**RU**B

# Next generation sequencing in Virology

Daniel Todt

TRACiR

TRANSLATIONAL AND COMPUTATIONAL
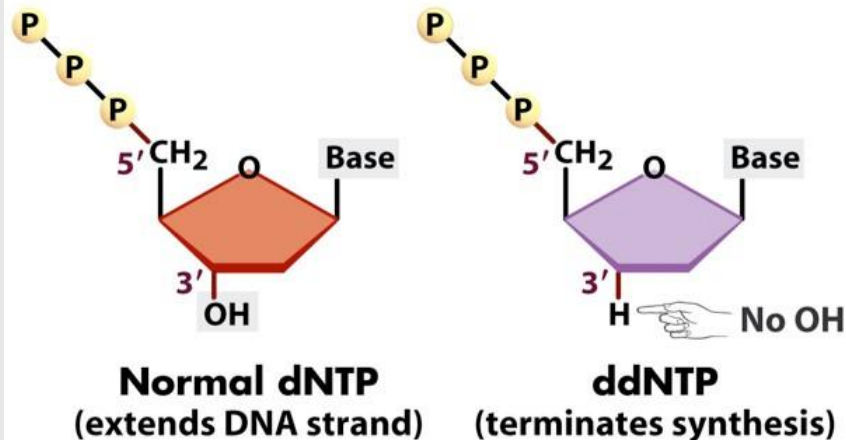INFECTION RESEARCH

Phylogenomics

Quasispecies
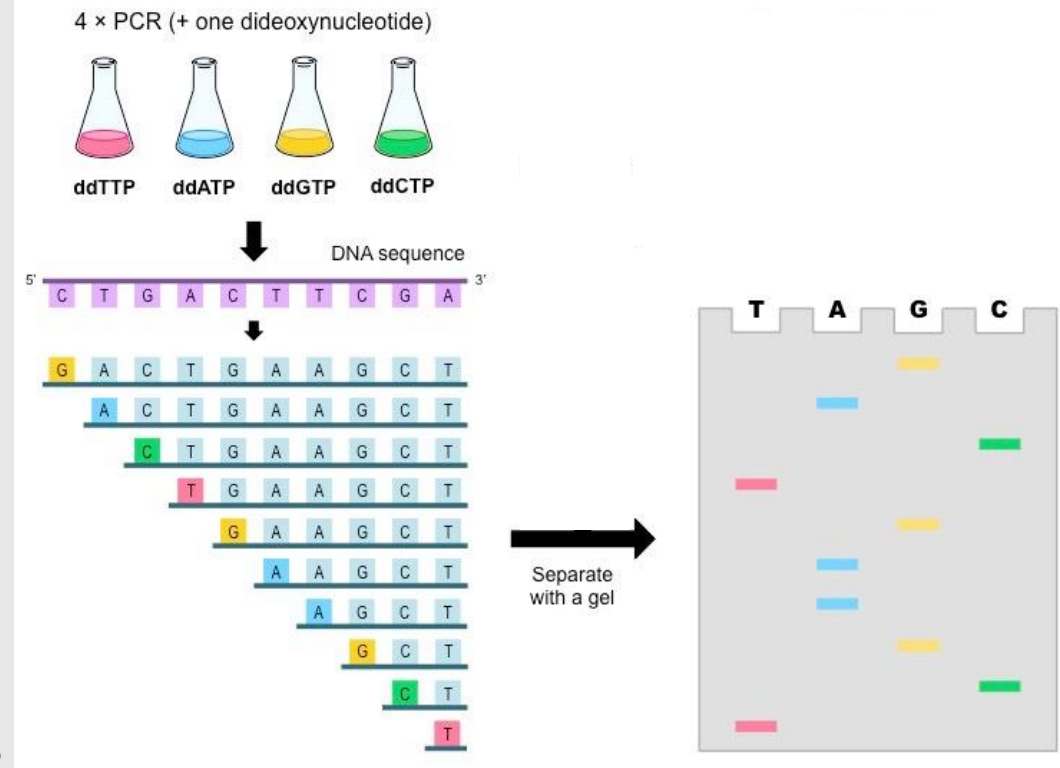
Virocentric Evolution

Algorithms

*NGS*

Recombination

Databases

Selection Pressure

Genotyping

Metagenomics

Population genetics

# First generation sequencing

# Sanger Sequencing
# (dideoxy chain termination)

ddNTPs terminate DNA synthesis.

Normal dNTP (extends DNA strand)

ddNTP (terminates synthesis)

4 × PCR (+ one dideoxynucleotide)

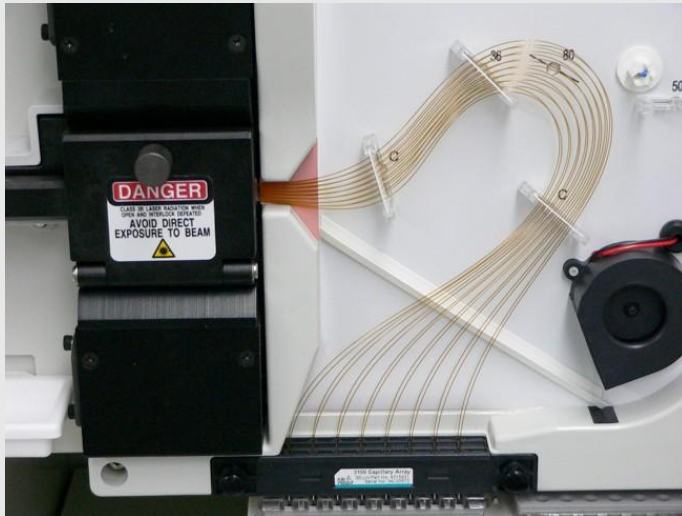ddTTP   ddATP   ddGTP   ddCTP

DNA sequence

Separate with a gel

- Developed in 1977 by Fred Sanger

- DNA extended from **radiolabelled** primers using a mix of dNTP and ddNTP nucleotides

- Random **chain termination** upon ddNTP incorporation

- **Separate reaction** for each terminator (ddC-ddT-ddA-ddG)

- DNA fragments resolved on large **polyacrylamide gels** and detected on film by **autoradiography**

- Sequence **read by hand** and typed in

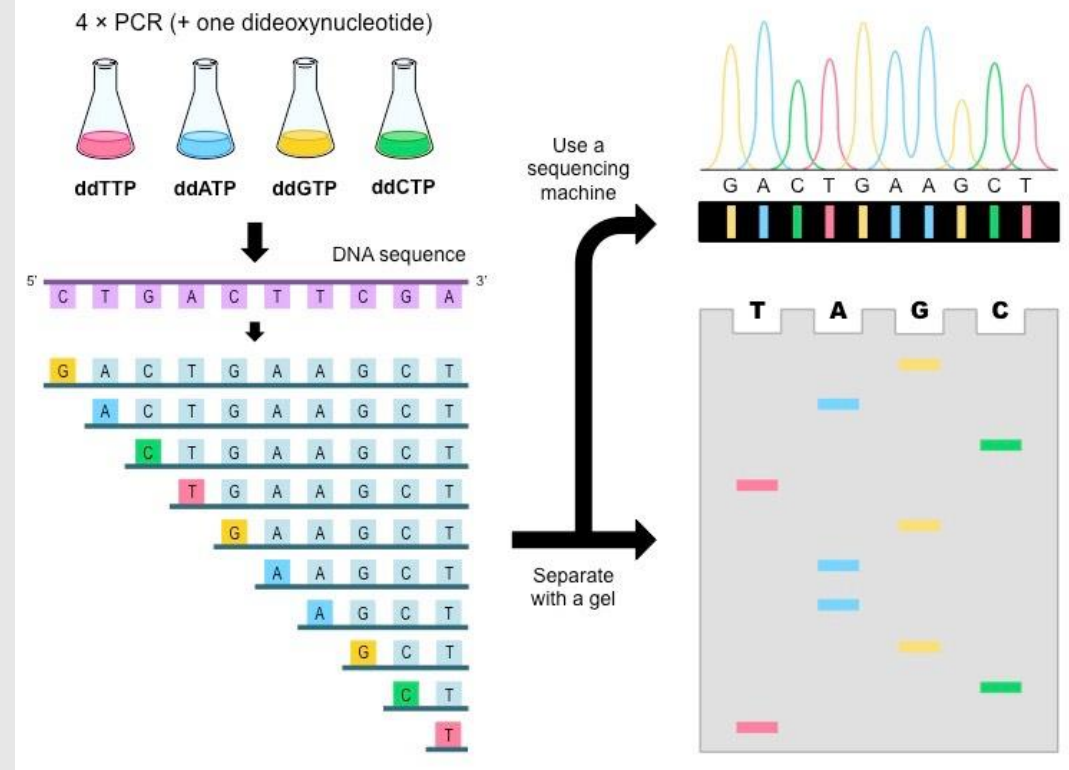- Labour intensive, slow and expensive

3

# Sanger Sequencing

- Automation developed by Leroy Hood and **Applied BioSystems**

- Improved using **fluorescently-labelled** ddNTP terminators (**one reaction** per sequence)

- Separated by **capillary electrophoresis** in automated sequencing machine

- **Long reads (up to 1100 bp) and low error rate**

- Limited throughput, **expensive** per base but still in wide use.

# Second generation sequencing

- General characteristics:

  - DNA molecules from a single library are clustered on a planar substrate (**bridge PCR**), or to the surface of micron-scale beads (**emulsion PCR**).

  - Sequencing by **synthesis** or by **ligation**.

- Advantages over first generation sequencing:

  - *In vitro* **clonal amplification** circumvents time consuming steps such as ligation of DNA fragments into a plasmid, transformation of E. coli and colony picking.

  - Array-based sequencing enables higher degree of **parallelism** than conventional capillary-based sequencing.

# Principal characteristics of the four most used deep sequencing platforms
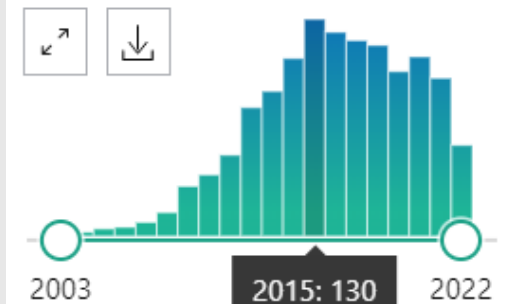
## 454 — GS Jr. / GS FLX+

| | GS Jr. | GS FLX+ |
|---|---|---|
| Amplification Method | Emulsion PCR on beads | |
| Chemistry | Synthesis (pyrosequencing) | |
| Read length (bp) | 400 | 700 |
| Yield/run (Gb) | 0.05 | 0.9 |
| Primary Error | Indel | |
| Error rate | ~1% | |
| Run time (hours) | 10 | 20 |
| Virus-related Publications | 187 | |
| Advantage(s) | Long reads, maturity | |
| Disadvantage(s) | Homopolymer misreads, high cost/Mb | |

## Illumina — MiSeq / HiSeq

| | MiSeq | HiSeq |
|---|---|---|
| Amplification Method | Bridge PCR in situ | |
| Chemistry | Synthesis (reversible termination) | |
| Read length (bp) | 250 | 125 |
| Yield/run (Gb) | 8 | 1,000 |
| Primary Error | Substitution | |
| Error rate | ~0.1% | |
| Run Time (hours) | 39 | 276 |
| Virus-related Publications | 129 | |
| Advantage(s) | Easy work flow, maturity | |
| Disadvantage(s) | Shortest reads, long run | |

## Ion Torrent — PGM / Proton

| | PGM | Proton |
|---|---|---|
| Amplification Method | Emulsion PCR on beads | |
| Chemistry | Synthesis ($H^+$ detection) | |
| Read length (bp) | 400 | 200 |
| Yield/run (Gb) | 2 | 10 |
| Primary Error | Indel | |
| Error rate | ~1% | |
| Run time (hours) | 7 | 4 |
| Virus-related Publications | 13 | |
| Advantage(s) | Low cost, fast run | |
| Disadvantage(s) | Homopolymer misreads | |

## PacBio — RS II

| | RS II |
|---|---|
| Amplification Method | No PCR |
| Chemistry | Single-molecule real-time sequencing |
| Read length (bp) | 8,500 |
| Yield/run (Gb) | 0.15 |
| Primary Error | Indel |
| Error rate | ~13% |
| Run time (hours) | 2 |
| Virus-related Publications | 6 |
| Advantage(s) | Longest reads |
| Disadvantage(s) | High error rate, expensive |

RESULTS BY YEAR

2003 — 2015: 130 — 2022

RESULTS BY YEAR

2012 — 2021: 54
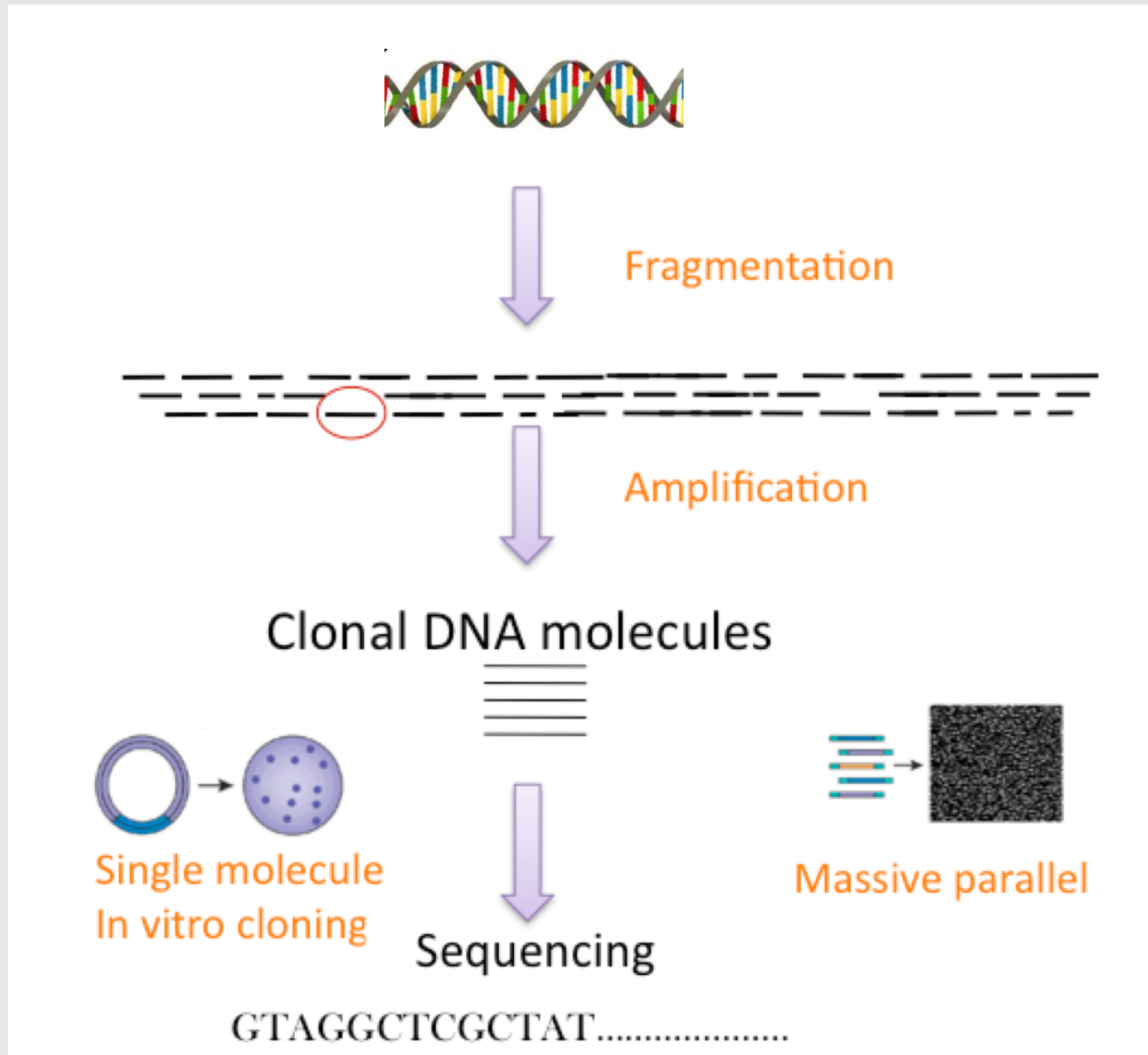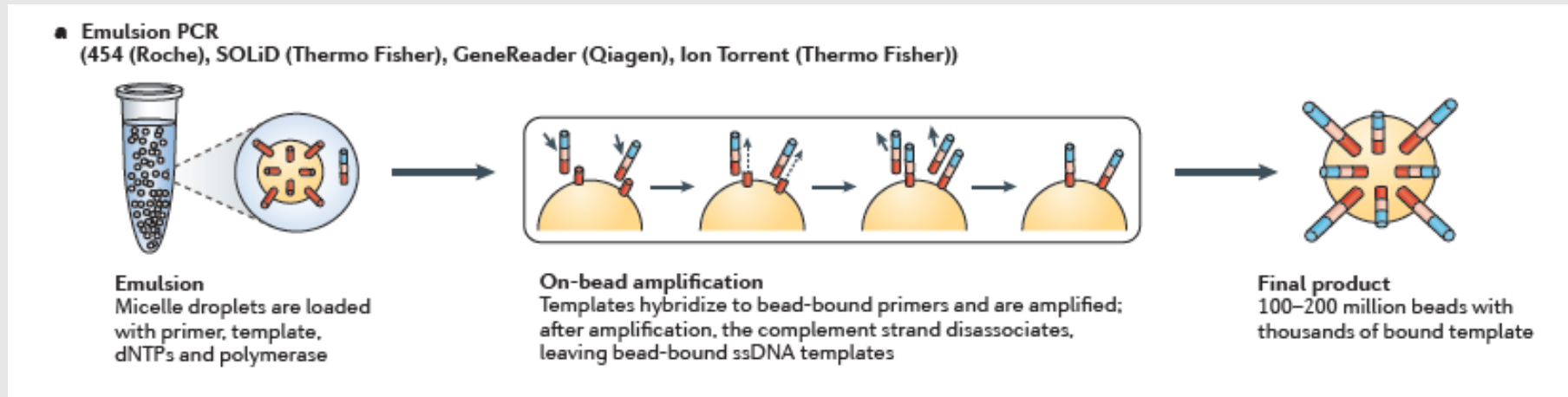
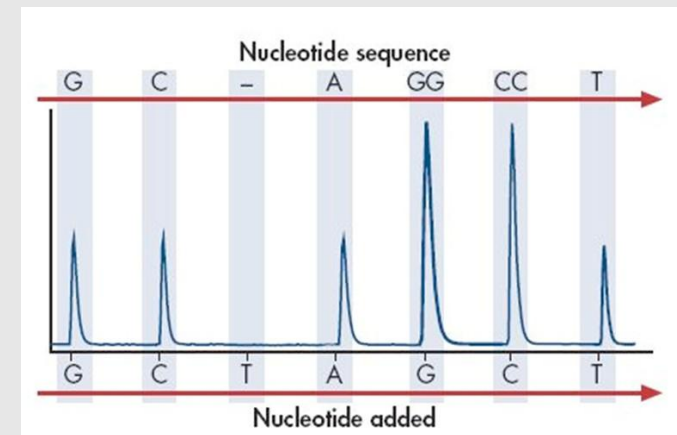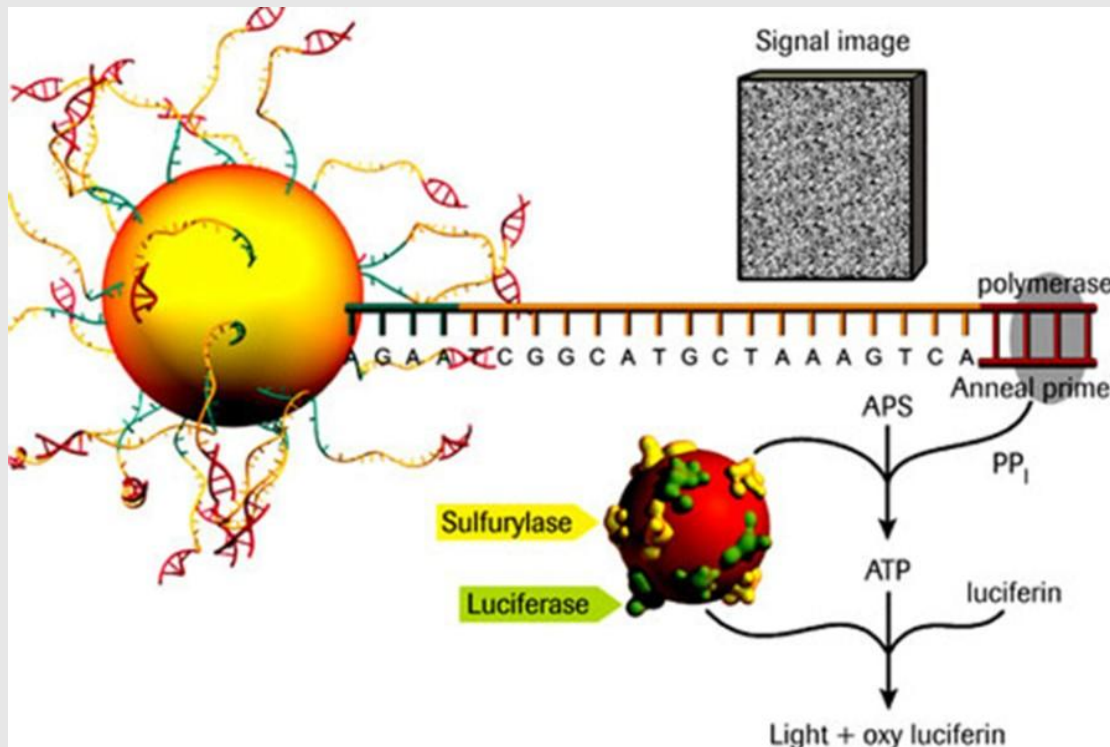# Principal characteristics of the four most used deep sequencing platforms

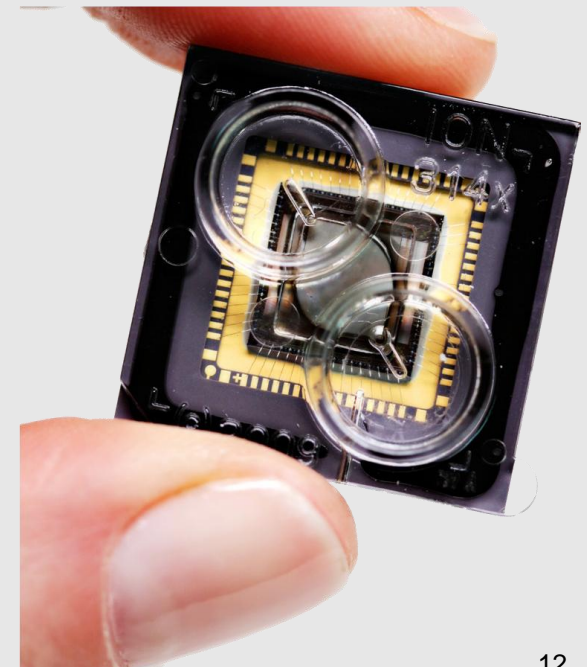| Method | Read length | Accuracy (single read not consensus) | Reads per run | Time per run | Cost per 1 million bases (in US$) | Advantages | Disadvantages |
|---|---|---|---|---|---|---|---|
| Single-molecule real-time sequencing (**Pacific Biosciences**) | 10,000 bp to 15,000 bp | 87% single-read accuracy | 500–1000 megabases | 30 minutes to 4 hours | $0.13–$0.60 | Longest read length. Fast. Detects 4mC, 5mC, 6mA. | Moderate throughput. Equipment can be very expensive. |
| Ion semiconductor (**Ion Torrent** sequencing) | up to 400 bp | 98% | up to 80 million | 2 hours | $1 | Less expensive equipment. Fast. | Homopolymer errors. |
| Pyrosequencing (**454**) | 700 bp | 99.9% | 1 million | 24 hours | $10 | Long read size. Fast. | Runs are expensive. Homopolymer errors. |
| Sequencing by synthesis (**Illumina**) | MiSeq: 50-600 bp HiSeq: 50-500 bp | 99.9% (Phred30) | MiSeq: 1-25 Million; HiSeq: 300 million - 2 billion, | 1 to 11 days, depending upon sequencer and specified read length | $0.05 to $0.15 | Potential for high sequence yield, depending upon sequencer model and desired application. | Equipment can be very expensive. Requires high concentrations of DNA. |
| Sequencing by ligation (**SOLiD** sequencing) | 50+35 or 50+50 bp | 99.9% | 1.2 to 1.4 billion | 1 to 2 weeks | $0.13 | Low cost per base. | Slower than other methods. Has issues sequencing palindromic sequences. |
| Chain termination (**Sanger** sequencing) | 400 to 900 bp | 99.9% | N/A | 20 minutes to 3 hours | $2400 | Long individual reads. Useful for many applications. | More expensive and impractical for larger sequencing projects. This method also requires the time consuming step of plasmid cloning or PCR. |

TRACiR
TRANSLATIONAL AND COMPUTATIONAL
INFECTION RESEARCH



- **Emulsion PCR**
  (454 (Roche), SOLiD (Thermo Fisher), GeneReader (Qiagen), Ion Torrent (Thermo Fisher))

**Emulsion**
Micelle droplets are loaded with primer, template, dNTPs and polymerase

**On-bead amplification**
Templates hybridize to bead-bound primers and are amplified; after amplification, the complement strand disassociates, leaving bead-bound ssDNA templates

**Final product**
100–200 million beads with thousands of bound template

- Developed in Uppsala, Sweden. Later acquired by Qiagen, then licensed to Life Sciences (**454**)

  – DNA fragmentation

  – **Adapters** ligated to DNA fragments (**biotin** tag)

  – Bound to **Streptavidin** beads (each fragment, one bead)

  – Amplified by **emulsion PCR**

  – Beads deposited into **separate wells** on PicoTitrePlate with separate  pyrosequencing reaction in each well, in a large-scale parallel  pyrosequencing system.

# Pyrosequencing

Discontinued in 2014!!!

- **Sequencing by synthesis method**

-  G-C-T-A nucleotides are added sequentially, dNTP incorporation releases **pyrophosphate** (PPi)

- ATP **sulfurylase** converts dNTP to ATP, acts as a substrate for the **luciferase**

- Generates **light** in amounts that are proportional to the amount of PPi (homopolymer error!)

- Unincorporated nucleotides and ATP are degraded by the **apyrase**

- **Light signal** recorded on camera

11

# Ion Torrent



- Licensed from DNA Electronics Ltd, developed by **Ion Torrent Systems**. Later bought over by **Life Technologies**
  - **Adapters** ligated to DNA fragments
  - Bound to **beads**, amplified by **emulsion PCR**
  - Beads deposited into separate wells on **semiconductor chip** with A-T-C-G nucleotides are added **sequentially**
  - **Sequencing by synthesis**
  - Nucleotide incorporation **releases a proton** and the pH of the well changes. A sensing layer detects the change and translates the chemical signal to a digital signal. (Avoids using optical sensors or fluorescent nucleotides; still with homopolymer errors)

- Developed by Balasubramanian and Klenerman who founded Solexa, later acquired by Illumina
  - **Adapters** ligated to DNA fragments
  - **Flow cell** – glass slide with oligos matching adapters
  - Captured DNA replicated through **bridge amplification** to make identical 'colonies'
  - **Fluorescent reversible terminators** passed over flow cell
  - **Image** captured, terminator and dye removed (better performance with homopolymers)
  - barcoding and UMI for multiplexing

13

- Low error rate, read lengths have increased to ≥300 bp.
- Currently used for vast majority of sequencing
- Range of machines with different throughput and cost
- Run time is slower than Ion Torrent (days compared to hours)
- Low error rate – 0.1%
- Single or paired end reads



MiniSeq System   MiSeq Series   NextSeq Series



HiSeq Series   HiSeq X Series   NovaSeq Series

# Illumina

| | NextSeq System | HiSeq System | NovaSeq Series†† | |
|---|---|---|---|---|
| | NextSeq 500* | HiSeq 4000* | NovaSeq 5000* | NovaSeq 6000* |
| Output Range | 20–120 Gb | 125–1500 Gb | 167—2000 Gb | 167—6000 Gb |
| Run Time | 11–29 hr | <1–3.5 days | TBA | 19—40 hr |
| Reads per Run | 130–400 million | 2.5–5 billion | 1.4–6.6 billion | 1.4–20 billion |
| Max Read Length | 2 × 150 bp | 2 × 150 bp | 2 × 150 bp | 2 × 150 bp |
| Samples per Run† | 1 | 6–12 | 4—16 | 4—48 |
| Relative Price per Sample† | Higher Cost | Mid Cost | Lower Cost | Lower Cost |
| Relative Instrument Price† | Lower Cost | Mid Cost | Higher Cost | Higher Cost |
| Downloads | Spec Sheet | Spec Sheet | Spec Sheet | Spec Sheet |

**MGI**

High Speed: 22 hrs ~24 hrs for PE150 sequencing

High Flexibility: 4 Flow Cells, PE150, and PE100 at the same time

Ultra-high Throughput: 7 Tb per day; high quality data around the clock

# Illumina

# Illumina

# Illumina

# Illumina

# Instruments generate short reads that must be mapped to the reference

# Typical screenshot representing aligned HTS reads

# Human genome project



- 1984: Plan – Sanger sequencing

- 1990: Start at National Institutes of Health (NIH)

- 1998: Craig Venter and Celera Genomics

  shotgun sequencing

- 2001: Draft(s) published together with

  Francis Collins of NIH

- 2004: Final published

- Size: ~3 billion base pairs

- Cost: ~$3 billion



Cloned genomes

Multiple genomes are sheared into variable sized segments

Unordered sequenced segments

Computational automated assembly

Resulting overlapping sequence segments. (The higher the coverage the better the quality of the sequencing.

ATGTTCCGATTAGGAAACCTATCTGTAACTGTTTCATTCAGTAAAAGGAGGAAATATAA

Overlapping sequence segments combined to construct the genome consensus.

# Human genome XPRIZE

- **XPRIZES**: intended to encourage technological development that could benefit mankind

- 1996: Ansari XPRIZE for suborbital spaceflight. Claimed by SpaceShipOne in 2004 ($10 million)

- **2006: Archon Genomics XPRIZE: $10 million will be awarded to the first team to rapidly, accurately and economically sequence 100 whole human genomes to an unprecedented level of accuracy** ~~cancelled~~

- 2007: Google Lunar XPRIZE: $20 million to land a rover on the moon, move more than 500 m, and transmit HD images and video back to earth

- 2011: Tricorder XPRIZE: $10 million for a mobile device that can diagnose patients as accurately as a panel of board-certified physicians

# Sequencing costs

| | Sanger 3730xl | 454 GS FLX | Ion Proton | Illumina HiSeq | Oxford Nanopore |
|---|---|---|---|---|---|
| Generation | 1 | 2 | 2 | 2 | 3 |
| Max read length | 1,100 bp | 700 bp | 200 bp | 150 bp | 15,000 bp |
| Max output | 0.1 Mb | 700 Mb | 1,000 Mb | 1,800,000 Mb | 400 Mb |
| Error rate | 0.1% (1:1000) | 1% (1:100) | 1% (1:100) | 0.1% (1:1000) | 25% (1:4) |
| Cost per Mb | £1000 | £5 | £0.05 | £0.005 | £0.50 |



Cost per Human Genome

# 1000 genomes project



1000 Genomes
A Deep Catalog of Human Genetic Variation

- 2008: Launched

- Establish a detailed catalogue of **human genetic variation** correlated with ethnicities

- Sequence 1000 anonymous participants from various ethnic groups within 3 years

- **2012**: 1092 genomes announced

- Each person carries **250-300 loss- of- function variants** in annotated genes

- **50-100 variants** previously implicated in inherited disorders

- Mutation rate of **10-8 per bp per generation** (based on mother-father- child trios)

- 1000 nematode genomes, 1000 plant genomes, Genome 10K project, etc.

**RUHR-UNIVERSITÄT BOCHUM**

**TRACiR**
TRANSLATIONAL AND COMPUTATIONAL
INFECTION RESEARCH

**RUB**

## RNA Sequencing

- mRNA Sequencing
- Targeted RNA Sequencing
- Ribosome Profiling
- RNA Exome Capture Sequencing
- Total RNA Sequencing
- Small RNA Sequencing
- Ultra-Low-Input and Single-Cell RNA-Seq

## Methylation Sequencing

## DNA Sequencing

- Whole-Genome Sequencing
- Targeted Sequencing
- ChIP-Seq
- ATAC-Seq

- **Metagenomics** can be defined as the sequenced-based analysis of the **whole collection of genomes** isolated directly from a sample
- The advantage is that **isolation is not needed** – only extraction and sequencing (although there's more to it than that!)
- Bacteria and archaea: 16S rRNA gene, relatively short, often conserved within species, and generally different among species
- **Viruses**: often present with a large excess of host DNA, making their efficient and reliable detection problematic

# Metagenomics methods

extraction

Sequencing
quality control

*De novo*
assembly

BLAST

Species
A

Species
B

Species
C

Species
D

Species
E

Species
F

Species
G

# Metagenomics – detection



Analysis of samples collected at Penn Station on one day, compared at each hour

- Saunders et al. (2012): Geospatial resolution of human and bacterial diversity with city-scale metagenomics

# Metagenomics – detection



Geospatial analysis of the most prevalent genus, *Pseudomonas*, across the subway system

- Saunders et al. (2012): Geospatial resolution of human and bacterial diversity with city-scale metagenomics

# Metagenomics – detection



**Scientists Basically Just Discovered Alien Life — In The NYC Subway**

Nearly half of the germs on the train are unrecognizable even to the experts.

- Saunders et al. (2012): Geospatial resolution of human and bacterial diversity with city-scale metagenomics

# Metagenomics – virus discovery



- Lauber & Seitz et al. (2017): Deciphering the Origin and Evolution of Hepatitis B Viruses by Means of a Family of Non-enveloped Fish Viruses

# Epidemics – Ebola

- The 2013-2015 West Africa  Ebola epidemic, 26648 cases, 11017 deaths
- HTS used throughout the  epidemic to sequence Ebola  virus genomes from patient  samples
- Used to monitor viral  evolution: how fast is it  mutating, where is it  mutating, which selection  pressures are operating

Can then combine with epidemiological information – date, location, etc.

# Epidemics – who infected whom?

- Identify source of infection

- Identify long transmission events

- Identify super-spreaders – individual  or hub level

- Identify new incursions or spillovers

# Viral populations

- Viruses **mutate rapidly**

- A single virus can enter a cell, and output tens of thousands of virions within hours

- Every time the genome is copied, **mutations** are introduced

- Enables viruses to **adapt** to change rapidly

- New environments

- New hosts

- Drug and vaccine treatment

- Viruses exist as a large and constantly and rapidly evolving swarm – the **quasispecies**

*swarm theory*

# Bottlenecks and Founder Effect



Genetic Drift—Bottleneck Effect

Parent population → Bottleneck (drastic reduction in population) → Surviving individuals → Next generation

Initial founder virus

Expanded Quasi-species

# Viral mutation tracking

- Ability to detect mutations at **low levels** in a sample

- Can then examine samples for the presence of important mutations: e.g. **drug resistance**

  ➢ Hepatitis E virus



- clinical application in HIV diagnostics!

# Third generation sequencing

# Third generation sequencing



- Advantages over second generation sequencing:
  - Very long reads (Oxford nanopore)
  - Real time output
  - scRNAseq (10x Genomics)

# Oxford nanopore sequencing (ONT)

2016 - Kate Rubens becomes the first person to ever sequence in space

investigated the effects of microgravity on RNA isolation and PCR analysis



**+ she is a virologist!!!**

# Oxford nanopore sequencing (ONT)



A brief history of nanopore sequencing

# Oxford nanopore sequencing (ONT)

Professor David Deamer's initial sketch for sequencing DNA using a nanopore

# Oxford Nanopore

- **'Strand sequencing'** is a technique that passes intact DNA polymers through a protein nanopore, sequencing in real time as the DNA translocates the pore
- Simple sample preparation
- Nucleotide base detected as passes through pore (median kmers 5nt)
- Very long reads, up to 15,000 bp
- Small and portable devices useable in field studies (**MinION**), benchtop system for high throughput (**PromethION**) and for use with mobile devices (**SmidgION**)
- High error rate



**Ab** Oxford Nanopore Technologies

**Leader–Hairpin template**
The leader sequence interacts with the pore and a motor protein to direct DNA, a hairpin allows for bidirectional sequencing

**Motor protein**

**Alpha-hemolysin**
A large biological pore capable of sensing DNA

**Current**
Passes through the pore and is modulated as DNA passes through

**ONT output (squiggles)**
Each current shift as DNA translocates through the pore corresponds to a particular k-mer

# Oxford Nanopore



https://nextstrain.org

49

# Oxford Nanopore

RUB



https://nextstrain.org

# Single Cell sequencing





**Timeline of single cell RNA-seq**

| 2012 | 2013 | 2014 | 2015 | 2017 | 2019 | 2021 |
|------|------|------|------|------|------|------|

| 100 cells | 1,000 cells | 10,000 cells | 100,000 cells | Human Cell Atlas |
|-----------|-------------|--------------|----------------|------------------|
| Smart-seq | MARS-seq | Drop-seq | Drop-seq | FET Flagship |

# A single cell and single nucleus atlas of COVID-19 lung

# A single cell and single nucleus atlas of COVID-19 lung

Biological systems are complex – Tissue Heterogeneity

~250,000 single cells from >40 mouse tissues

# Spatial Transcriptomics



2022
Spatial Multiomics

2020
Spatially Resolved
Transcriptomics

2022
1st Visium Cover

2020
Visium Spatial Gene
Expression

## 200+ Visium Publications and Preprints

# Spatial Transcriptomics

## analysis of Human Lymph Nodes

# Spatial Transcriptomics

## analysis of Human Lymph Nodes

Organ-specific gene expression in mouse embryo

# Thank you for your attention !!