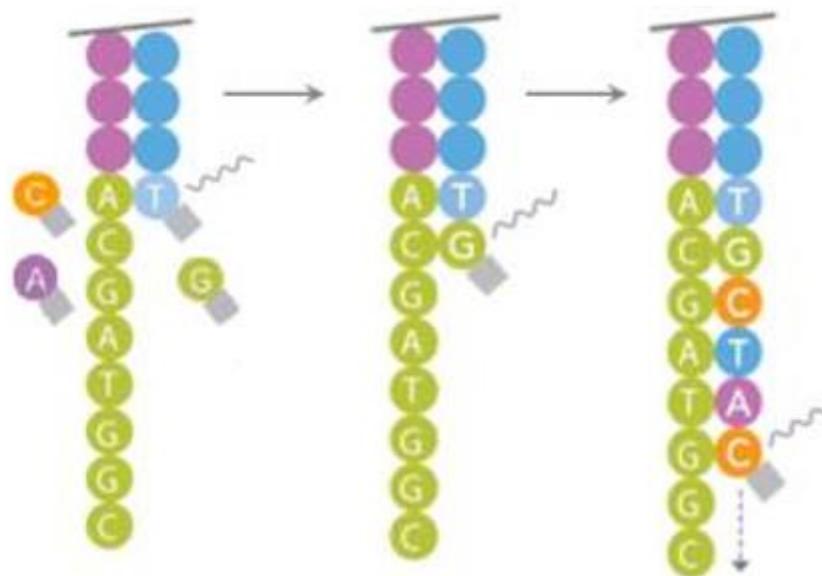
## 2. Cluster Generation

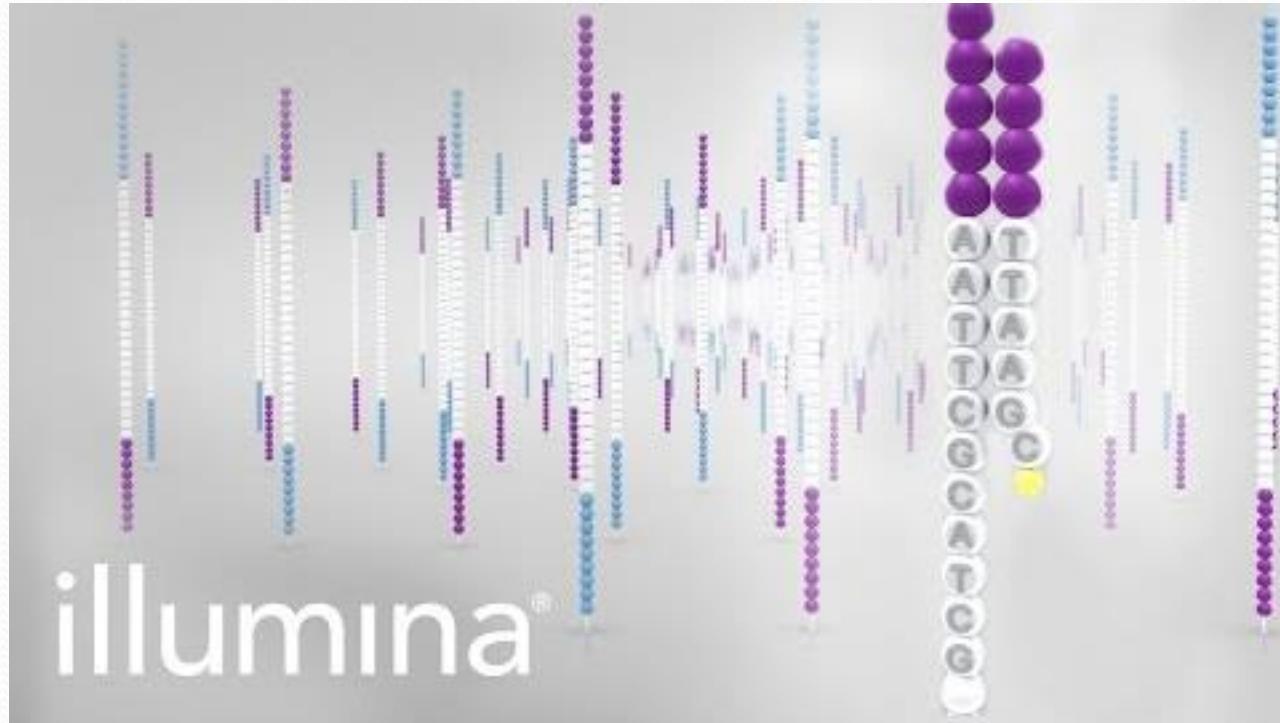
- Library loaded onto flow cell
- Fragments bind to surface bound oligos
- Amplified into clusters through bridge amplification
- Clusters templates ready for sequencing



# 3. Sequencing

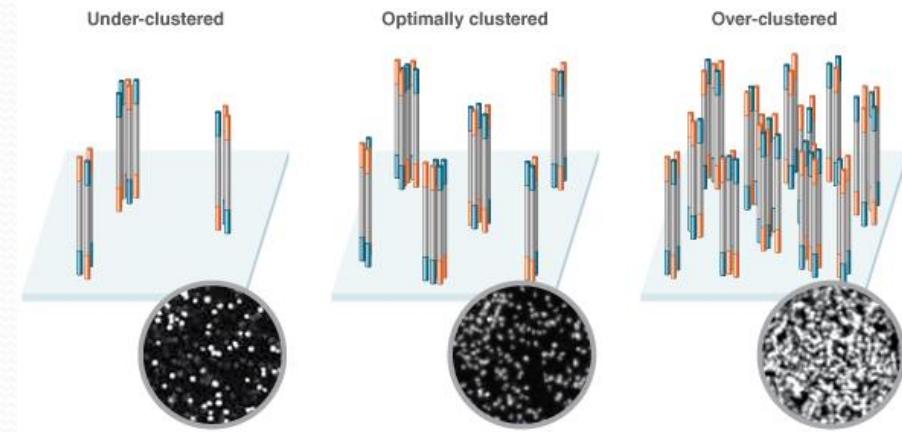
- Sequencing by Synthesis



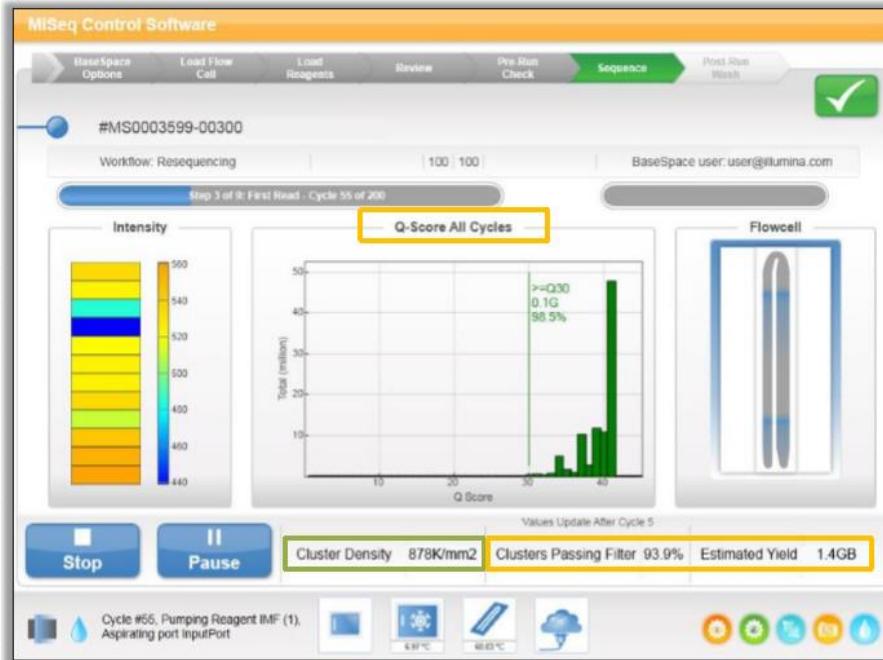


# Run Monitoring

- Cluster Density - under clustering can lead to insufficient data, while over clustering can lead to diminished data quality



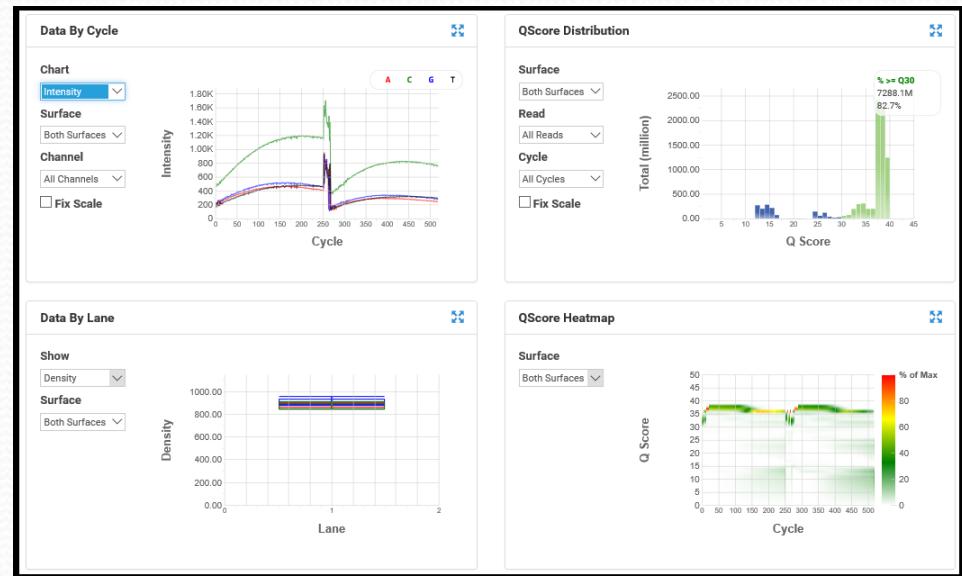
# Run Monitoring (continued)



- Clusters Passing Filter
- Q<sub>30</sub> Score
- Estimated Yield

# Run Analysis

- Sequence Analysis Viewer
  - 1. Q-score
  - 2. Data by Lane
  - 3. Intensity
  - 4. Full Width/  
Half Max
  - 5. Summary Tab

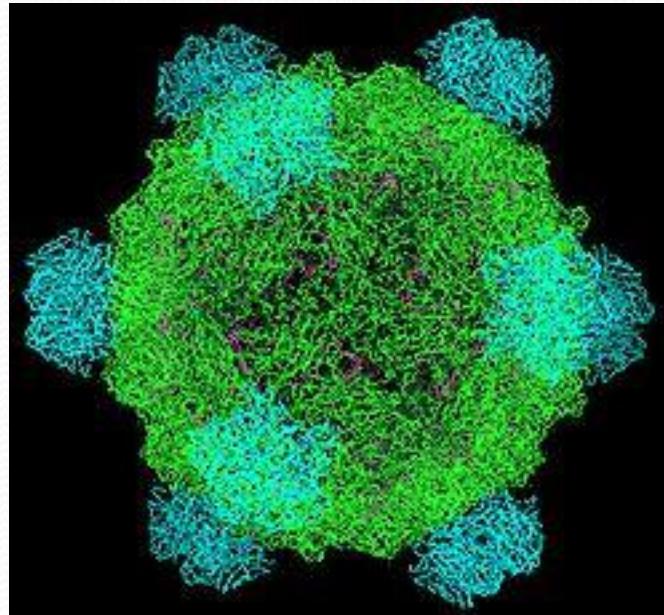


# GHOST MiSeq Run Analysis

<b>GHOST Run No.</b>	<b>Cluster Density (K/mm<sup>2</sup>)</b>	<b>Clusters Passing Filter</b>	<b>Estimated Yield (MB)</b>	<b>Q Score &gt;=Q30</b>	<b>MiSeq Reagent Kit</b>
MI-GHOST04	855	90.3	11547.6	78.4	600 cycles v3
MI-GHOST06	981	93.9	13471	9.7G 87.7%	600 cycles v3
MI-GHOST07	957	92.4	12893.8	9.4G 89.0%	600 cycles v3
MI-GHOST08	726	93.2	6624.5	5.9G 90.4%	500 cycles v2
MI-GHOST09	883	91.5	561.8	0.5G 86.1%	500 cycles v2 Nano
MI-GHOST10	840	94.6	552.7	0.5G 88.8%	500 cycles v2 Nano
MI-GHOST11	844	89.9	523.9	0.5G 90.7%	500 cycles v2 Nano
MI-GHOST12	870	91.9	7697.1	6.8G 90.1%	500 cycles v2

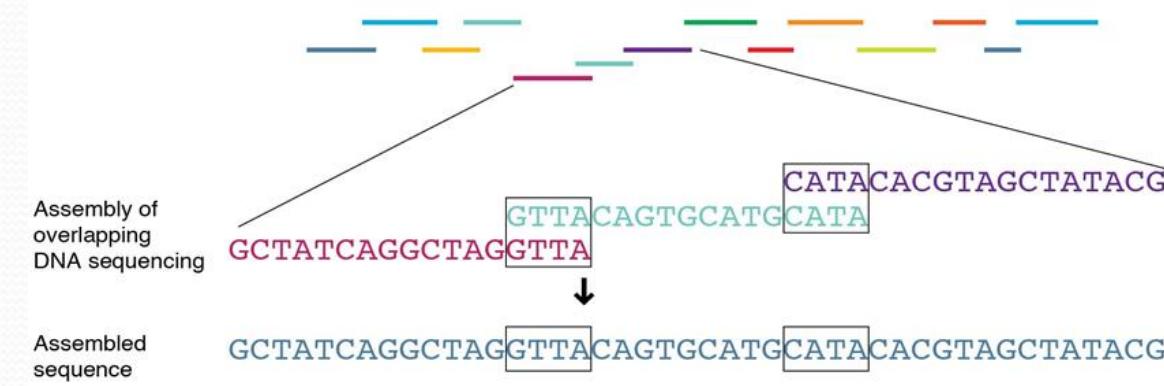
# Troubleshooting

- PhiX – acts as positive control for instrument and reagent performance



# 4. Data Analysis

- Newly identified reads are aligned to reference genome
- After alignment, many analysis options: SNP (single nucleotide polymorphism), Indel (Insertion/Deletion), Read Counting for RNA, Phylogenetic Analysis, Metagenomic Analysis



# Determining Coverage

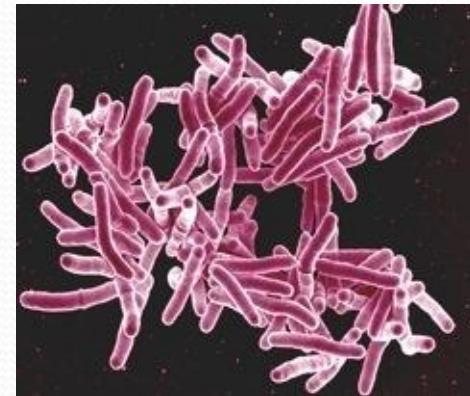
- Fast QC

The screenshot shows the FastQC software interface. The left sidebar lists various quality control metrics with corresponding icons: Basic Statistics (green checkmark), Per base sequence quality (green checkmark), Per tile sequence quality (green checkmark), Per sequence quality scores (green checkmark), Per base sequence content (red X), Per sequence GC content (red X), Per base N content (green checkmark), Sequence Length Distribution (orange exclamation mark), Sequence Duplication Levels (green checkmark), Overrepresented sequences (orange exclamation mark), Adapter Content (green checkmark), and Kmer Content (orange exclamation mark). The main panel displays "Basic sequence stats" for the file "16RF8216-MI-TBWGS-2018-13\_S14\_L001\_R1\_001.fastq.gz". The stats include:

Measure	Value
Filename	16RF8216-MI-TBWGS-2018-13_S14_L001_R1_001.fastq.gz
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	530389
Sequences flagged as poor quality	0
Sequence length	35-251
%GC	65

# Determining Coverage

- Example:
- 530389 sequences
- 2x250 cycle kit
- TB genome – 4.5 million base pairs
- $(530389 \text{ sequences} * 500 \text{ cycles}) / 4500000 \text{ base pairs}$   
= 58.9x theoretical coverage



# Sequencing Platforms



Method	Generation	Read length (bp)	Single pass error rate (%)	No. of reads per run	Time per run	Cost per million bases (USD)
Sanger ABI 3730×1	1st	600–1000	0.001	96	0.5–3 h	500
Ion Torrent	2nd	200	1	$8.2 \times 10^7$	2–4 h	0.1
454 (Roche) GS FLX+	2nd	700	1	$1 \times 10^6$	23 h	8.57
Illumina HiSeq 2500 (High Output)	2nd	2 × 125	0.1	$8 \times 10^9$ (paired)	7–60 h	0.03
Illumina HiSeq 2500 (Rapid Run)	2nd	2 × 250	0.1	$1.2 \times 10^9$ (paired)	1–6 days	0.04
PacBio RS II: P6-C4	3rd	$1.0\text{--}1.5 \times 10^4$ on average	13	$3.5\text{--}7.5 \times 10^4$	0.5–4 h	0.40–0.80
Oxford Nanopore MinION	3rd	$2\text{--}5 \times 10^3$ on average	38	$1.1\text{--}4.7 \times 10^4$	50 h	6.44–17.90

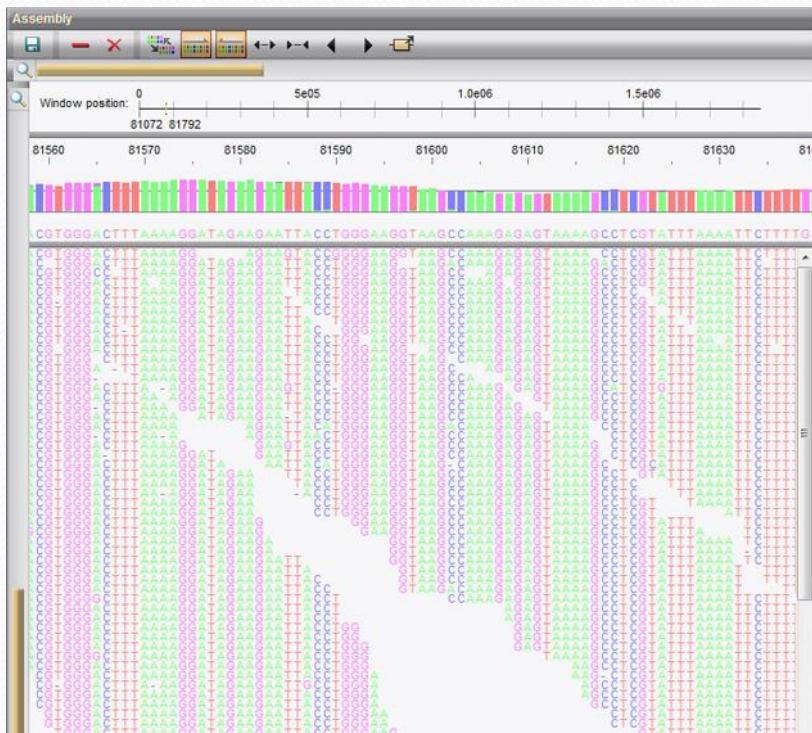


# Software Analysis

## Best Alignment Software?

There is no best read aligner! It depends on each laboratories goals and testing needs. What is the application? What sequencing technology? What is the species?

What are the computational constraints, etc.? One must take into account these questions, each software's features, and determine the suitable read mapper software.



# Any Questions?

