

Introduction to Next-Generation Sequencing Technologies

Javier Santoyo-Lopez

javier.santoyo@ed.ac.uk

22nd October 2015

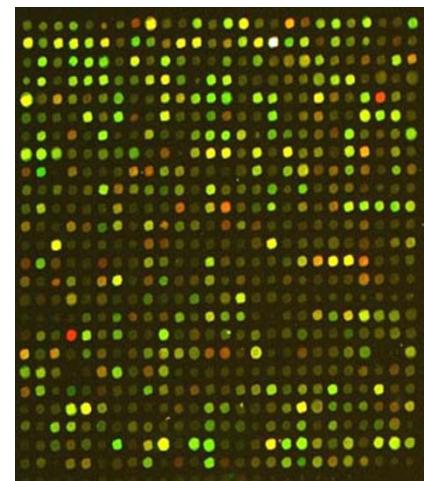
Outline

- ❖ Introduction to high-throughput technologies
- ❖ Description of the NGS technologies
- ❖ Anatomy of an NGS library
- ❖ NGS for variant calling

Introduction to high-throughput technologies

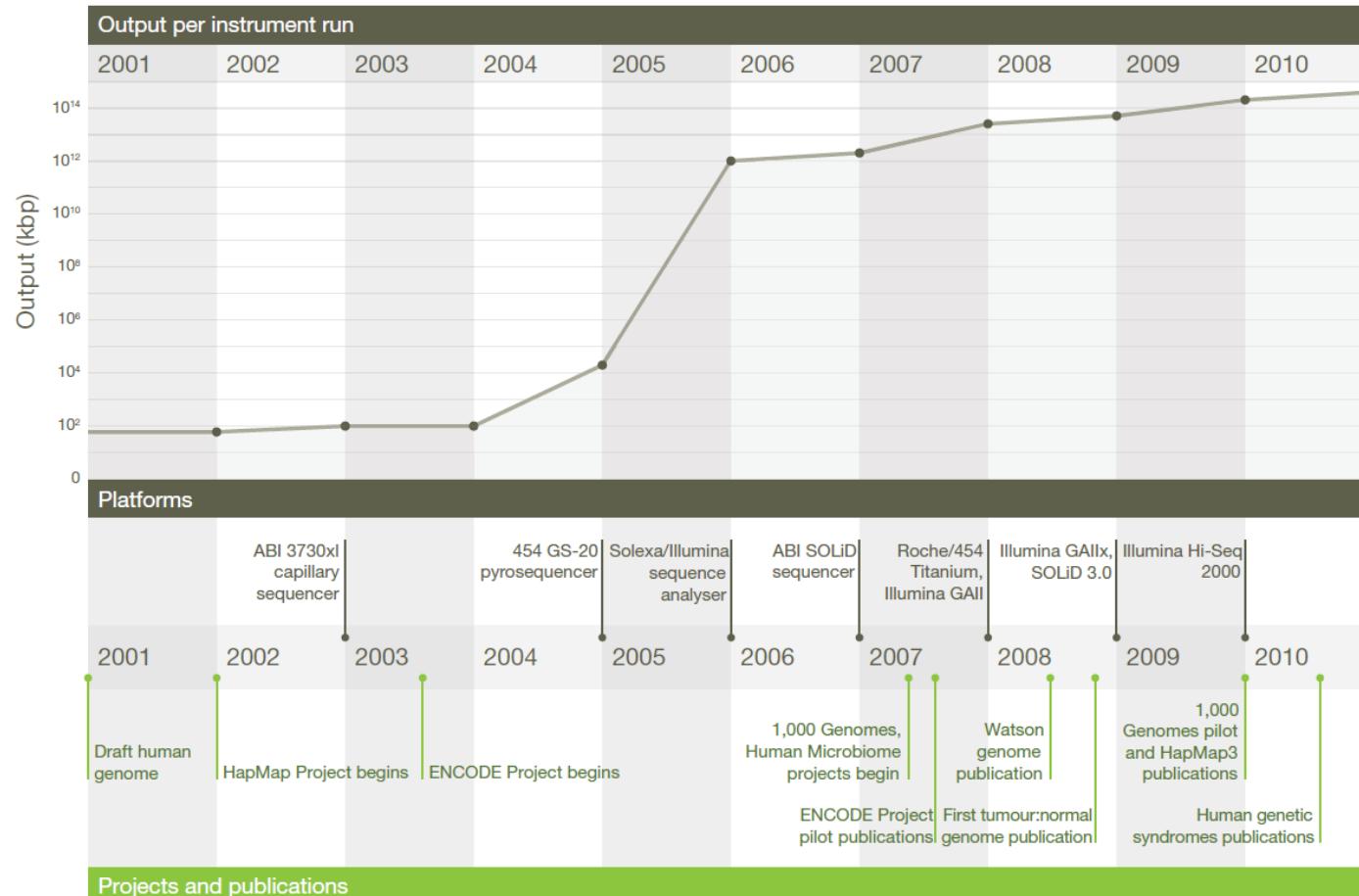
First HT Technology

- 1988 arrayed DNAs were used
- 1991 oligonucleotides are synthesized on a glass slide through photolithography (Affymax Research Institute)
- 1995 DNA Microarrays
- 1997 Genome wide Yeast Microarray



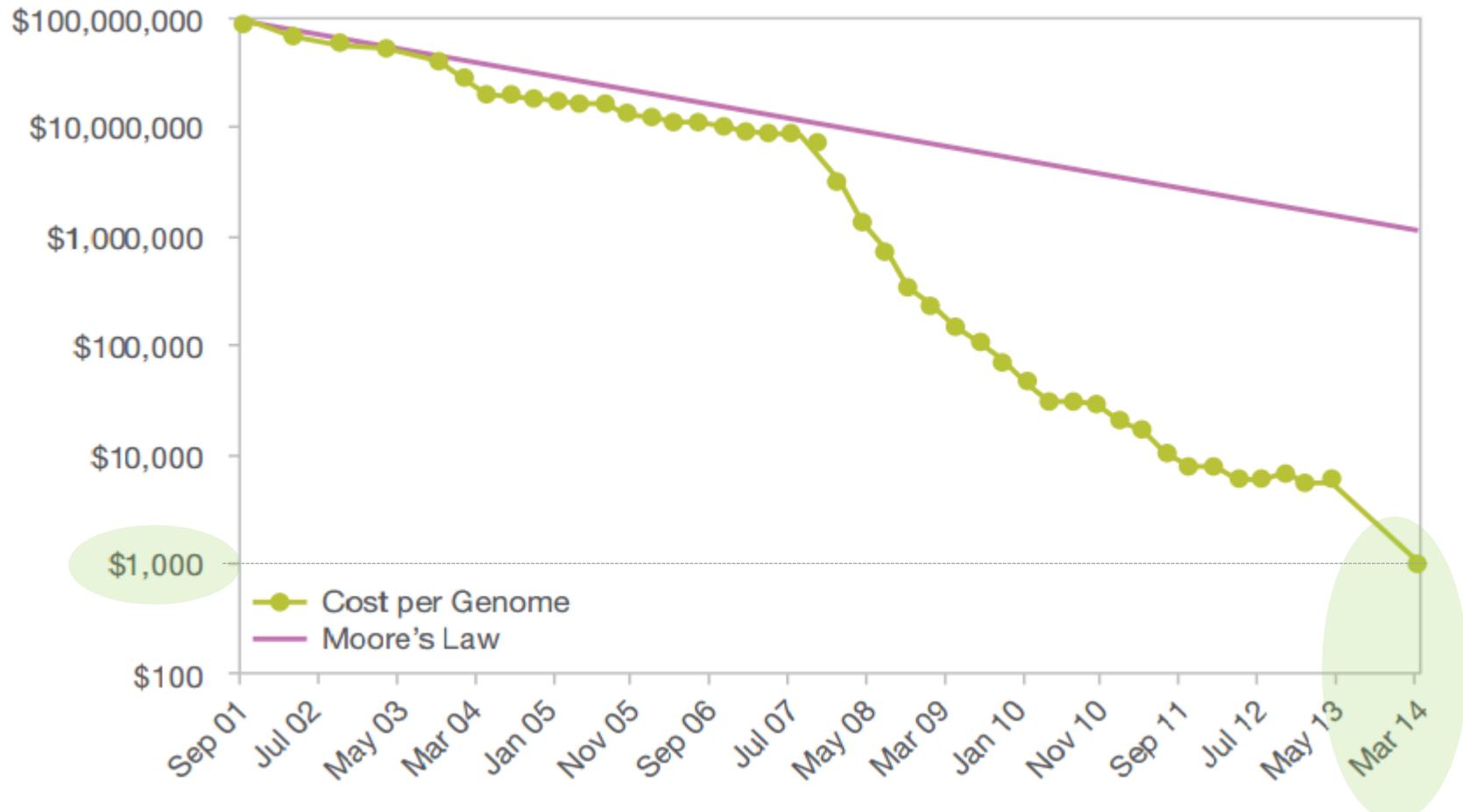
Milestone of DNA Technologies

Projects & sequence output



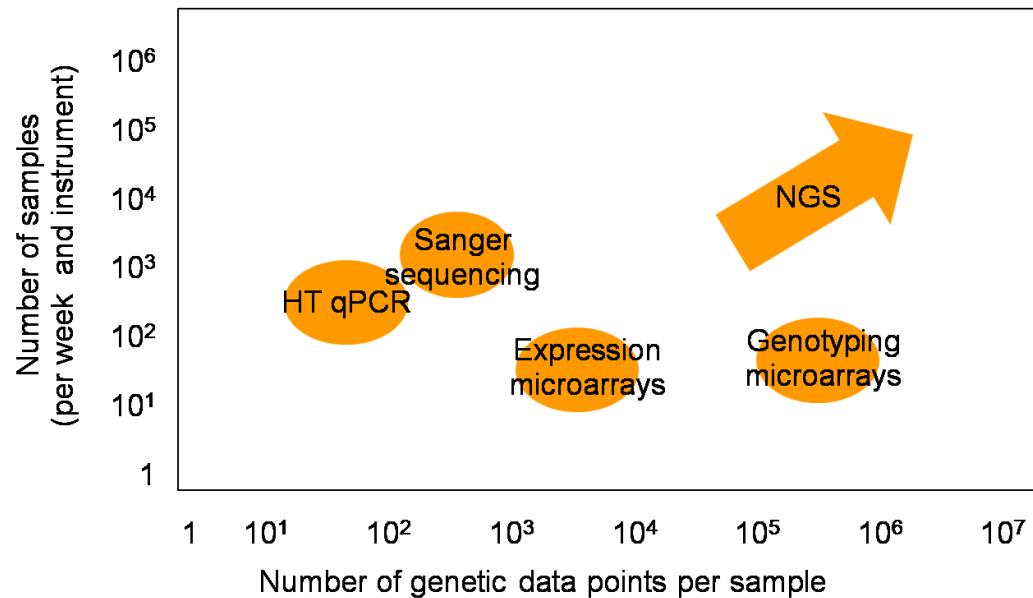
E.R. Mardis, Nature (2011) 470:198 - 203

Genome Sequencing Cost per Mb (30x)



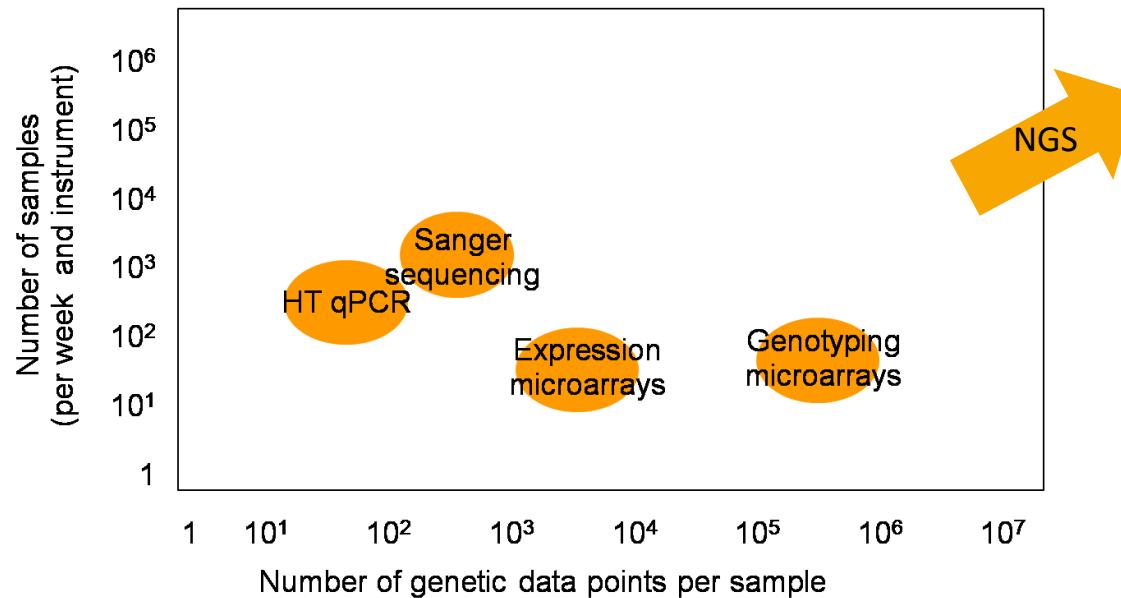
Relative throughput of HTT

Next Generation Sequencing emerges with a potential of data production that will, eventually wipe out conventional HT technologies in the years coming



Relative throughput of HTT

Next Generation Sequencing emerges with a potential of data production that will, eventually wipe out conventional HT technologies in the years coming



NGS: Too many sequences to be handled in standard hardware

Many Gbs (Tbs) of Sequences

- Data management becomes a challenge.
 - Moving data across file systems takes time (several hundred Gbs)
- What quality does the raw data have?
 - Sequencers provide quality values for each bp
- How to do Analysis of the DATA
 - Primary data analysis (QC)
 - Secondary data analysis



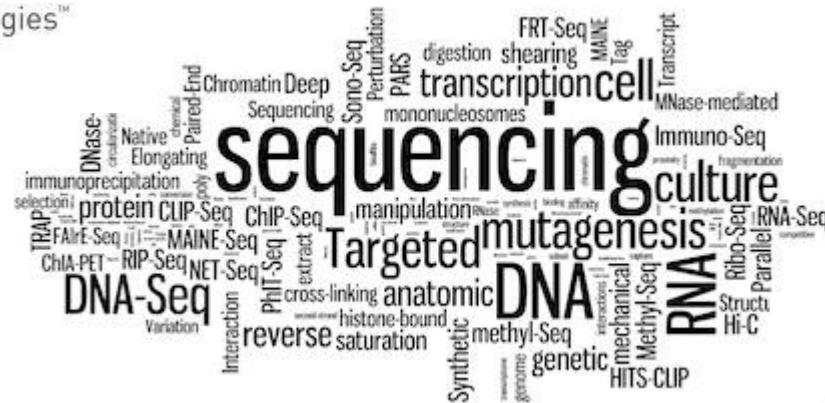
Description of NGS technologies

NGS Technologies

edinburgh
genomics.

illumina[®]

AB applied biosystems™
part of *life* technologies™



ion torrent

by life technologies™

 Oxford
NANOPORE
Technologies

PACIFIC
BIOSCIENCES™

1

NGS sequencers

edinburgh
genomics.



Roche 454 FLX+



Illumina GAIIx



Life Tech SOLID 5500



Life Tech Ion Torrent



Helicos Heliscope



Roche 454 Junior



Illumina MiSeq



NextSeq



Illumina HiSeq



Life Tech Ion Proton



Pacific Biosciences RS



Oxford Nanopore GridION



Oxford Nanopore MinION



Oxford Nanopore PromethION



Complete Genomics RevoloCity

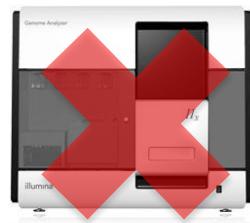


PacBio Sequel

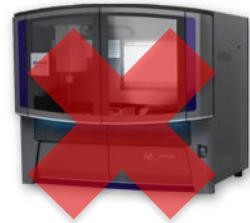
NGS sequencers



Roche 454 FLX+



Illumina GAIIx



Life Tech SOLID 5500



Life Tech Ion Torrent



Helicos Heliscope



Roche 454 Junior



Illumina MiSeq



NextSeq



Illumina HiSeq



Life Tech Ion Proton



Pacific Biosciences RS



Oxford Nanopore GridION



Oxford Nanopore MinION



Oxford Nanopore PromethION

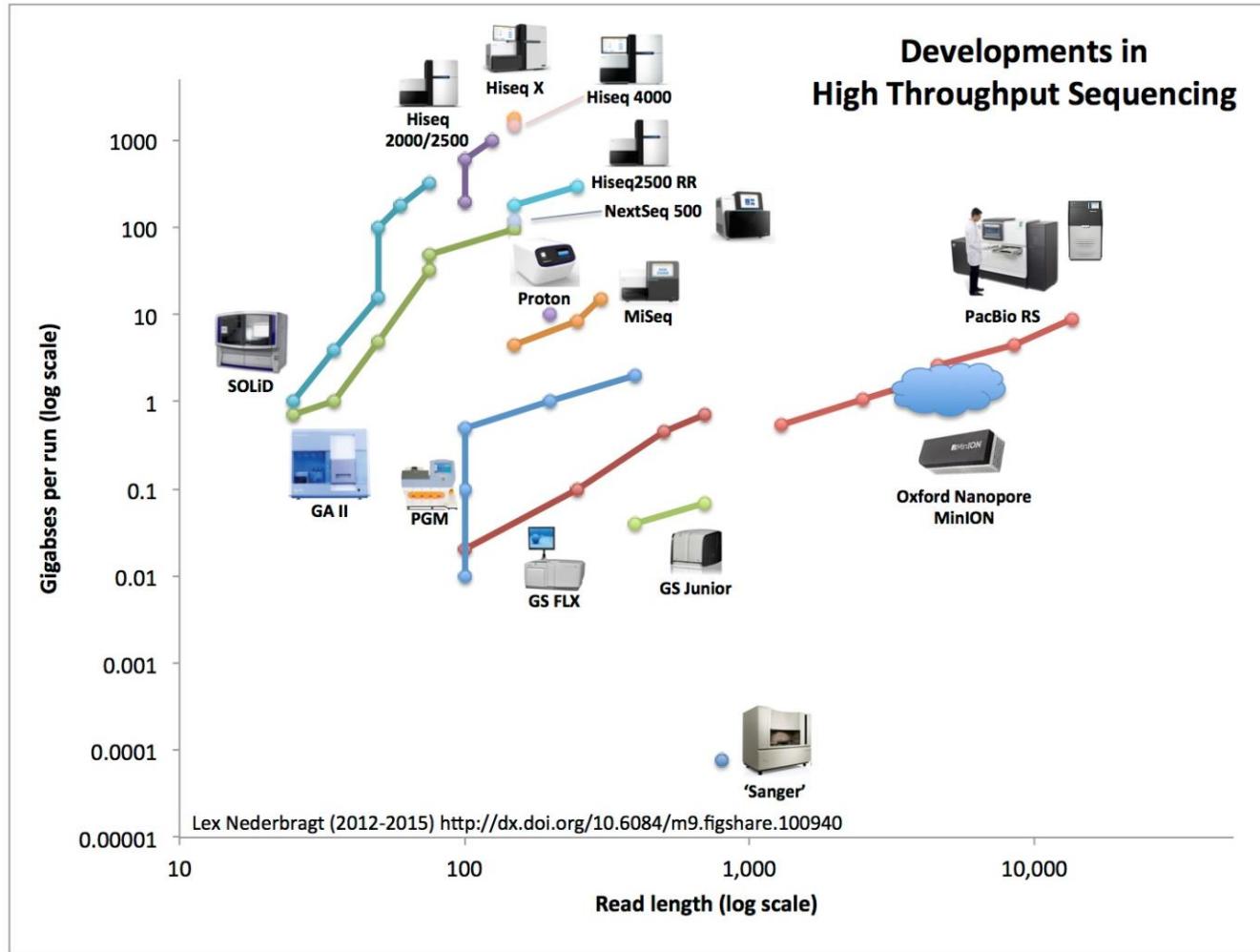


Complete Genomics RevoloCity

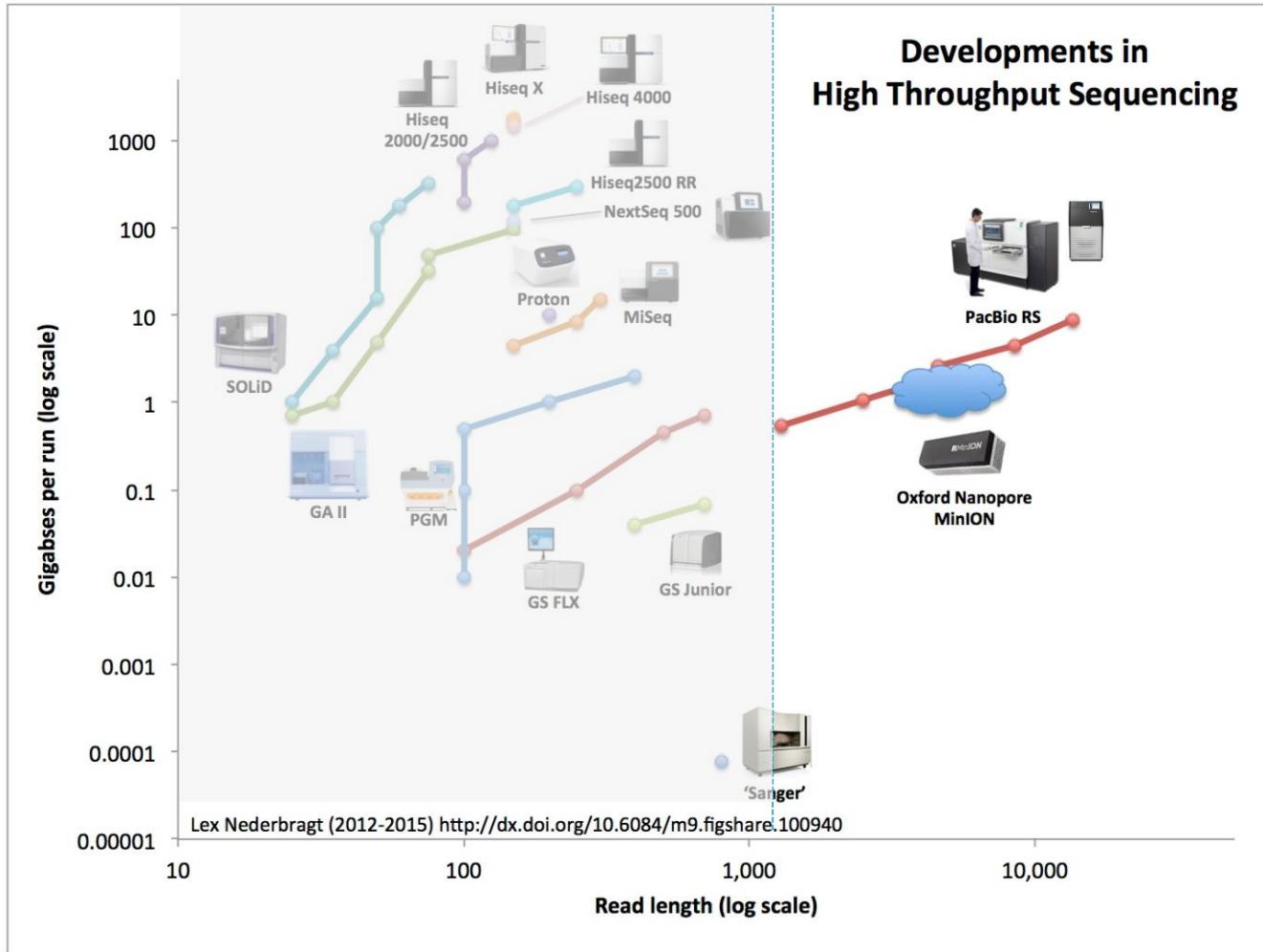


PacBio Sequel

Length & throughput

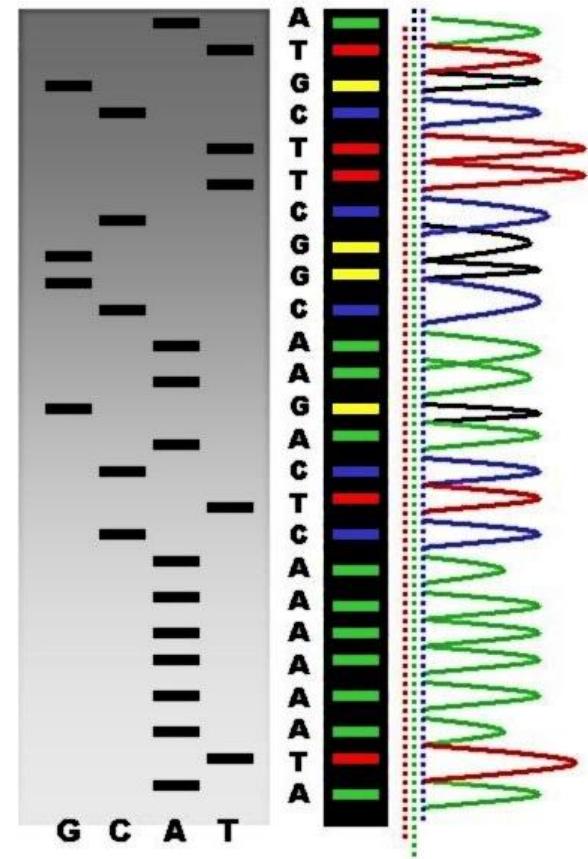
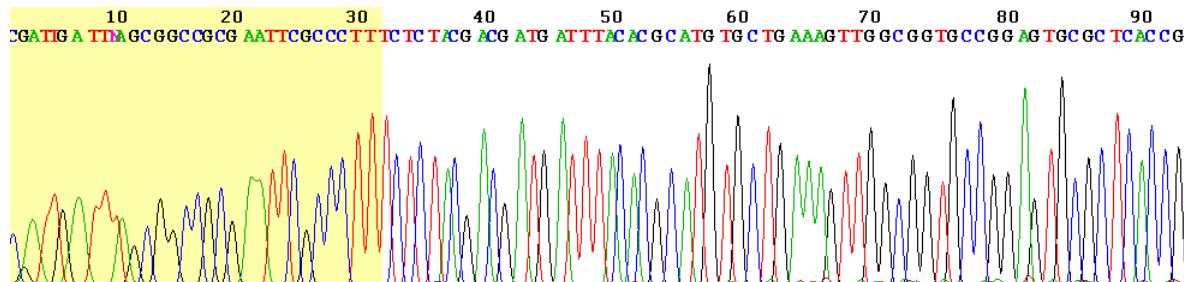


Length & throughput

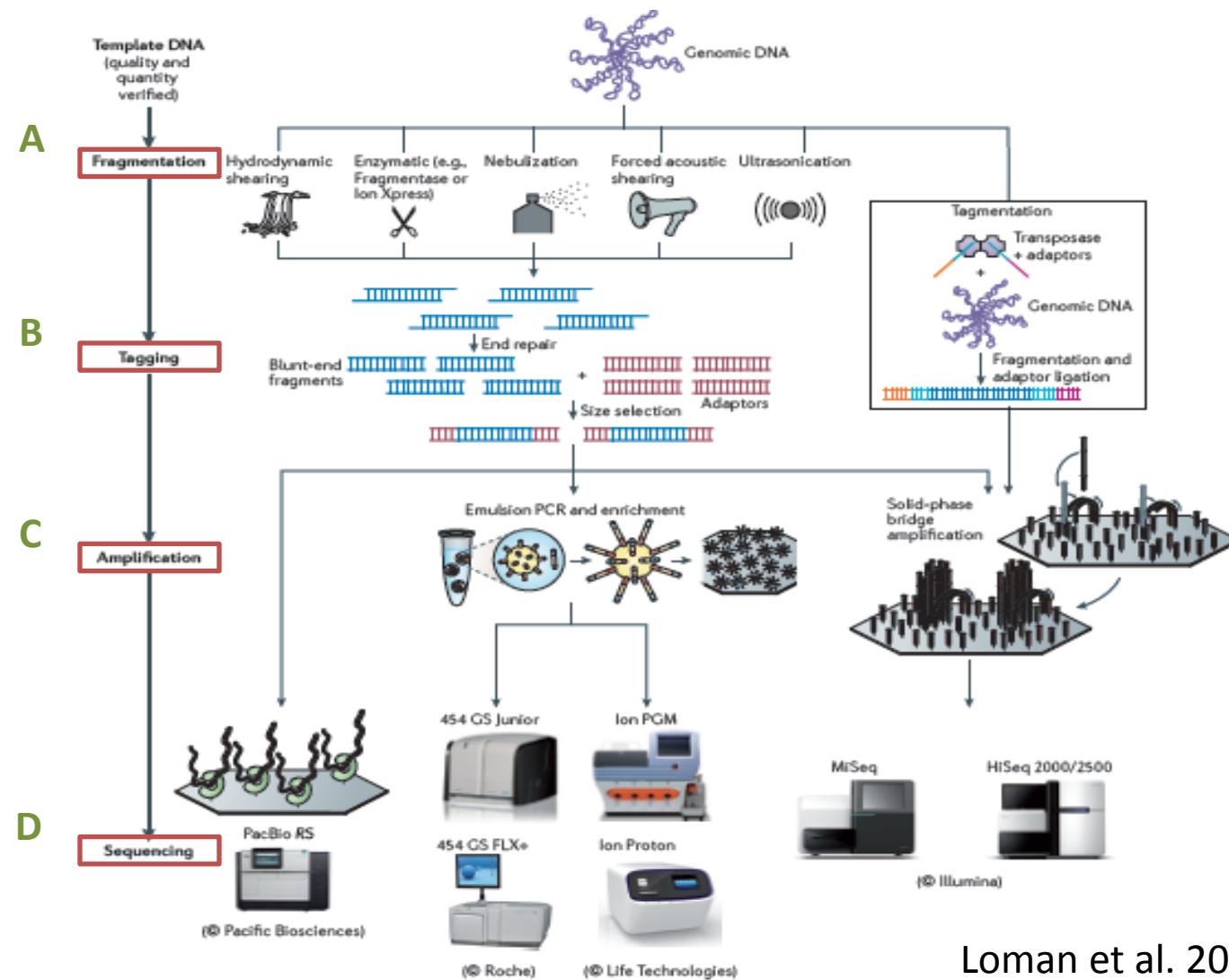


Basics of Sanger Sequencing

- ↳ Clone the DNA.
- ↳ Generate a ladder of labeled (colored) molecules that are different by 1 nucleotide.
- ↳ Separate mixture on some matrix.
- ↳ Detect fluoroscope by laser.
- ↳ Interpret peaks as string of DNA.
- ↳ Strings are 500 to 1,000 letters long
- ↳ 1 machine generates 57,000 nucleotides/run

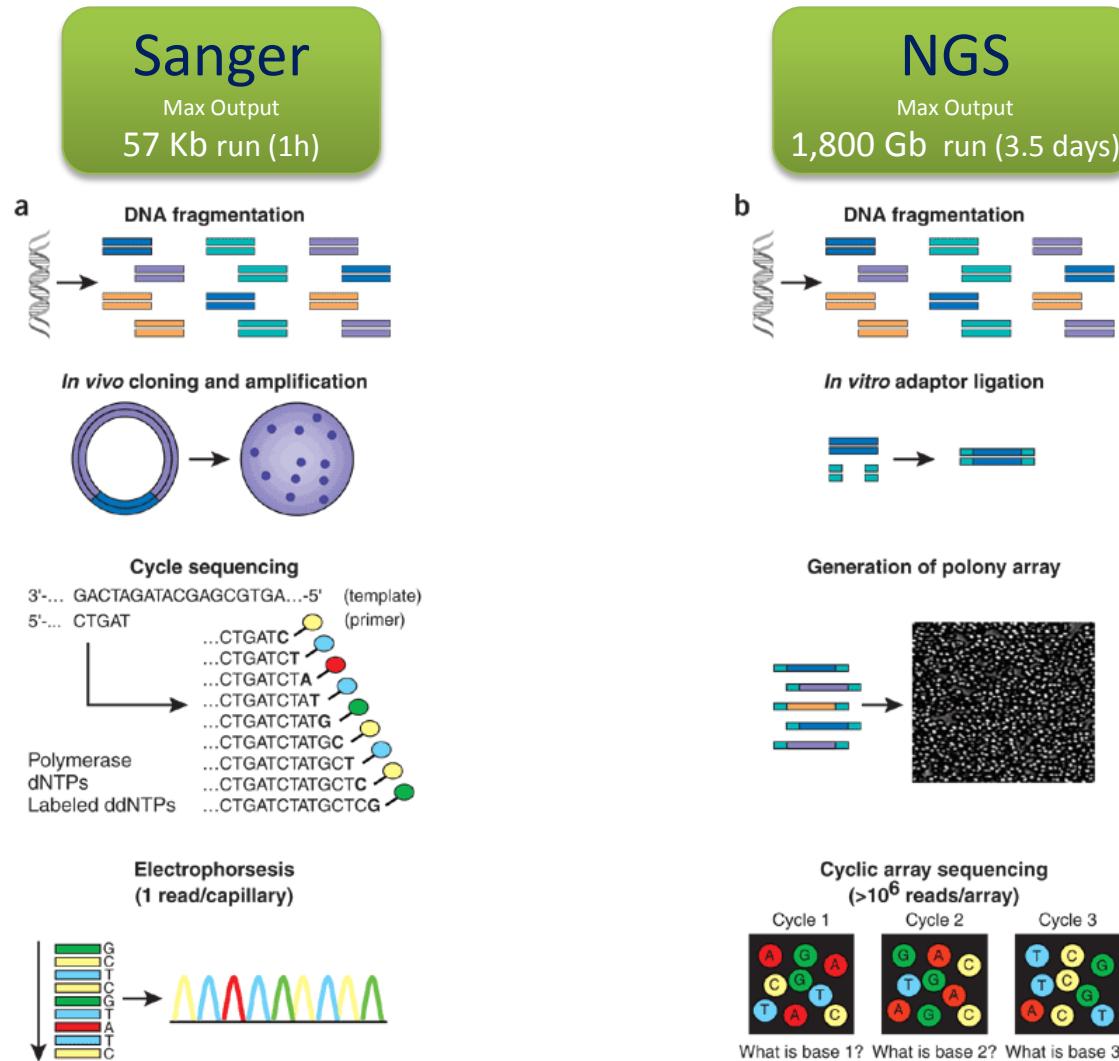


Basics of NGS Sequencing



Loman et al. 2012

Comparison of Technologies



Jay Shendure & Hanlee Ji. Nature Biotechnology 26, 1135 - 1145 (2008)

Illumina HiSeq/MiSeq



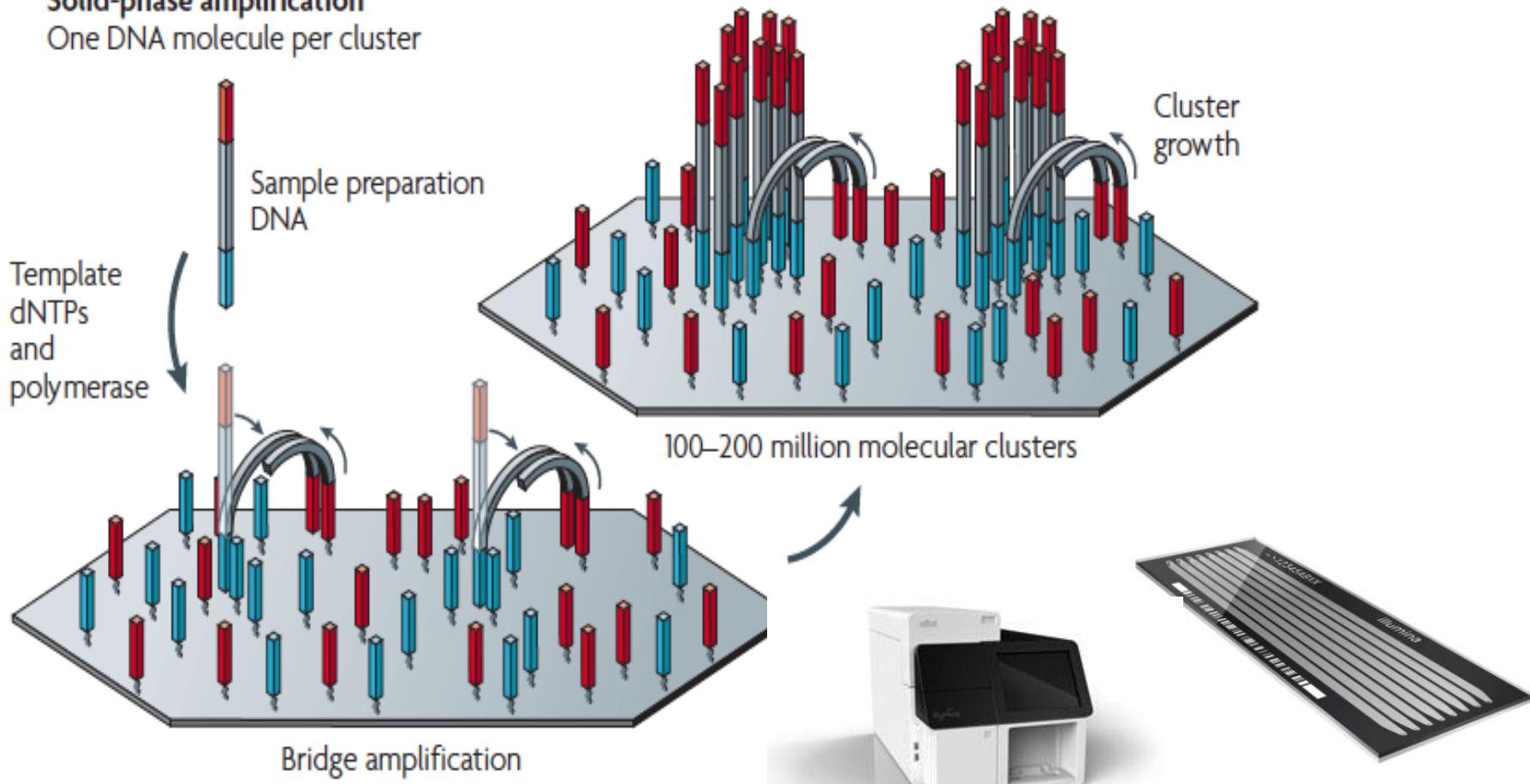
Over 90% of all sequencing data is produced on Illumina systems.

Uses a “sequencing by synthesis” approach:

- **Library:** DNA is broken into small fragments and ligated to adaptors.
- **Amplification:** The fragments are attached to the surface of a flow cell and amplified.
- **Sequencing:** DNA is sequenced by adding polymerase and labeled reversible terminator nucleotides (each base with a different color).
 - The incorporated base is determined by fluorescence.
 - The fluorescent label is removed from the terminator and the 3' OH is unblocked, allowing a new base to be incorporated
- Started with 35 bp, increased now to up to 300 bp
- One run can give up to 10-1,800 Gb, 300-6000 million paired-end reads
- 75-85% of bases at or above Q30
- Substitution errors

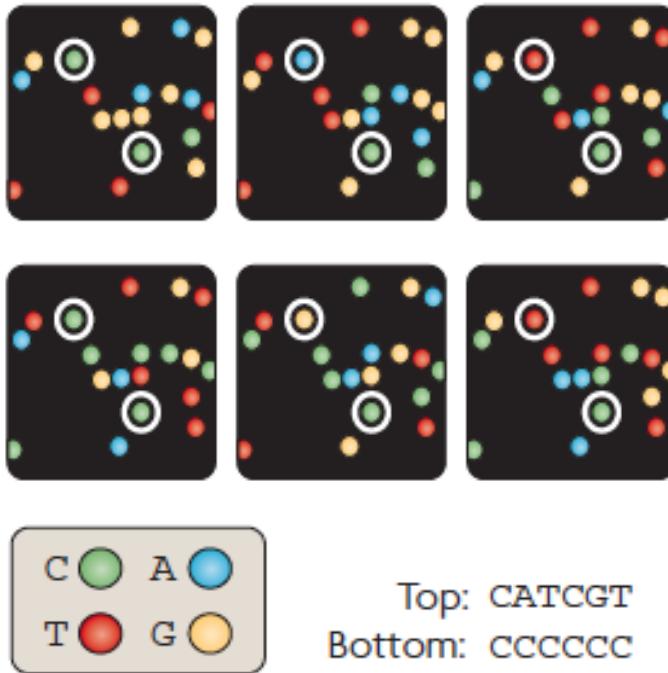
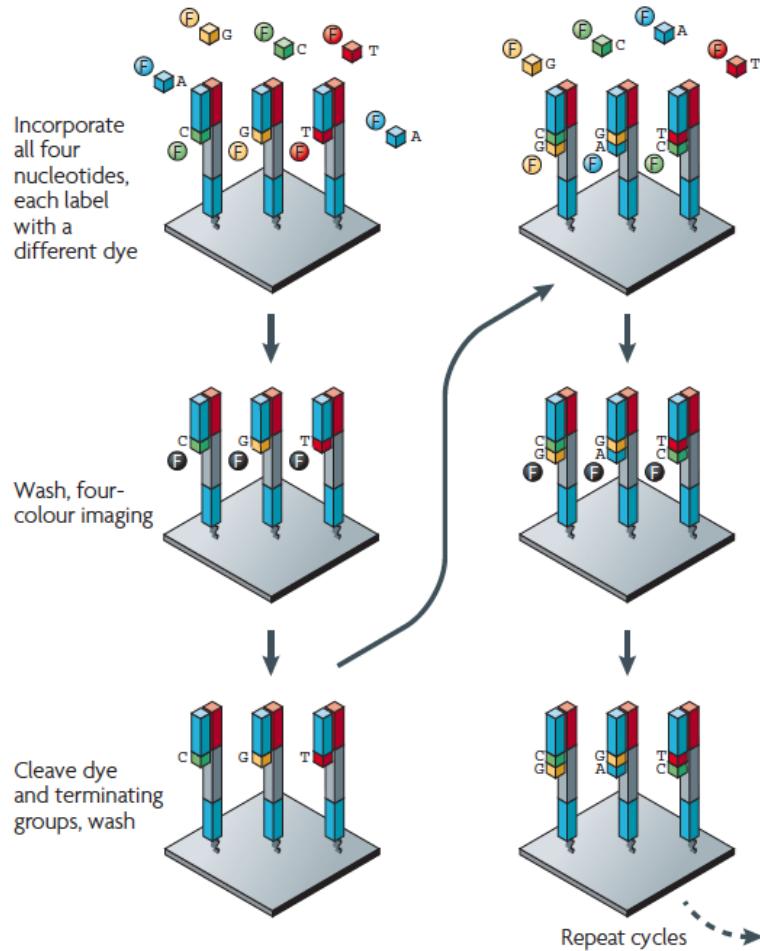
PCR bridge amplification

b Illumina/Solexa
Solid-phase amplification
One DNA molecule per cluster



Adapted from Metzker 2010

Sequencing-by-synthesis



- ❖ Fixed length reads
- ❖ Reversible terminators
- ❖ The identity of each base of a cluster is read off from **sequential images**
- ❖ Illumina sequencing [video](#)

From Michael Metzker, <http://view.ncbi.nlm.nih.gov/pubmed/19997069>

Illumina Sequencers

edinburgh
genomics.

MiSeq



NextSeq 500/550



Max Output
15 Gb

Max Read Number
25 M

Max Read Length
2x300 bp

Max Output
120 Gb

Max Read Number
400 M

Max Read Length
2x150 bp

www.illumina.com

Illumina Sequencers



HiSeq 2500*/3000/4000

* Max Output
1,000* Gb

Max Read Number
4,000* M

Max Read Length
2x125* bp



HiSeq X Ten/ X Five

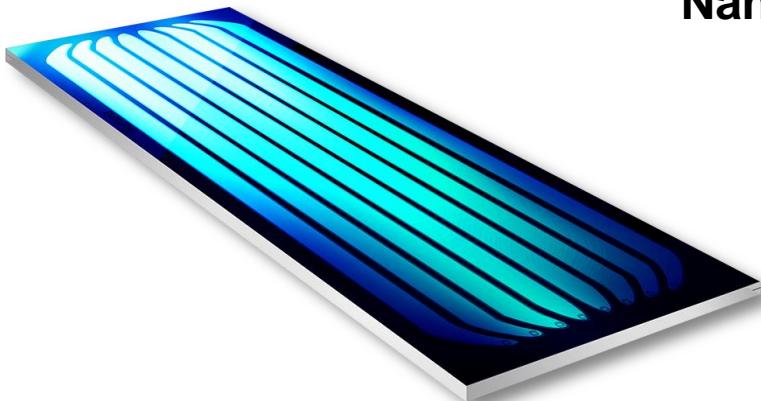
Max Output
1,800 Gb

Max Read Number
6,000 M

Max Read Length
2x150 bp

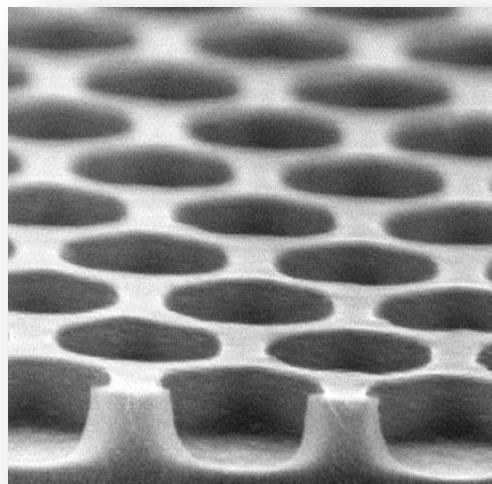
www.illumina.com

Patterned flow cells

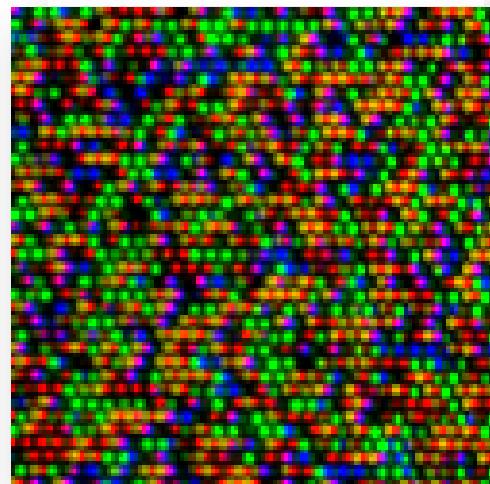


Nanowell substrate | billions of ordered wells

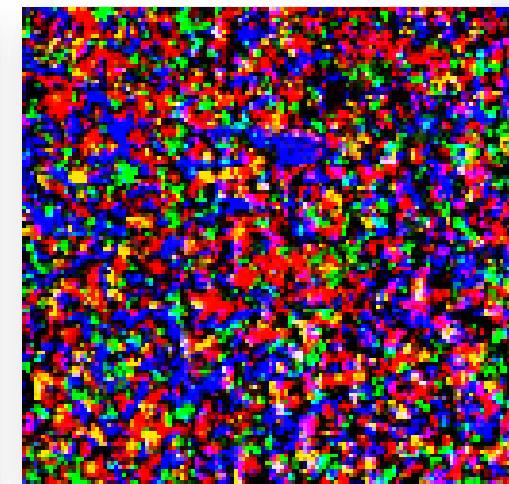
- Defined feature size
- Optimal “fixed” cluster spacing
- Increased cluster density
- Simplified imaging



Ordered spacing



Random spacing

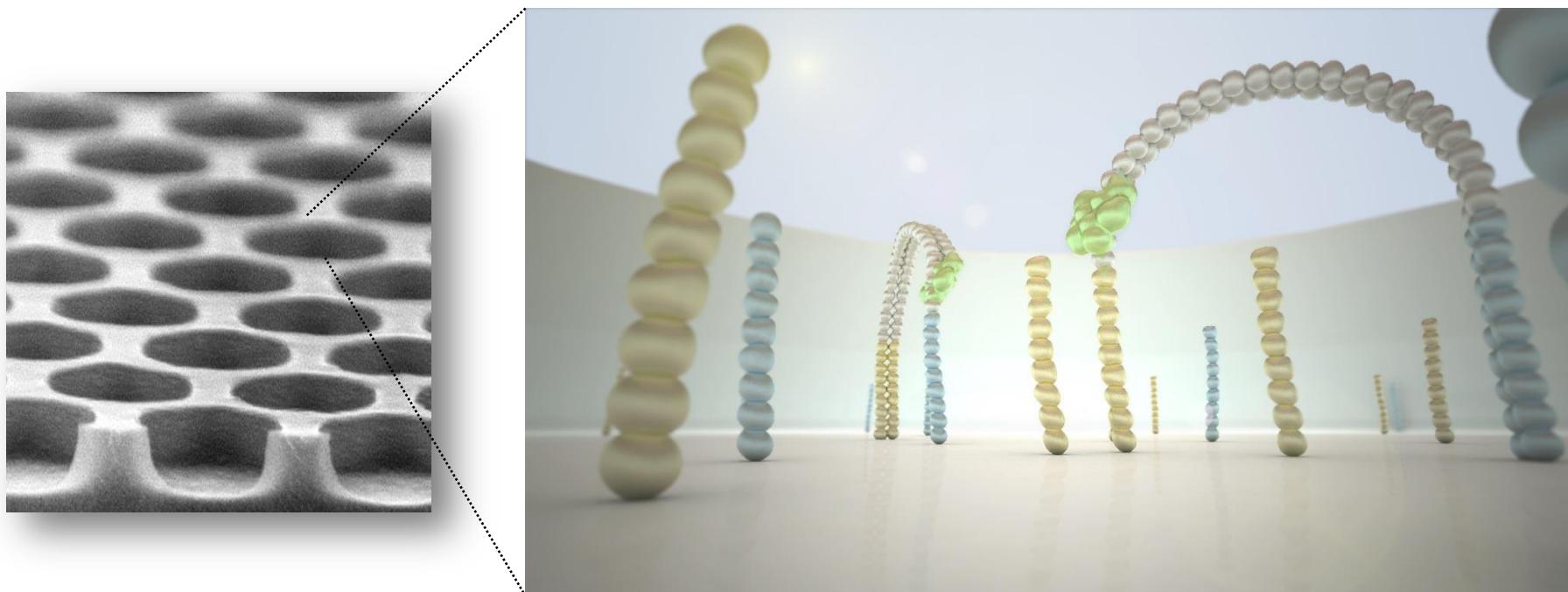


Kinetic Exclusion Amplification (single template per well)

simultaneous template hybridization and amplification

amplification occurs at 20x the rate of template hybridization

- [Patterned cells and KEA video](#)



Other species 30x also welcome



HiSeq X Ten/ X Five

Max Output
1,800 Gb

Max Read Number
6,000 M

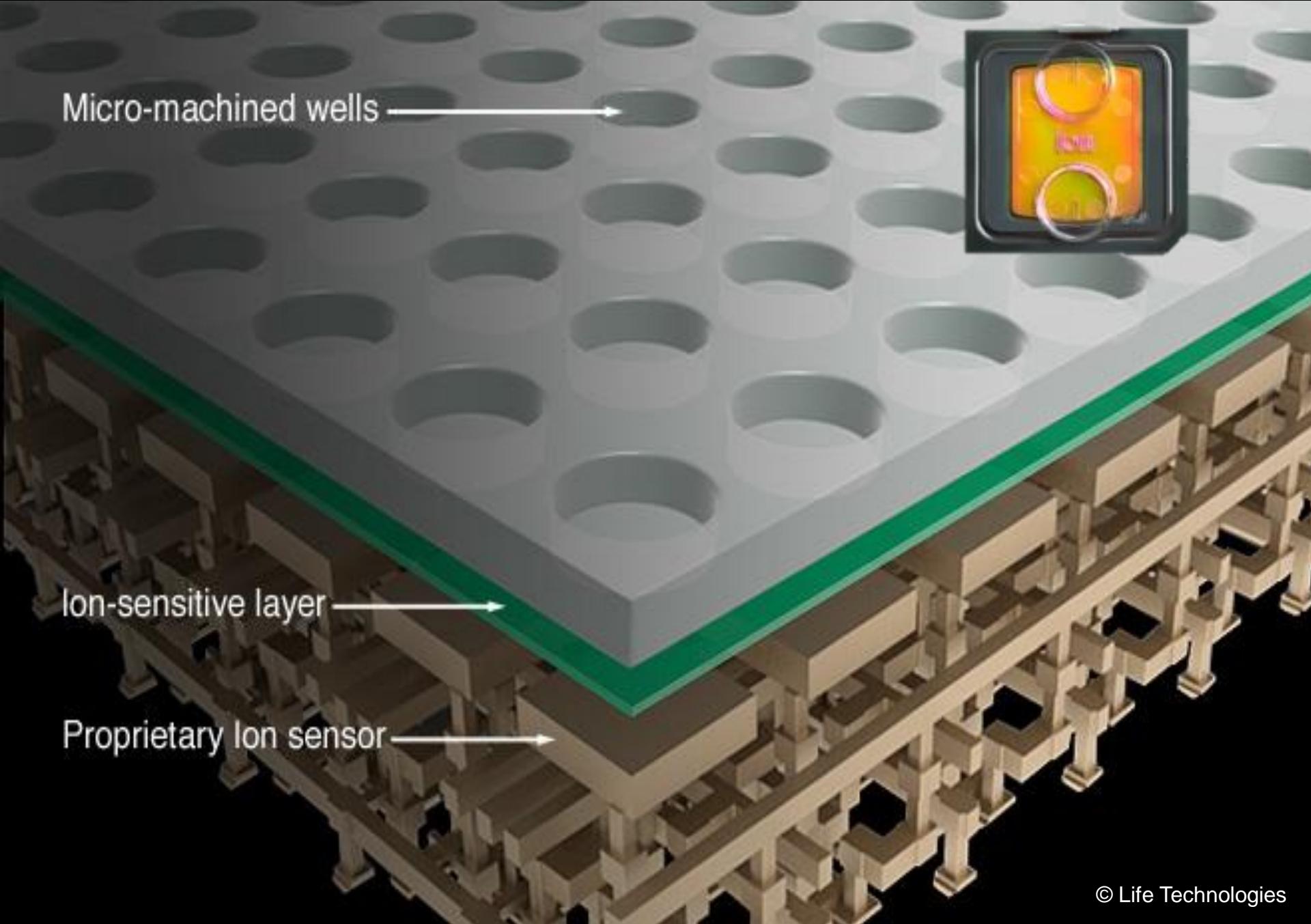
Max Read Length
2x150 bp

www.illumina.com



Ion Torrent/Proton

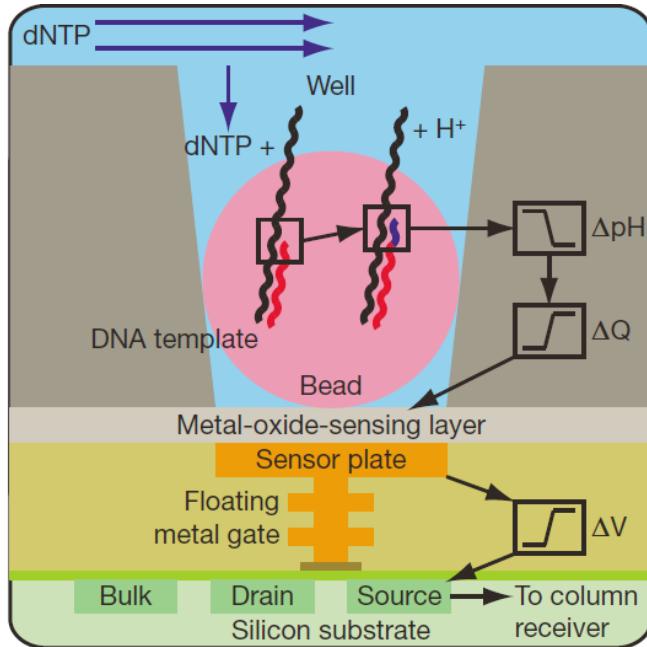




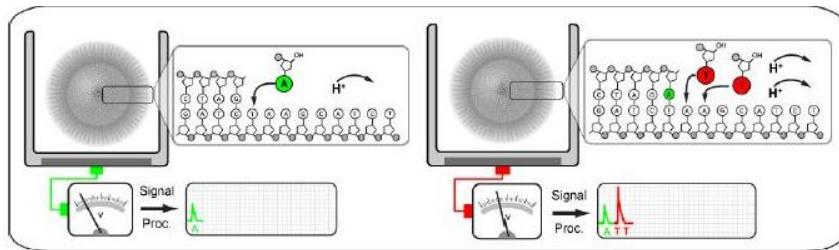
© Life Technologies

Ion Torrent (Life Technologies)

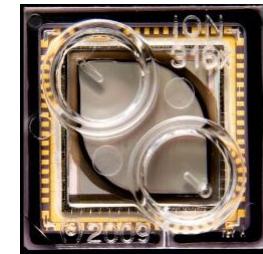
edinburgh
genomics.



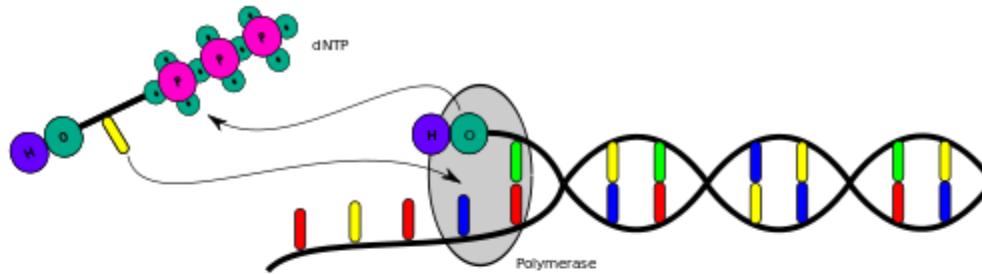
- Similar to pyrosequencing but uses semiconducting chip to detect dNTP incorporation.
- The chip measure differences in pH.
- Different types of chips (throughput/length)
- Shown to have problems with homopolymer reads and coverage bias with GC-rich regions.



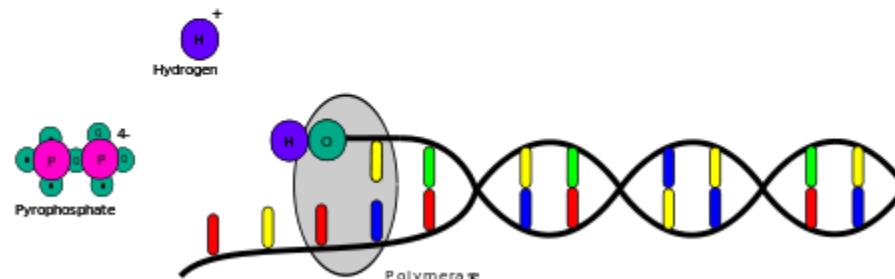
J. M. Rothberg, et al. Nature (2011) 475:348-352



Ion Semiconductor Sequencing



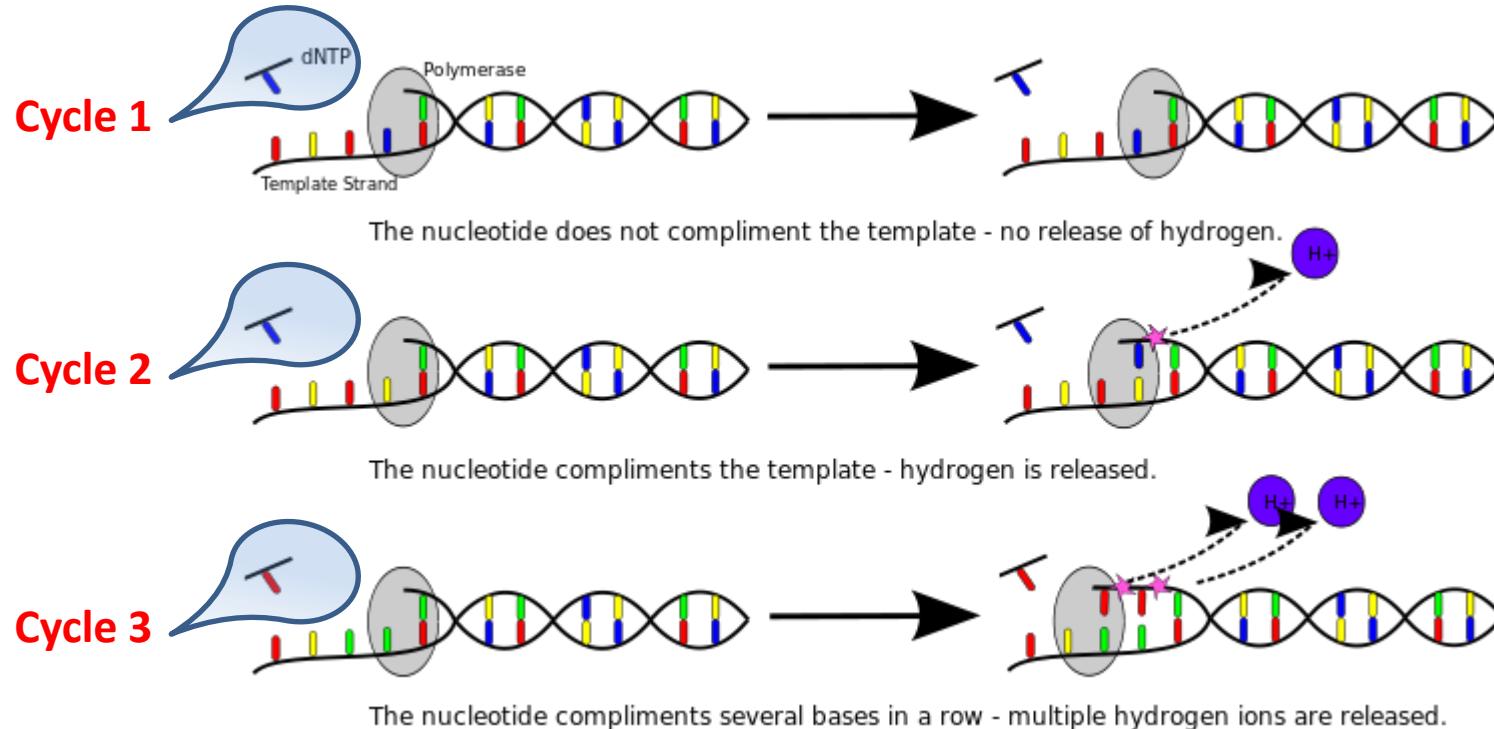
Polymerase integrates a nucleotide.



Hydrogen and pyrophosphate are released.

by David Tack - Own work. Licensed under CC BY-SA 3.0 via Wikimedia Commons

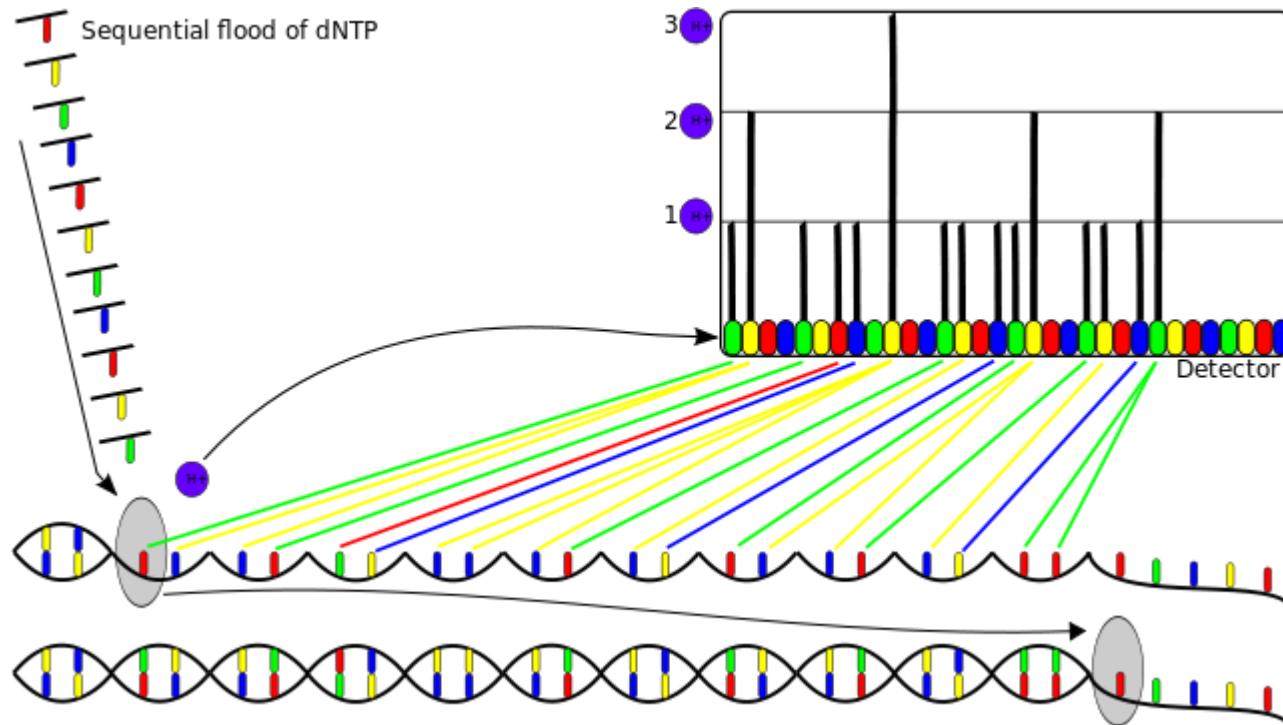
Ion Semiconductor Sequencing



by David Tack - Own work. Licensed under CC BY-SA 3.0 via Wikimedia Commons

Ion Semiconductor Sequencing

edinburgh
genomics.



by David Tack - Own work. Licensed under CC BY-SA 3.0 via Wikimedia Commons

Ion PGM Torrent - Proton

edinburgh
genomics.



10 Mb to 1 Gb

1 hour/run, > 200 nt lengths

Reads H⁺ released by DNA polymerase

Chips: 314, 316, 318

Up to 10 Gb

Up to 200 nt, 2-4 h

Reads H⁺ released by DNA polymerase

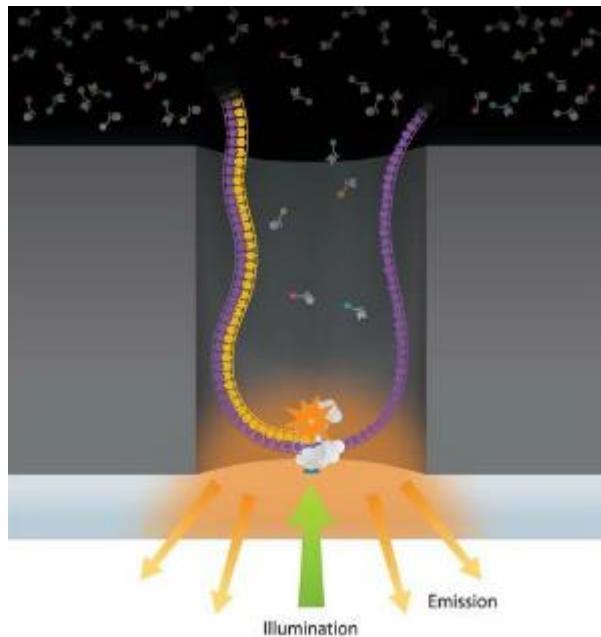
Chips: Proton I, Proton II

Pacific Biosciences



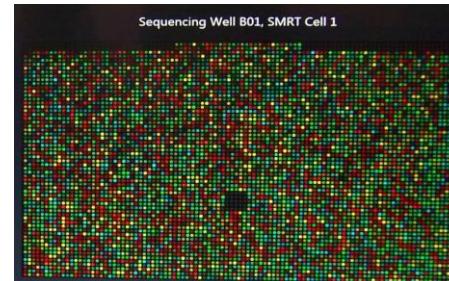
Third Generation: PacBio RS

- **SMRT: Single Molecule Real time DNA synthesis.**
- Single Molecule Sequencing – instead of sequencing clonally amplified templates from beads (Pyro) or clusters (Illumina)
DNA synthesis is detected on a single DNA strand.
 - Up to 15,000 nt, 50 bases/second



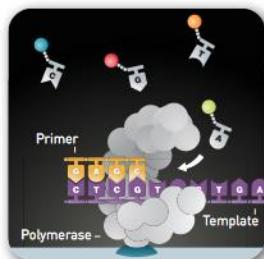
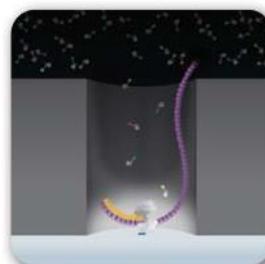
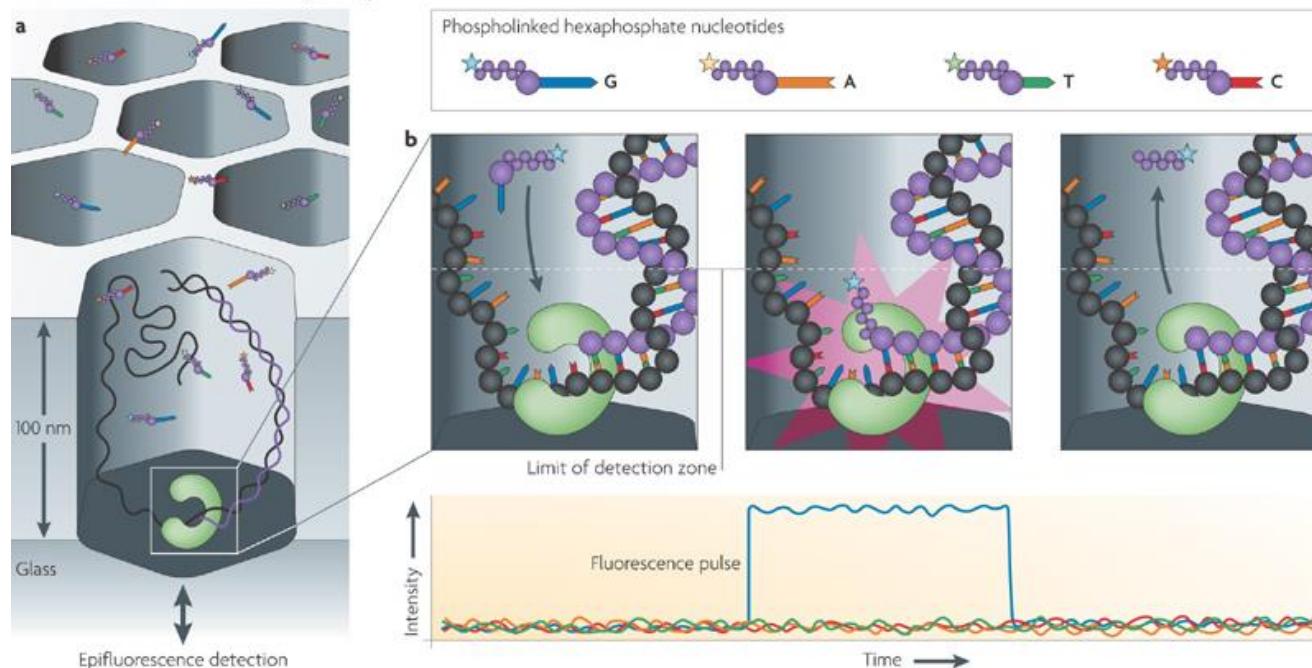
Zero-mode waveguide (ZMW)

- DNA polymerase is affixed to the bottom of a tiny hole (~70nm).
- Only the bottom portion of the hole is illuminated allowing for detection of incorporation of dye-labeled nucleotide.



Third Generation: PacBio RS

edinburgh
genomics.

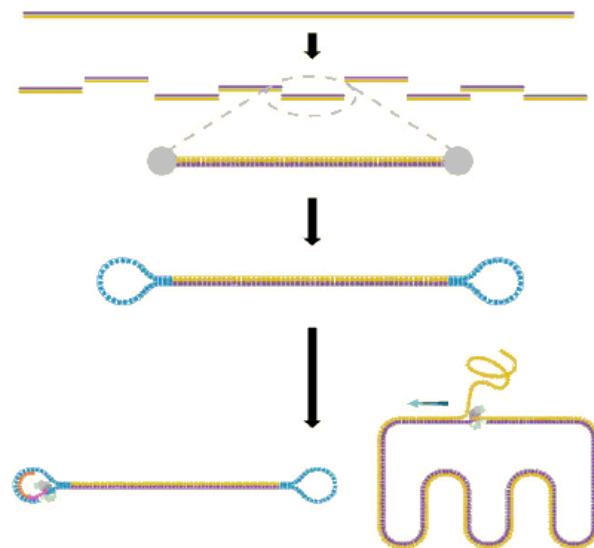


[Link to PacBio movie](#)

From Michael Metzker, <http://view.ncbi.nlm.nih.gov/pubmed/19997069>

Third Generation: PacBio RS

- Real-time Sequencing
 - Unlike reversible termination methods (Illumina) the DNA synthesis process is never halted. Detection occurs in real-time.



Library Prep.

- DNA template is circularized by the use of “bell” shaped adapters.
- As long as the polymerase is stable this allows for continuous sequencing of both strands.

Third Generation: PacBio RS

edinburgh
genomics.



Advantages

- No amplification required.
- Extremely long read lengths.
- Average 2500 nt. Longest 15,000 nt.

Disadvantages

- High error rates.
- Error rate of ~15% for Indels.
1% Substitutions.

Third Generation: Sequel



Advantages

- No amplification required.
- Extremely long read lengths.
- Average 2500 nt. Longest 15,000 nt.
- **7 times more throughput than an RS II**
- **Half price of an RS II.**

Disadvantages

- High error rates.
- Error rate of ~15% for Indels. 1% Substitutions.

Oxford Nanopore MinION



Oxford Nanopore: MinION



The MinION is a memory key–sized disposable unit that can be plugged into a laptop for under \$1,000, according to the company.

- ❖ Announced Feb. 2012 at ABGT conference.
- ❖ Portable device
- ❖ Disposable device with a sensor chip and nanopores
- ❖ Plugging directly into a USB port
- ❖ Real-time sequencing data
- ❖ DNA sequencing or protein sensing
- ❖ MiniON Access Programme (Nov 2013-Jan 2014)



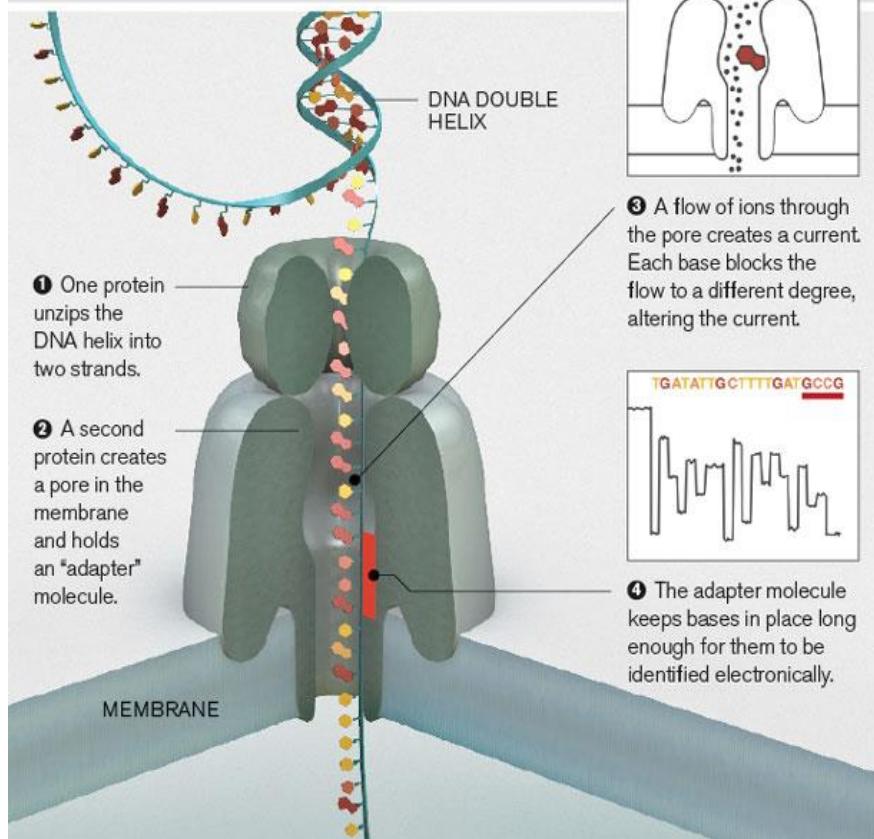
Tricorder

It is a multifunction hand-held device used for sensor scanning, data analysis, and recording data
[Wikipedia]

Third Generation: Oxford Nanopore

edinburgh
genomics.

DNA can be sequenced by threading it through a microscopic pore in a membrane. Bases are identified by the way they affect ions flowing through the pore from one side of the membrane to the other.



<http://www2.technologyreview.com>



- Measure changes in ion flow through nanopore.
- Potential for long read lengths and short sequencing times.

[Link to MinION movie](#)

PromethION system

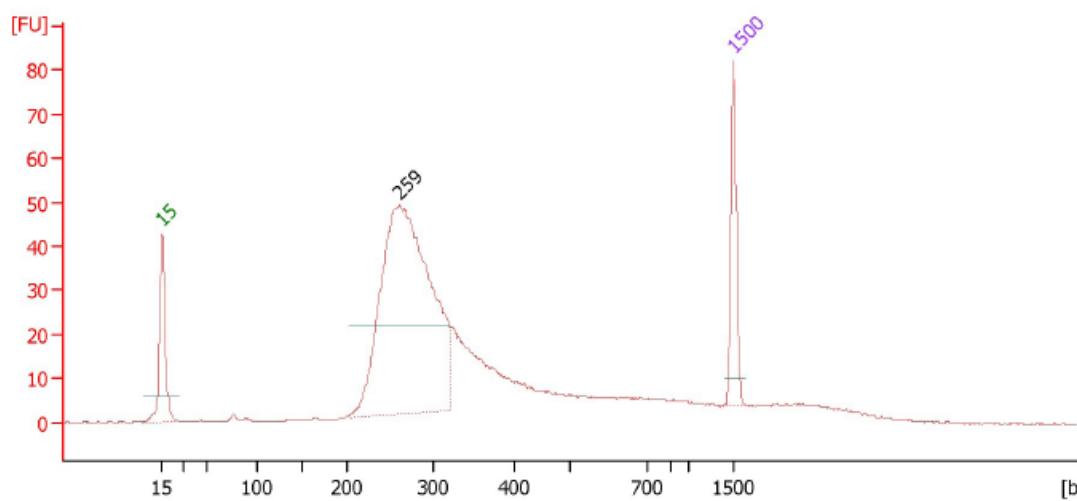


- Scale up system
- Like 96 minION in parallel
- You can use multichannel to load samples into flowcells
- Same chemistry as MinION
- Same software to analyse data
- PromethION early access program PEAP
 - \$ 75K deposit

Anatomy of an NGS Library

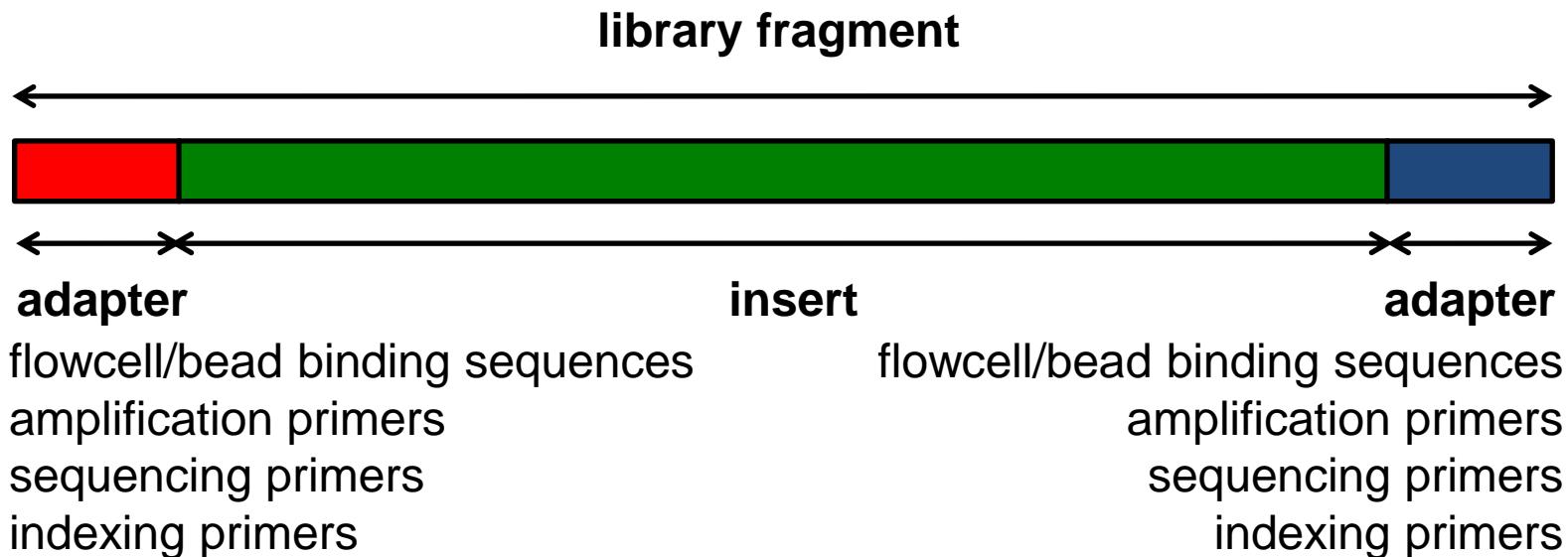
Illumina NGS Library

edinburgh
genomics.



library fragment

Illumina NGS Library



Illumina NGS Library

edinburgh
genomics.

Read



Illumina NGS Library

edinburgh
genomics.

Read 1



Read 2

Illumina NGS Library

Read



- single-end read
- partial sequence from library fragment

Read



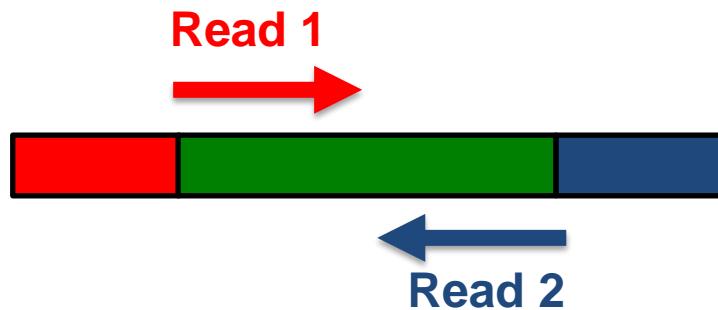
- single-end read
- complete sequence from library fragment
- partial (or complete) adapter sequence

Read

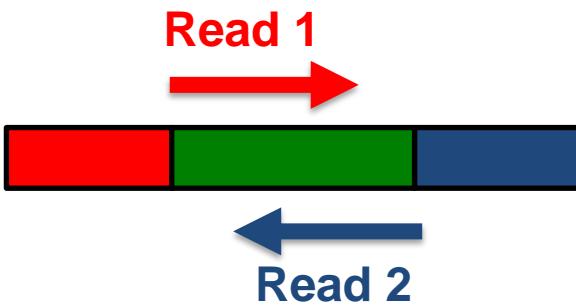


- single-end read
- no library fragment
- partial (or) complete adapter sequence

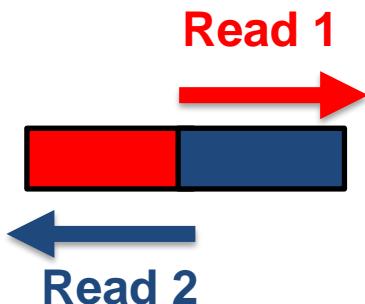
Illumina NGS Library



- paired-end read
- non-overlapping reads
- partial sequence from library fragment



- paired-end read
- complete sequence from library fragment
- overlapping reads

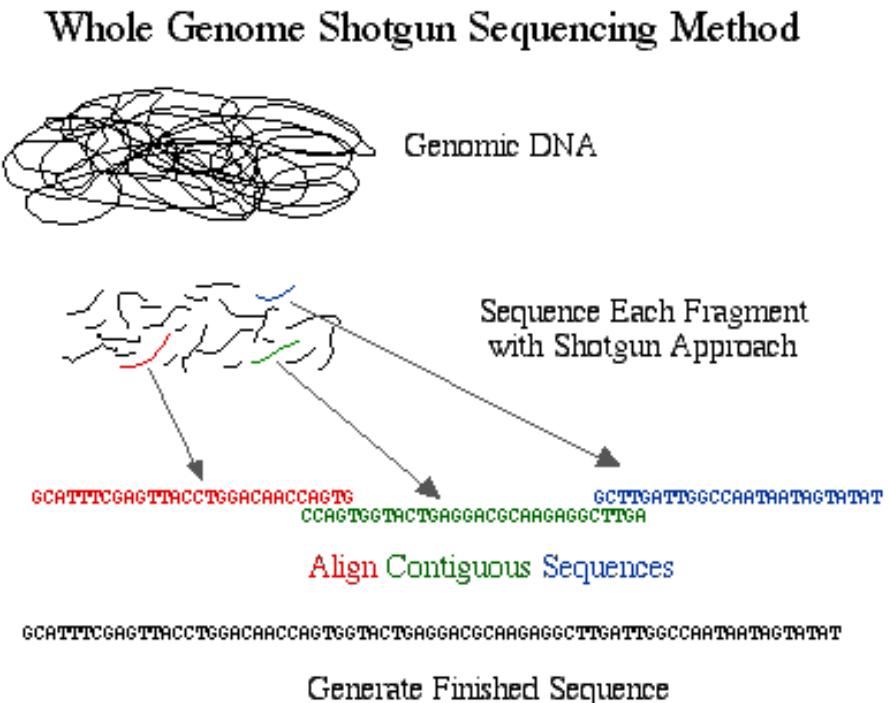


- paired-end read
- no library fragment
- non-overlapping reads
- partial (or) complete adapter sequence

NGS for variant calling

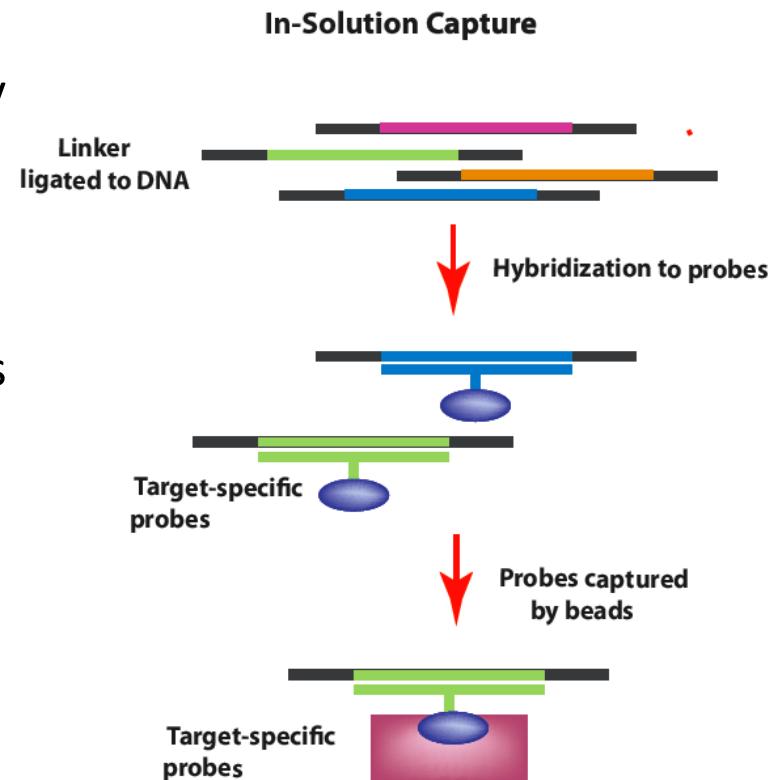
DNA Sequencing - 1

- **Whole GENOME Resequencing**
 - Need reference genome
 - Variation discovery



DNA Sequencing - 2

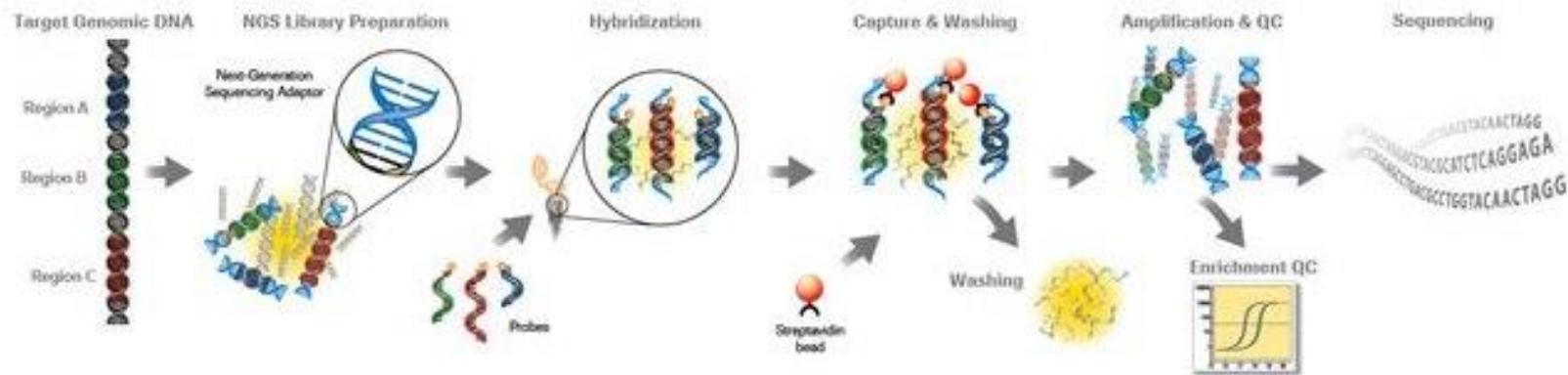
- **Targeted Resequencing**
 - Specific regions in the genome
 - Need reference genome
 - Need custom probes complementary to the genomic regions
 - Nimblegen
 - Agilent
- **Custom genes panel sequencing**
 - Allows to cover high number of genes related to a disease
 - Low cost and quicker than capillary sequencing
 - *E.g. Disease gene panel*
- **Whole EXOME Resequencing**
 - Available for Human and Mouse
 - Variation discovery on ORFs
 - 2% of human genome (lower cost)
 - 85% disease mutation are in the exome



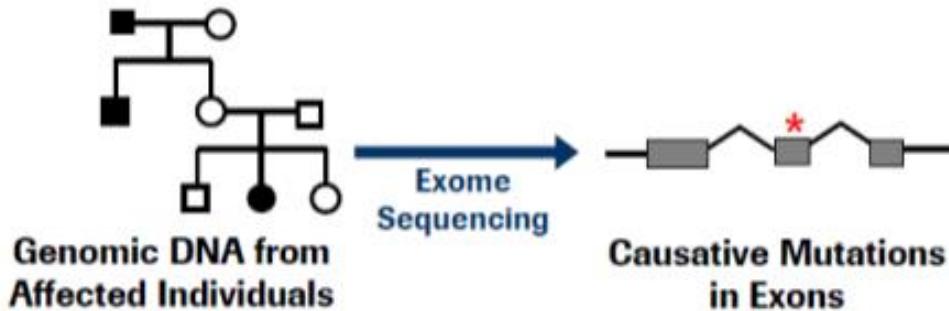
DNA Sequencing - 3

Don't sequence all, just what you need

- Focus on the Most Relevant Portion of the Genome



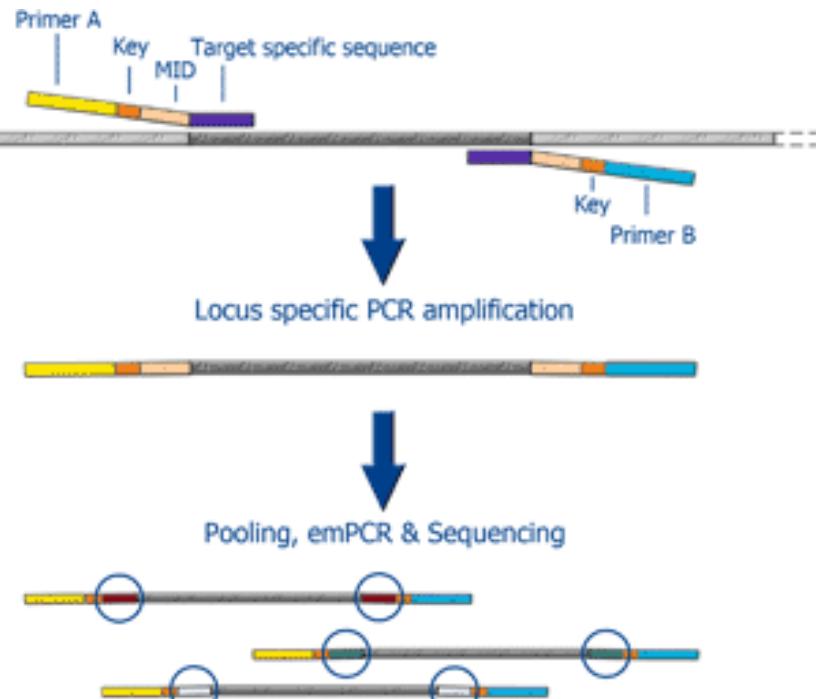
- Capture all exons in the genome: EXOME
 - the most functionally relevant ~2% of the genome.
 - where the majority of known inherited disease-causing mutations reside.



DNA Sequencing - 4

- **Amplicon sequencing**

- Sequencing of regions amplified by PCR.
- Shorter regions to cover than targeted capture
- No need of custom probes
- Primer design is needed
- High fidelity polymerase
- Multiplexing is needed
- Low complexity. Lower quality



NGS Data Analysis

- Similar pipeline -

DNA Sample



NGS Instrument



Data

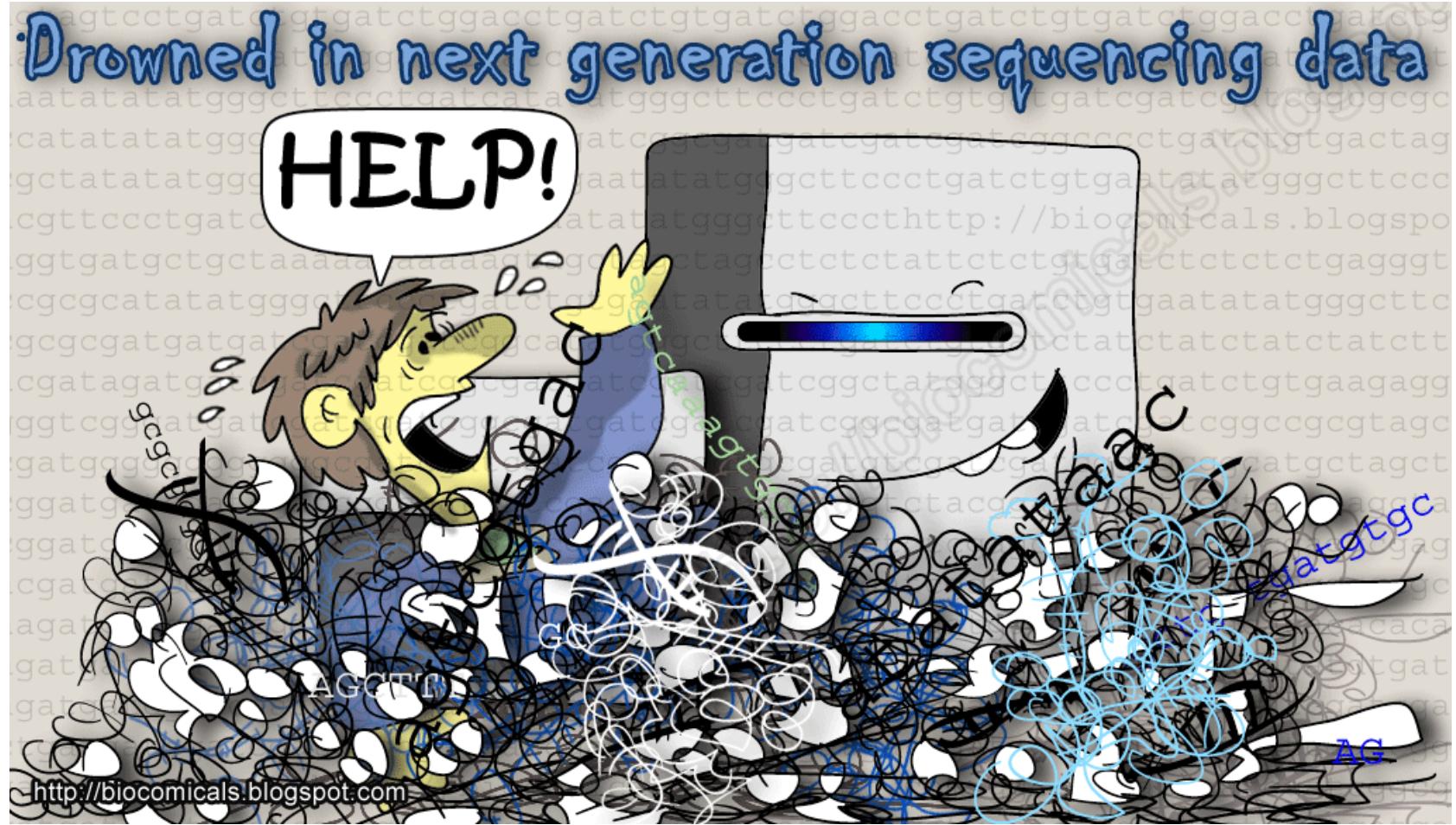
gctacaccttaag
acttcgtcaaa
acttcgtcaaa
acgtaccgtaa
gctacaccttaag
accttaggcctt
acgtaccgtaa
gctacaccttaag
accttaggcctt

Library Preparation

Sequencing

Data Analysis

NGS is relatively cheap but you have to think how you are going to analyze **HUGE AMOUNTS of DATA**



The screenshot shows the Edinburgh Genomics website homepage. At the top, there is a dark header bar with the "Log-in" link, social media icons for YouTube and LinkedIn, and navigation links for "WORK WITH US", "NEWS", "EVENTS", "RESOURCES", and "CONTACT US". Below the header is a search bar with a magnifying glass icon. The main content area features four circular icons representing different services: "Sequencing" (showing a DNA sequencing gel), "Genotyping-Arrays" (showing a person holding a smartphone), "Bioinformatics" (showing a fly being analyzed), and "Clinical Genomics" (showing a human figure composed of a DNA helix). Each service section includes a brief description and a "LEARN MORE" button.

Sequencing

We are skilled in delivering high-quality next-generation sequencing data from RNA and DNA samples, using our advanced Illumina HiSeq2500 and MiSeq platforms. We also regularly carry out Sanger dideoxy sequencing, using AB3730 instruments.

[LEARN MORE](#)

Genotyping-Arrays

High-density genotyping of pedigrees and populations, and microarray analyses of expression, binding and copy number variation are mainstays of genetic and genomic analysis. We offer genotyping and microarray analyses using Illumina and Affymetrix platforms.

[LEARN MORE](#)

Bioinformatics

Bioinformatic analysis is at the heart of what we do, and our dedicated team of research informaticians is skilled in all aspects of next generation genomics and genetic analyses. We have an extensive, secure compute install for high-throughput data processing.

[LEARN MORE](#)

Clinical Genomics

Whole Genome Sequencing of human samples using the Illumina HiSeq X 10 platform produces high quality data at an unprecedented speed and low cost. Laboratory automation using Illumina SeqLab platform allows very high efficiency and low variability.

[LEARN MORE](#)

GRACIAS ARIGATO SHUKURIA JUSPAXAR TASHAKKUR ATU SUKSAMA EKHMET YAOHANYELAY TINGKI BIVAN SHUKRIA

MAKES PARDIES LAH FAKAUE GAEJTHO SAICO MERASTAWHY GOZAIMASHITA EFCHARISTO BAIKA TAVATAPUCH MEDAWAGSE SPASSIBO NUHUN SHACHALHYUA CHALTU WABEEJA MAITEKA DHARYVAAD YUSPIGARATAM HUI GU HATIR EKOAU SIKOMO MAKETAI MINMONCHAR

GRAZIE MEHRBANI KOMAPSUMNIDA MAKE ATTO UNALCHEESH

MAKES DENKAUJA HENACHALHYA